1    **Benchmarked approaches for cell lineage reconstructions of *in vitro* dividing cells and *in***

2    ***silico* models of *Caenorhabditis elegans* and *Mus musculus* developmental trees.**

3

4

5    Wuming Gong[1]\*, Alejandro Granados[2]\*, Jingyuan Hu[3]\*, Matthew G Jones[4,5]\*, Ofir Raz[6]\*,

6    Irepan Salvador-Martínez[7]\*, Hanrui Zhang[8]\*, Ke-Huan K. Chow[2], Il-Youp Kwak[9], Renata

7    Retkute[10], Alidivinas Prusokas[11], Augustinas Prusokas[12], Alex Khodaverdian[4], Richard Zhang[4],

8    Suhas Rao[4], Robert Wang[4], Phil Rennert[13],Vangala G. Saipradeep[14], Naveen Sivadasan[14], Aditya

9    Rao[14], Thomas Joseph[14], Rajgopal Srinivasan[14], Jiajie Peng[15], Lu Han[15], Xuequn Shang[15], Daniel

10   J. Garry[1], Thomas Yu[16], Verena Chung[16], Michael Mason[16], Zhandong Liu[3], Yuanfang Guan[8],

11   Nir Yosef [4], Jay Shendure[17,18,19,20], Maximilian J. Telford[7], Ehud Shapiro[6], Michael B. Elowitz[2],

12   Pablo Meyer[21+]

13

14   [1] Lillehei Heart Institute, University of Minnesota, 2231 6th St S.E, 4-165 CCRB, Minneapolis,

15   MN 55114, USA.

16   [2] California Institute of Technology, Pasadena, CA, USA

17   [3] Program in Quantitative and Computational Biosciences, Baylor College of Medicine Houston,

18   TX, USA

19   [4] Department of Electrical Engineering & Computer Science, University of California, Berkeley,

20   CA, USA

21   [5] Integrative Program of Quantitative Biology, University of California, San Francisco.

22   [6] Department of Computer Science and Applied Mathematics, Weizmann Institute of Science,

23   Rehovot 761001, Israel.

24   [7] Centre for Life's Origins and Evolution, Department of Genetics, Evolution and Environment,

25   University College London, Gower Street, London, WC1E 6BT, UK.

26   [8] Department of Computational Medicine and Bioinformatics, University of Michigan, Ann

27   Arbor, Michigan, USA

28   [9] Department of Applied Statistics, College of Business & Economics, Chung-Ang University,

29   84, Heukseok-ro, Dongjak-gu, Seoul, Republic of Korea

30   [10] Department of Plant Sciences, University of Cambridge, Downing Street, Cambridge CB2

31   3EA, UK

32   [11] School of Natural and Environmental Sciences, Newcastle University, Newcastle, NE1 7RU,

33   UK

34   [12] Department of Life Sciences, Imperial College London, London, SW7 2AZ, UK

35   [13] EC Wise Inc.

36   [14] TCS Research and Innovation, Tata Consultancy Services Ltd., Hyderabad, INDIA

37   [15] School of Computer Science, Northwestern Polytechnical University, Xi'an, China

38   [16] Sage Bionetworks, Seattle, WA, USA

39   [17] Department of Genome Sciences, University of Washington, Seattle, WA, USA

40   [18] Allen Discovery Center for Cell Lineage Tracing, Seattle, WA, USA

41   [19] Brotman Baty Institute for Precision Medicine, Seattle, WA, USA

42   [20] Howard Hughes Medical Institute, Seattle, WA, USA

43   [21] T.J. Watson Research Center, IBM, Healthcare & Life Sciences, 1101 Kitchawan road 10598,

44   Yorktown Heights, NY, USA

45

46   *Contributed equally and sorted in alphabetical order

47   +Corresponding author

48

49

## Abstract

The recent advent of new CRISPR-based and other molecular tools now enables the reconstruction of cell lineages based on DNA mutations induced by CRISPR and promises to solve the lineage of complex model organisms at single-cell resolution. To date, however, no lineage reconstruction algorithms have been rigorously examined for their performance and robustness across datasets, diverse molecular tools, and most importantly the number of cells in the lineage tree. In order to benchmark methods of cell lineage reconstruction we decided to organize the Allen institute lineage reconstruction DREAM challenge where we rigorously examined multiple methods using experimental and *in silico* data. On one hand, we took advantage of intMEMOIR recordings, a recently developed synthetic image-readable lineage tracing technology, and asked participants to reconstruct the lineages for 30 *in vitro*-grown mouse embryonic stem cell colonies. We also provided *in silico* datasets for a *C. elegans* lineage tree of about 1000 cells and a simulation of one year of *Mus musculus* development down-sampled to 10,000 cells upon which we simulated CRISPR-based GESTALT-like recordings. For these three lineage reconstruction tasks we provided training data with the ground true trees, as one of the goals of this challenge was to encourage machine-learning approaches different from the ones used for phylogenetics. The challenge was successful in its main goal of attracting a variety of successful approaches and teams: twenty-two full submissions were received and scored using two different metrics. The availability of a training set allowed not only the development of a successful machine-learning decision-tree based approach, but also the optimization of accurate distance-based algorithms and maximum parsimony approaches. This DREAM challenge was a first attempt to rigorously examine the performance and robustness of various reconstruction algorithms under varying conditions and underlies the importance of using several metrics when evaluating reconstruction accuracy. For the experimental dataset, we found that while some trees were reconstructed perfectly, the overall scores were far from the theoretical maximum, mainly due to the structural features of the trees and not the high degeneracy in recorded states across cells. On the other hand, the *in silico* results showed that using smaller subtrees as training sets is a good approach for tuning the algorithms to reconstruct larger trees. Together, these results and the availability of tools for generating and solving lineage trees delineate a potential way forward for solving larger cell lineage trees such as for mouse and human.

81

## Introduction

83

**Lineage inference for understanding Development**

85    A fundamental challenge in biology is the reconstruction of the developmental histories

86    of cells as they divide and progress through differentiation into different cell types. Indeed,

87    multicellular organisms can be composed of billions or trillions of cells that derive from a single

88    cell through repeated rounds of cell division. Knowing the lineage relationships between the cells

89    of a fully developed organism -its cell lineage- would provide a framework to understand when,

90    where and how cell fate decisions are made. Further, it can also be useful to understand the

91    progression of disease such as in tumor subclonal reconstruction (Salcedo et al., 2020) or the

92    development of an organ such as the brain (Evrony et al., 2015; Lodato et al., 2015). Historically,

93    lineages of individual cells have only been fully reconstructed by their direct observation through

94    microscopy as for the nematode *Caenorhabditis elegans* (Sulston and Horvitz, 1977). This direct

95    observation approach is however not possible for most animals as the cells are not visible (Livet

96    et al., 2007). In the 1980's new methods allowed marking all the descendants of a single cell by

97    the injection of a dye or the expression of a marker gene. Since then, many new methods have

98    been devised to improve cell lineage tracking, including inducible recombinases (Kretzschmar

99    and Watt, 2012), fluorescent or genetic reporters (Kebschull and Zador, 2018; Weissman and

100    Pan, 2015), or a combination of both (Garcia-Marques et al., 2020). However, these approaches

101    come at the cost of resolution, meaning that lineage relationships of individual cells are not fully

102    recovered.

103    Recent advances in sequencing technologies have enabled a variety of RNA-based

104    methods to infer differentiation trajectories in multiple organisms and cell types by ordering the

105    changes in single-cell gene expression along a pseudo-time axis representing the progression

106    through differentiation (Wagner and Klein, 2020). However, these methods focus on the

107    expression profiles of cells but do not have access to their genealogical relationships. In this

108    regard, somatic mutations accumulated during normal development have been used to

109    reconstruct genetic lineages (Behjati et al., 2014; Frumkin et al., 2005) and for example trace

110    mosaicism in the brain (Evrony et al., 2015; Lodato et al., 2015). Deep sequencing of cDNA

111    from T cell receptors has also been used to establish clonal development of T cells (Becattini et

112    al., 2015). Cell lineage inference has also been done using copy-number variations, structural

113    markers such as SNVs, indels, retrotransposon elements, microsatellite repeats, as well as

114    epigenetic markers such as DNA methylation (Kester and Oudenaarden, 2018).

115    **New lineage recording technologies**

116          Recently, the advent of CRISPR-based molecular tools have produced a new generation

117    of lineage reconstruction approaches inspired by principles of phylogenetic inference using

118    naturally occurring DNA mutations. The DNA-editing technologies have been applied to

119    introduce mutations in the genetic material of cells such that a registry of their genetic

120    relationships is recorded and available for readout by sequencing (Alemany et al., 2018; Chan et

121    al., 2019; McKenna et al., 2016a; Perli et al., 2016; Spanjaard et al., 2018). Indeed, the inserted

122    synthetic construct can accumulate stochastic mutations upon induction of CRISPR-Cas9 activity

123    as cells differentiate during development with the goal of resolving cellular lineages of complex

124    model organisms (McKenna et al., 2016b; Wagner and Klein, 2020). Different versions of

125    CRISPR-based methods such as scGESTALT, LINNAEUS and ScarTrace techniques have been

126    successfully used to investigate cellular lineages in various animal models (Alemany et al., 2018;

127    McKenna and Gagnon, 2019; Raj et al., 2018; Spanjaard et al., 2017). At the same time, other

128    types of lineage recording techniques have been applied to allow readout by *in situ* imaging

129    which enables lineage analysis through the maintenance of the spatial information (Chow et al.,

130    2021; McKenna and Gagnon, 2019).  Some of these approaches have applied phylogenetic

131    reconstruction algorithms to infer the cell lineage, whilst others developed ad-hoc cell lineage

132    reconstruction algorithms, but this explosion of lineage tracing technologies has increased the

133    urgency for new reconstruction methods (Salvador-Martínez et al., 2019).

134          In principle, as in phylogenetic tree reconstruction (Frieda et al., 2017; McKenna et al.,

135    2016a), the recorded mutations should encode enough information enabling inference of the

136    likely tree structures that could represent the actual lineage relationships. However, there are

137    significant challenges for tree-inference when applying standard phylogenetic methods to lineage

138    recordings. The main limitations include noise from the experimental readout, restrictions in the

139    total available 'DNA memory' for recording, and the random convergence of identical edit

140    patterns in non-related cells, or homoplasy (Salvador-Martínez et al., 2019). It also remains

141     unclear whether machine-learning algorithms that go beyond classical phylogenetic methods,

142     such as Neighbor-Joining or Maximum Parsimony, could consistently reconstruct cell lineages

143     with higher accuracy. While phylogenetic methods typically analyze a relatively small number of

144     species and many more DNA sites, genes or even whole genomes (McKenna and Gagnon,

145     2019), CRISPR-based lineage recording aims to capture hundreds to thousands of cells with the

146     compromise of limited numbers of editable sites. Additional limitations include variability in

147     mutation rates for each site, large nucleotide deletions resulting in sequence dropouts, and single

148     deletions that can erase previous mutations or ablate multiple targets. Although maximum

149     parsimony-based methods have shown initial success when applied to lineage tracing (McKenna

150     and Gagnon, 2019; McKenna et al., 2016a; Price et al., 2010), the key differences discussed

151     above make it challenging to directly apply phylogenetic methods to lineage tracing data.

152         After having performed lineage tree inference one would ideally like to evaluate the

153     reconstruction accuracy, however for most of these technologies the ground truth is inaccessible,

154     meaning that we do not know the actual lineage relationships. Indeed, with rare exceptions

155     (Sugino et al., 2019), to date no lineage reconstruction approach has been rigorously examined

156     for its performance/robustness across diverse molecular tools, DNA-based recording methods,

157     datasets, number of cells, topology of lineage trees and diverse metrics used for evaluation.

158     Given the lack of benchmarking, there is still no agreement regarding the best practices for

159     inferring cellular lineages from the recording datasets generated with these recently developed

160     molecular tools.

**The DREAM initiative**

161

162         To catalyze the development of new methods to perform lineage reconstruction, we

163     organized the Allen institute lineage reconstruction DREAM challenge, which ran from October

164     2019 through February 2020. DREAM challenges are a platform for crowdsourcing collaborative

165     competitions where a rigorous evaluation of each submitted solution allows for objective

166     comparison and assessment of their performance (Saez-Rodriguez et al., 2016). The value of

167     DREAM resides not only in the acceleration of research through the participation of many teams

168     while solving a common problem, but just as importantly, in the diversity of approaches used

169     and the quality and reproducibility of each provided solution to problems in emerging areas of

170     biology. The aggregation of the individual solutions, *i.e.*, the different approaches and insights to

171 a common problem, namely the 'wisdom of the crowds', leads to a generally superior

172 performance than any individual solution, from where collective insights can be garnered.


173 **The DREAM challenge for lineage reconstruction**

174       The lineage reconstruction DREAM challenge aimed to provide a new perspective on

175 lineage inference by enabling participants from diverse fields to submit their reconstruction of

176 trees for which the ground truth, *i.e.* the actual lineage, existed but was not provided. It consisted

177 of three challenges with lineages of increasing numbers of cells. The first challenge leveraged a

178 then unpublished experimental dataset of 106 trees recorded with intMEMOIR in mouse

179 embryonic stem cell colonies of less than 100 cells (Chow et al., 2021). This technique was

180 chosen as it has the key advantage of readout by imaging which can be coupled with a time-lapse

181 movie of the cells as they divide to provide a ground truth lineage tree (**Fig 1A**). In the second

182 challenge participants had to reconstruct an *in silico* tree of 1,000 cells, whose topology was

183 derived from the *Caenorhabditis elegans* developmental cell lineage tree by removing a few

184 clades in order to mask its identity to the participants. A general framework for simulation of

185 CRISPR-based lineage recording (**Fig 1B**) (Salvador-Martínez et al., 2019) was used to simulate

186 mutations in a recording array on top of the resulting topology (see **Fig 1C**). In the third

187 challenge, participants had to infer the lineage of cells in a simulated tree of ~10,000 cells (**Fig**

188 **1D**) representing 11 different cell types after one year of *M. musculus* development (**Fig 1E**).

189 Simulating such a large tree was made possible by applying the Environment-dependent

190 Stochastic tree Grammars (eSTGt), a programming and simulation environment for population

191 dynamics (Spiro and Shapiro, 2016) adapted to simulate cell lineages (see STAR methods).

192 While the size of the actual simulated tree is estimated to be about $10^{12}$ or a trillion cells, the final

193 sub-sampled lineage stored information for only 10,000 cells (see **Fig S1**).

194

195 **Experimental *in vitro* dataset**

196       intMEMOIR is a synthetic image-readable lineage recording system that has been

197 recently developed and tested in mouse embryonic stem cells and the brain of *Drosophila*

198 *melanogaster* (Chow et al., 2021). This technology builds upon a previously developed recording

199 system named MEMOIR (Memory by Engineered Mutagenesis with Optical In situ Readout)

200 (Frieda et al., 2017). In its current implementation, intMEMOIR consists of a multi-state

201  memory DNA array that can be edited irreversibly by serine integrases and integrated at defined

202  genomic sites. While MEMOIR's design enabled 2 different states for each recording unit in the

203  memory array, intMEMOIR enables 3 different states. Upon induction by doxycycline, the serine

204  integrase Bxb1 can bind to the editable character array elements or barcodes, and by DNA-

205  recombination mutate the recording element ground state (represented as '1') into either two

206  possible states, a deletion (represented as '0') or an inversion (represented as '2') of the DNA

207  sequence. The recording process is fully stochastic and happens irreversibly at a constant rate, as

208  any element in the array can be edited at any moment. On mouse embryonic stem cells, Chow *et.*

209  *al* showed that lineage information can be recorded irreversibly and stored in the intMEMOIR

210  array, while also read-out using microscopy. From the recorded data, the lineage history can then

211  be inferred (**Fig 1A**).

212      In the experiment, the growth of 106 cell colonies was traced, each one started from an

213  individual cell carrying an unedited 10-character array. Recording was induced for the first 36

214  hours of growth (approximately 3 cell divisions) and cells were then allowed to grow with no

215  further recording for an additional 24 hrs. At this point the arrays for each cell in the colony were

216  read-out using single molecule fluorescent *in situ* hybridization (smFISH). For each colony, the

217  ground truth lineage was obtained from time-lapse movies. As cells grow at different speeds and

218  some of them die, the resulting colonies had a distribution of sizes, from 4 to 39 cells (see **Table**

219  **1**).


220  **Simulated *in silico* datasets**

221      To complement the challenge datasets, data from simulated recording arrays, with

222  respectively 200 Cas9 targets in each cell for *C. elegans* and 1000 targets for *M. musculus*, were

223  generated. Inspired by the GESTALT technique (McKenna et al., 2016a), in the simulations,

224  every cell is represented as a vector of 200 (or 1000) characters, each character representing one

225  Cas9 target. The simulations started with one cell, the fertilized egg, and all its targets in an

226  unmutated ground state represented with "0" (see **Fig 1C**) had the possibility to change to either

227  of 30 different mutational outcomes stochastically as cells divide (see **Box 1**). The initial cell

228  then undergoes a series of cell divisions growing into a population of ~1,000 cells for *C. elegans*

229  and about a trillion cells from which ~10,000 cells are preserved for *M. musculus* (see STAR

230  Methods). The recording array accumulates independent and irreversible CRISPR-induced

231  mutations with a constant probability per time unit, inherited in subsequent cell divisions (see
232  **Box 1**).

233      When a Cas9-induced mutation occurs, the double strand of DNA is broken, which is
234  eventually repaired by the cell. However, in cases where two or more relatively close double
235  strands break before the cell repair machinery can act, the DNA between these breaks can be lost
236  and such events are called an "inter-target deletions". To make these simulations more realistic,
237  we included inter-target deletions affecting 5-10% of the mutation events (see STAR Methods
238  and **Box 1**). We also introduced different probabilities for the different mutational outcomes, in
239  agreement with experimental evidence (McKenna and Gagnon, 2019). Additionally, for the *M.*
240  *musculus* simulations we implemented a 20% data acquisition dropout to reflect the fact that the
241  data acquisition from single cells is rarely perfect (Qiu, 2020) (see **Box 1**). In summary, we
242  introduced experimental parameters where possible in the simulation in order to approximate
243  realistic recording assays.

244  **Training data**

245      As the goal of these challenges was not only to benchmark cell lineage reconstruction
246  algorithms, but also to mobilize a larger community for evaluating new optimal tree-building
247  methods, we provided training data for each challenge. In the *in vitro* challenge, participants
248  were asked to reconstruct the test dataset consisting of 30 cell colonies using only the
249  intMEMOIR array readout, as the ground truth for these lineages was not accessible to the
250  participants. As training set, participants were given array readout data from 76 colonies along
251  with the corresponding ground truth lineages (**Box 1**).

252      For the *in silico* challenges, the training data included the ground truth simulations of 100
253  lineage trees and their mutated array states. These trees comprised 100 cells for *C. elegans* and
254  1000 cells for *M. musculus* generated with the same simulation scheme as for the whole *C.*
255  *elegans* and *M. musculus* trees. The rationale was to test whether training sets composed of
256  smaller trees could still be helpful to fine-tune algorithms then used to reconstruct larger
257  lineages. The *C. elegans* training set tree topology was generated by 100 iterations of pruning
258  and regrafting sub-trees of 100 cells from the whole animal lineage tree (**Box 1**), to preserve
259  some of the initial topology without giving away the origin of the tree. We indeed verified that
260  the aggregation of the 100 trees given for training showed no direct similarity to the 1000 cells

261   *C. elegans* tree. The *M. musculus* training set was obtained using the same eSTG algorithm used

262   for the test dataset but ran for a shorter time in order to obtain smaller trees of 1000 cells.

263   Importantly, the *M. musculus* challenge also had an intermediate step where participants could

264   submit solutions to a ~6000 cell tree and obtain their scoring results on a leaderboard in real

265   time. The leaderboard encouraged participation through competition and provided a way of

266   testing the scalability of the approaches. For scoring, the submitted lineage tree inferences for the

267   test dataset were then compared to their corresponding ground truth using two different metrics

268   (see **Box 2**).

269

270   **Results**

271   **Best performing methods**

272       Overall, the challenge was successful in its main goal to attract a variety of approaches

273   and teams, as twenty-two submissions were received in total for the three challenges. **Figure 2A-**

274   **C** shows the score rankings by both the RF and triplet distances. For the *in vitro* challenge,

275   where nine teams participated, it is clear that the diverse set of approaches reached a plateau in

276   performance for both metrics which suggests that participants successfully extracted and used all

277   available information in the data (**Fig 2A, Fig S2** and **Fig S3A & B** fitted blue line to the

278   medians). We found that the top three teams performed equally well even when calculating the

279   Bayes Factor and an additional quartet metric (**Fig S2**). Interestingly, the two distance-metrics

280   generated different rankings, showing that while correlated the two metrics are not identical. We

281   noted that in general teams performed better on the RF distance compared to the triplet metric

282   (**Fig 2A** and **Fig S3C**). This indicates that for trees less than 100 cells, the triplet metric is more

283   stringent than the whole-tree partitions measured by RF.

284       Five teams submitted solutions for the *C. elegans* and three teams for the *M. musculus*

285   challenge. In both challenges, the distance-based *DCLEAR* method outperformed all other

286   participants. In general, *DCLEAR*'s performance in both challenges and under both metrics was

287   excellent (**Fig 2B** and **2C**) and although the *M. musculus* tree was ten times larger, *DCLEAR*

288   scored higher compared to the *C. elegans* tree.

289   **Summary statistics for the *in vitro* challenge**

290    Given that the *in vitro* challenge predictions consisted of 30 trees of different sizes, we
291    were able to further analyze the results. When considering only perfectly reconstructed trees,
292    defined by a distance value of 0, *AMbeRland\** performed better as we see a larger number of
293    perfect trees when considering triplets (28 trees across teams **Fig. 2D** *top*) than when using RF
294    (21 trees **Fig. 2D** *bottom*). This discrepancy indicates that even when all triplets from a tree are
295    correctly inferred, there might still be incorrect clades in the tree as measured by RF. We then
296    asked whether the different teams performed better depending on the size of the tree, a main
297    constraint for inference accuracy. Larger trees were defined as having more than 8 leaves/cells
298    and small trees as having less or equal to 8 leaves/cells. Irrespective of the tree size,
299    *AMbeRland\** also performed better (see **Fig 2E**). To visualize that indeed tree size has an overall
300    effect on reconstruction accuracy, we plotted the accuracy of individual trees in both metrics
301    colored by the number of cells per tree (**Fig 2F**). Across all trees and submissions, the two
302    metrics correlation is overall high $r$=0.77, but it becomes clear that larger trees generally have a
303    larger triplet distance compared to RF. A total of six trees were reconstructed perfectly by at
304    least one of the teams (**Fig S4**) and we noted that these perfect trees consisted of small trees of
305    less than 9 cells. For these small trees, edit patterns can be slightly redundant without affecting
306    accuracy (e.g. Tree 1 in **Fig S4**) indicating that the size of the tree is a dominant factor in
307    reconstruction accuracy. The largest perfect tree (Tree 20, **Fig S4**) comprises 9 cells with
308    redundant mutations in two array states across cells, despite this, the tree can still be perfectly
309    resolved. More generally, higher redundancy in array states effectively decreases the information
310    that can be used for lineage reconstruction and we indeed observed high levels of redundancy in
311    several trees with an average of 65% ± 20% of cell arrays being unique (**Table S1**). However,
312    tree reconstruction was not affected by this (**Fig S3D & E**). Considering non-perfect trees, the
313    largest tree with the highest score was reconstructed by *AMberRland* (29 leaves/cells, 55%
314    unique arrays RF distance = 0.44 and triplet distance = 0.40, **Fig S5**). The second largest tree
315    with high score was reconstructed by *Cassiopeia* (23 leaves/cells, 71% unique arrays, RF = 0.48,
316    Triplets = 0.70, **Fig S5**). In tree 29 we noted that some cells with identical array states were
317    placed correctly in the reconstruction, this is due to the fact that *AMberLand\** and *Jasper06*
318    decided to leverage the biological restriction that lineage trees must be binary. Therefore, they
319    imposed a binary structure even when cells had identical array states, reaching slightly higher
320    accuracy (**Fig. S5**).

## Methods summary

321

322      The best performing methods across challenges can be roughly divided into three groups:

323      (1) distance-based methods such as the best performers *Liu*'s method, *Guan's* method and

324      *DCLEAR* (2) a machine learning based method to predict probabilities of sister cells using a

325      Gradient Boosting Machine *AMbeRland*, and (3) a maximum parsimony-based method

326      *Cassiopeia-ILP* and *Cassiopeia-Greedy*. The distance-based methods reconstruct the lineage

327      trees by first defining a distance to build a matrix between all pairs of cells as the distance

328      between mutated characters in two cells'arrays should be proportional to the time since they split

329      from a common ancestor. Therefore, distance matrices are commonly used in phylogenetic

330      inference and clustering (Jones et al., 2020) or by hierarchical algorithms that represent the

331      distance matrix as a tree such as in Neighbor-Joining (NJ)(1987). Conversely, the machine

332      learning approach learns from the training set the importance of features/mutations to predict

333      whether two cells are sisters. *Cassiopeia*'s maximum parsimony method reconstructed trees by

334      minimizing the total number of steps required to explain a given configuration of the leaves.

335      Distance-based methods combined with hierarchical clustering overall performed well

336      with the additional advantage of being scalable. Hamming distance is a metric used for

337      phylogenetic analysis where the distance between sequences from two taxa (or cells in this case)

338      is calculated as the number of different sites between the two sequences. While in the traditional

339      Hamming distance, every mutation is assigned the same weight, in lineage recording

340      technologies the editing rates of each array character are generally not uniform (**Fig 3A** and

341      **Box1**), and so, mutations that occur with higher frequency are likely to arise independently in

342      non-related cells, confounding the analysis. Conversely, some edit patterns are unlikely to

343      happen independently and could be informative of a true inheritance event. Therefore, the

344      uneven frequency of array edits suggests that each array element could potentially bring different

345      information about the underlying lineage relations. To calculate the weighted Hamming

346      distances between cells, several teams transformed the initial edited array sites of all cells in the

347      lineage to their observed mutation frequencies and calculated the absolute difference between the

348      arrays of two cells (**Fig 3B**). Tables S2 and S3 include a concise summary of all methods. For the

349      *in vitro* challenge we included the type of parameters or features that different teams estimated

350      from the data, how was the tree built from their estimations and how did they use the training

351      dataset to estimate or learn the different features and parameters (**Table S3**).  For the *in silico*

352 challenges, given the larger scale of the trees, we also show the CPU running time as well as the

353 code accessibility (**Table S3**).

354

355 *Liu***: Inference of internal states.**

356 In all three challenges Team *Liu's* method reconstructed internal nodes to represent the

357 ancestral nodes that likely gave rise to the leave cells. For the *in vitro* challenge, the state of

358 every internal node is inferred using the states of its children by applying the following rule for

359 each site: the parent node gets the state of the children nodes if both children states are the same,

360 alternatively it gets the unedited state if its two children states are different. Next, for each array

361 element, the transition rate from state '1' to state '0' or '2' is calculated as the probability of

362 parent node having state 1 and child node having the mutated state (**Fig 3C** *top*). Finally, the

363 pairwise distance between two cells is considered to be the probability of two cell states arising

364 from independent events, that is, the product of the transition rate of shared states between the

365 two cells. In a similar way for the *in silico* challenges, team *Liu* estimated the character array of

366 the internal nodes based on the fact that a target can only mutate once (**Fig 3C** *middle*). Deletions

367 or d*ropouts* were replaced by the initial character "0". After inferring all the internal nodes, *Liu's*

368 method derived the empirical transition probability from the ground state to the 30 possible

369 mutated states, 'A-Z' and 'a-c' or deletion '-'. This empirical distribution was then used to

370 calculate the probability of two cells arising from two independent events, assuming that each

371 target was independent of the other. The log likelihood of the transition probability for shared

372 states was considered as the cell-to-cell distance. Finally, the distance matrix was clustered using

373 Unweighted Pair Group Method with Arithmetic Mean Algorithm (UPGMA) (**Fig 3C**

374 *bottom*). For the *M. musculus* challenge *Liu's* method added an extra step for clustering taking

375 into consideration the 11 different types of cells.

376 *Guan:* **weighted Hamming distance.**

377 For the *in vitro* challenge *Guan Lab*'s method first designed a rule-based hierarchical

378 clustering method using weighted Hamming distances between cells **(Fig. S6A** for frequency and

379 weight values**)**. *Guan Lab* transformed the initial edited array sites of all cells in the lineage to

380 their observed mutation frequencies while retaining the mutation directions by mathematical

381   signs ($+/-$ see **Fig S6A**) and calculated the weighted distance as the absolute difference

382   between the arrays of two cells. Finally, the lineage was reconstructed using a rule-based

383   hierarchical clustering method (**Fig S6B**). For the *C. elegans* challenge they first replaced all gap

384   mutations with the mutation types at both ends, since gaps even at the same sites could be the

385   result of simultaneous mutation incidents (**Fig 3D**). The mutation weights were defined for each

386   of the 200 characters in the *C. elegans* array as $1\text{-log}_{10}(P)$, where $P$ is the observed probability of

387   the mutation at that site. An iterative bifurcate clustering process was performed to combine the

388   nearest cells based on matrix calibration, until there was only one pair of cells left and their

389   parent cell was defined as the root of the tree (see **Fig 3D**).

*Cassiopeia***: Combinatorial optimization.**

391       *Yosef Lab* was the only team that did not opt for hierarchical clustering but instead, they

392   used combinatorial optimization. For the *C. elegans* challenge, the team adapted the previously

393   published *Cassiopeia-ILP* (Jones et al., 2020) an integer linear programming (ILP) which takes

394   as input a "character matrix," summarizing the mutations seen at heritable target sites across

395   cells (**Fig 3E** *Top*). It then infers a Steiner Tree, finding the tree of minimum weight connecting

396   all observed cell states across all possible ancestral states' histories and maximizes the

397   parsimony over all possible trajectories that could have generated the observed barcode states

398   which consistently finds a near-optimal solution.  Importantly, the edges connecting cell states

399   can be weighted by the number of mutations along that edge or the log-likelihood of these

400   mutations. A derived method *Cassiopeia-Greedy* was implemented for the *M. musculus*

401   challenge also adapted a different maximum parsimony-based strategy to infer the phylogeny

402   from a set of observed character-states across all cells summarized in a cell's x cut-site

403   "character-matrix" (Jones et al., 2020). To do so, the algorithm recursively applied a heuristic to

404   split cells into two groups based on the frequency of a given state at a character and the

405   likelihood of that state arising, taking into account mutations that occurred earlier in the tree (**Fig**

406   **3E** *Bottom*). This procedure was applied until a full lineage tree was resolved.

407

**Usage of the Ground Truth**

For the *in vitro* challenge, several teams computed the calculated transition rates across the 76 trees in the training data and found striking variability across the array element identities and positions (**Fig 3A)**. It is possible to assess in several ways how much information regarding the correct lineage of a cell is contained in the transition rate of a particular mutation. For example, given a tree in the training set it is possible to assess whether cells having the same mutation in an array element are in the same subtree branch (see diagram **Fig 3F**). To obtain the percentage of correct branch positioning associated to this mutation, this process can be repeated for all trees. It can then be expanded to all ten elements in the arrays, and for the two types of mutations (1 to 0 or 1 to 2). This information was used to quantify how for a given mutation and array position there is a negative correlation between the state transition rate and how well it can establish the correct relationships between four cells in a subtree ($R^2=0.58$, see plot **Fig 3F**). This observation is in line with teams assigning the observed mutation frequencies to the Hamming distance weights of different array elements, but also shows that weight values can be further refined when training data is available.

Participant teams used this type of information differently as *Cassiopeia-ILP* (*Yosef Lab*) used the average across sites of the transition probabilities for each type of mutation to weight the edges of their Steiner-Tree search (**Fig 3E *top***). Additionally, for this team the training data also proved useful in choosing a model as they were able to compare the performance of different algorithms and select the one that performed the best (**Fig 3G**). Team *Guan Lab* was able to use the ground truth for comparing several types of distance-based tree construction methods, including Neighbor-Joining (NJ) and UPGMA. This analysis showed that UPMGA performed similar to their rule-based hierarchical clustering whereas NJ was significantly outperformed (**Fig S6C**). Finally, *DCLEAR*(WHD) used the training set to weight the mutations for the *C.elegans* tree and *AMbeRland* used a Gradient Boosting Machine (GBM) to learn the relative importance of several features derived from the array states data and for determining the clustering thresholds for the tree reconstruction (see details below).

**DCLEAR estimates k-mer replacement distances by simulation**

*DCLEAR* (Distance based Cell LinEAge Reconstruction) implemented two best performing

437   strategies to compute the cell distances. A weighted Hamming distance strategy (WHD) that

438   requires a training set for optimizing each mutation weight for the *C. elegans* tree, and a *k*-mer

439   replacement distance (KRD), that does not require training data, for the *M. musculus*

440   tree. *DCLEAR* (KRD) first looks at mutations in the character arrays to estimate the parameters

441   of the generative process associated with the tree to be reconstructed. With these parameters,

442   they repetitively simulated trees with a size and mutation distribution similar to the *M. musculus*

443   target tree (**Fig 4A**). The *k*-mer replacement distances were estimated from the simulated lineage

444   trees and used to compute the distances between input sequences in the character arrays of

445   internal nodes and tips. As a toy example, two cells in a simulated tree have respectively the

446   character arrays A00A and E00C, their *1*-mer nodal distance will be the distance between A and

447   C, their *2*-mer nodal distance will be the distance between 0A and 0C while the whole sequence

448   nodal distance will be between A00A and E00C (see red cells in **Fig 4B**). Specifically, by

449   examining the simulated lineage trees, *DCLEAR* (KRD) estimated the expected *1*-mer

450   replacement distance between characters in the array (including ground state '0' and deletion

451   state '-') in the lineage trees (**Fig 4C**) and the probability for a given nodal distance of replacing

452   a character in a cell array (**Fig 4D and 4E**). To extend the *1*-mer replacement distance to the *k*-

453   mer replacement distance, the posterior probability distributions of *k*-mer replacement distance

454   were estimated by using a conditional model considering a dependance for the concurrence of

455   mutations (**Fig 4F** and **4G**). They found that by considering the neighboring characters, the

456   conditional model can more accurately estimate the nodal distance than an independent *1*-mer

457   model. The cell distance can then be readily computed as the mean expected *k*-mer replacement

458   distance (see STAR Methods). Similar to WHD, the lineage trees were reconstructed using the

459   Minimum Evolution (FastME) or Neighbor-Joining (NJ) algorithms (Gascuel and Steel, 2006;

460   Lefort et al., 2015). For both *DCLEAR* WHD and KRD, the deletions and dropouts were treated

461   differently. In WHD, the weight for deletion, dropout, regular state and ground state are 0.9, 0.4,

462   3 and 1, respectively. In KRD, deletion and dropout are treated as two different characters.

463   ***Amberland*, a decision tree-based method**

464   *AMbeRland's* approach relied on machine-learning to build a distance matrix between

465   cells through the calculation of the relative importance of features derived from the states of the

466   character arrays (**Fig 5**). In their approach for the *in vitro* challenge, they first defined four

467  features for every pair of cells consisting of whether two cells are both unedited at a given array
468  site (feature F1), a site has the same edits (feature F2), only one site is edited (feature F3) or if
469  both sites have different edits (feature F4) (**Fig 5A** *left*). Then, the prevalence of these four
470  features was extracted for a group of ~500 pairs of sister cells (label 1) and ~3000 non-sister
471  cells (label 0) using the 76 ground truth trees available in the training set. Finally, Gradient
472  Boosting (Friedman, 2001) was applied to learn from this data the relative weights of each
473  feature to predict whether two cells are actually sisters (see **Fig S7**). For the *C. elegans* challenge
474  *AMbeRland* applied a similar approach using the training set of 100 trees with 100 cells. They
475  similarly determined weights for features selected by counting pairwise positions in two cell's
476  arrays that were (1) not mutated, (2) had a single mutation, (3) both had different mutations, (4)
477  both had a missing record, (5) one had a missing record and the other not mutated, etc. (**Fig 5A**
478  *right*).
479      In both challenges, *AMberland* applied a custom hierarchical clustering method for
480  building the cell lineage tree from the predicted probabilities. During the tree construction, the
481  ground truth was used to evaluate a set of decreasing thresholds corresponding to how any two
482  individual clusters of cells were related at different levels of the lineage tree (see **Fig 5B** *left*).
483  The clustering starts at the lowest tree level, where all cell pairs are ordered according to the
484  predicted probability that they are sister cells, from here, cells with a probability higher than the
485  first threshold are assigned as pairs, while the rest are kept as a branch with a single cell. At each
486  consecutive level, pairwise comparison are performed between each lower level cluster by
487  calculating the maximum probability between any two elements of the two clusters. Pairs of
488  clusters were ordered again according to this probability and were assumed to have the same
489  parent node if their value was above the estimated threshold for this level. This process was
490  repeated until one or two clusters were left. The values for the thresholds at each level were
491  determined by performing a grid search minimizing the RF and triplet distance metrics (see
492  results for tree 29 **Figure 5B** *right*)**.** This procedure clearly helped obtaining better scores,
493  particularly regarding the triplet metric (see **Fig S8** for all trees in the *in vitro* challenge and **Fig**
494  **S9** for *C. elegans*).

**Consensus trees**

495

496     One advantage of having a set of different and diverse approaches trying to solve a

497     common problem is that it is possible to aggregate the solutions and gather collective insight.

498     Hence, we decided to test how a consensus tree of all teams would perform compared to

499     individual methods (**Fig 6A&D**). For the *in vitro* challenge*, w*e constructed the consensus tree

500     using the submissions from all teams (excluding *Bengal Tiger* because of their unusual number

501     of low-accuracy outliers, **Fig 2A**) by applying the majority-rule algorithm (Felsenstein, 1985).

502     Interestingly, we see that the consensus tree performs better than any individual team when

503     considering the RF distance, but this is not the case according to the triplet distance (**Fig 6B**). To

504     further understand this, we evaluated the agreement (or support) of each clade in a given tree

505     across teams using the Felsenstein's Bootstrap Proportion (FBP), which has been traditionally

506     used to assess the support of phylogenetic trees (Felsenstein, 1985). For FBP agreement, a

507     branch must match a reference branch exactly to be accounted for in the score, so we define FBP

508     as a strict agreement (**Fig 6A**). Alternatively, the Transfer Bootstrap Expectation (TBE) provides

509     higher resolution estimates of branch support and can be used to assess phylogenetic similarity

510     even when there is no strict majority consensus (Lemoine et al., 2018). The distribution of FBP

511     and TBE support scores at different normalized depths across all 30 trees in the test dataset

512     shows that the inference of earlier clades varies significantly across methods, whereas late splits

513     are resolved correctly by the majority (**Fig 6A and Fig S10**). The divergence for earlier clades

514     might explain the lower performance of the consensus tree under the triplet metric, given that for

515     these small trees more triplets are prone to include early divisions with wrong clade relationships

516     (see **Fig S3A&B**).

517     For the *in silico* challenges we also added for comparison the performance of the

518     algorithm *FastTree2*, a fast and reliable approximately-maximum-likelihood method (Price et al.,

519     2010) that performed better than neighbor joining or TripleMaxCut (Sevillya et al., 2016).

520     Interestingly, we observed that in the *C. elegans* challenge, *DCLEAR* outperforms *Fastree2* by

521     both metrics, which is not the case for the *M. musculus* challenge as *FastTree2* outperforms all

522     methods, with *DCLEAR* as a close second (**Fig 6C**). We also see that for the *C. elegans*

523     challenge, the consensus tree performs better than any individual team when considering the RF

524     distance, but under the triplet distance the consensus is nevertheless equivalent to a random

525  submission (**Fig 6C**). In the *M. musculus* challenge there were probably not enough submissions
526  to see a "wisdom of the crowds" effect as the consensus tree does not outperform *DCLEAR*. To
527  understand the difference between the RF and triplet distances, we evaluated the agreement of
528  each clade in the *C. elegans* tree across teams. Overall, as in the *in vitro* challenge we observed a
529  depth-dependent effect in the support between teams, as measured by TBE (**Fig 6D**) and the
530  divergence for earlier clades might explain the lower triplet metric performance in the consensus
531  tree solution but in this case probably due to the *C. elegans* tree topology having many internal
532  nodes.

## Discussion

534       The main goal of this DREAM challenge was to mobilize a larger community to generate
535  new methods for cell lineage reconstruction. This goal was catalyzed through the generation of
536  new *in silico* datasets and by the recent availability of *in vitro* datasets with an associated ground
537  truth. This study represents the first attempt to rigorously examine the performance of various
538  algorithms across diverse molecular tools and lineage trees. For the *in vitro* challenge a total of
539  nine approaches were submitted for which the maximum performance plateaued (see **Fig 2A** and
540  **Table S2**). While some trees were reconstructed perfectly, the scores were far from the
541  theoretical maximum. We thought this could be mainly due to the high degeneracy in cell arrays
542  where two or more cells show identical edit patterns, but further analysis showed that barcode
543  degeneration did not affect the performance of the teams (**Fig S3E**). This problem could be in
544  principle overcome by increasing the memory of the intMEMOIR system, as discussed by the
545  authors (Chow et al., 2021). On the other hand, the degeneracy problem was non-existent for the
546  *C. elegans* tree as all cells ended up with a different mutational character array and was minimal
547  for *M. musculus* with only ~2.7% of sister cells sharing exactly the same character arrays.
548  Indeed, the choice of the mutation rate and the diversity of mutations in the simulations has a
549  strong effect on the accuracy of cell lineage reconstruction as low diversity of possible
550  mutational outcomes generally gives poorer results. While too low mutation rates lead to more
551  unedited and therefore non-informative targets, too high mutation rates lead to most targets being
552  mutated during the early cell divisions, leaving few targets available for recording later events
553  (Salvador-Martínez et al., 2019). Hence, we tuned our *in silico* mutation rates and array sizes in
554  order to avoid cells having identical character arrays. As the performance of *DCLEAR* in the *in*

555 *silico* challenges was as good or even better than the results of the *in vitro* challenge (see **Fig 2**),
556 the limits of its performance must derive from the tree size or topology. We conclude that tree
557 topology was the most important parameter given that *DCLEAR M. musculus* lineage
558 reconstruction was more accurate than for the ten times smaller *C. elegans* tree. Given these
559 great performance, we also consider the *in silico* challenges a success despite not having as many
560 submissions, as the diversity and performance of the approaches was impressive (see **Fig 2** and
561 **Table S3**).

562       The implementation of several metrics to evaluate the participants was also an original
563 feature of the challenge as in general, lineage trees are evaluated with a single metric and no
564 comparison between metrics is systematically performed (Salvador-Martínez et al., 2019). This
565 aspect was essential not only to thoroughly evaluate participants (**Fig 2, S2** and **S3**) but also to
566 better understand their solutions. One of the striking observations was the disconnection in all
567 challenges of the performance as measured by the two metrics. Indeed, for the *in vitro* challenge
568 *AMberRland* optimized post competition their algorithm for the triplet distance and had the
569 overall best performance without compromising their RF performance (**Fig 3A**). Also, for larger
570 trees, team *AMbeRland\** had overall a similar performance than *Cassiopeia* relative to the triplet
571 distance (average triplets = 0.55) but scores better in the RF metric (RF = 0.57 and 0.65
572 respectively, see **Fig 2E**). We see the opposite for team *philrennert* although now the difference
573 for larger trees appears for the triplet distance (triplets = 0.57 and 0.72 respectively, see **Fig 2E**)
574 as the RF distance is similar. Such dissociation between metrics was also observed for the
575 majority-vote consensus solution which had the best score for RF but far from that for triplet
576 distance (**Fig 6B**). The analysis of the overall agreement between individual solutions at different
577 depths of the trees shows that indeed for earlier cell divisions agreement is low (**Fig 6A**). This
578 observation provides a possible explanation for the divergence between triplet and RF distances,
579 as in smaller trees such as the ones in the *in vitro* challenge, more triplets are prone to include
580 early divisions with wrong clade relationships, bringing down the triplet performance.
581 *AMberRland* was probably able to correct this by performing a grid search and changing the
582 thresholds for hierarchical clustering at higher levels of the tree. As *AMberRland* was also the
583 method that most consistently predicted smaller and larger trees (**Fig 2D, S4 and S5)**, this also

584 explains why we observed that overall the triplet distance is higher than RF in larger trees as
585 opposed to smaller trees (**Fig 2F** and **S3C**).

586       For the much larger trees in the *in silico* challenges the interpretation of the metrics is
587 different as the number of triplets included in the triplet distance grows cubically with the size of
588 the tree, while the number of partitions considered by the RF distance grows linearly. Hence, for
589 larger trees, the triplet distance will be dominated by the higher number of triplets close to the
590 tree leaves as the RF distance will be mostly measuring major branching events in the early cell
591 division stages. As *DCLEAR* was consistently better in both metrics, but scored less favorably in
592 RF distance, compared to the triplet distance, this suggests that *DCLEAR* is precisely having
593 trouble detecting those major branching events. Indeed, both WHD and KRD in *DCLEAR*
594 methods rely on the rare mutations to estimate the cell distances. During early cell division
595 stages, however, the rare mutations are significantly less likely to be present in the sequences and
596 result in difficulties for separating early branching events. Modeling the dependence between
597 multiple non-adjacent mutations in the sequences, on top of the neighboring *k-mers*, may be
598 necessary to more accurately evaluate the early branching events. It is also striking to see how
599 the maximum parsimony approach of *Cassiopeia* scored much better for the triplet distance for
600 larger trees in all challenges. Finally, the machine learning approach derived from the one
601 applied in the *C. elegans* challenge by *AMberRland* was able to perform acceptably in the RF
602 metric with much larger trees (see **Fig 2B**), but although the threshold optimization worked for
603 the training set of 100 cell trees (see **Fig S9**), it did not do well with the triplet distance of the *C.*
604 *elegans* tree probably due to the need to include many more thresholds given its 10 times larger
605 size.

606       The final observation regarding the metrics discrepancy is related to the performances in
607 the training and test sets of the *C.elegans* challenge, as all teams are similar regarding the RF
608 distance but with the exception of *DCLEAR* and *Cassiopeia*, the triplets performance is worse for
609 the test set than in the training set (see **Box 1**). Conversely, for the *M.musculus* challenge their
610 performances in the leaderboard tree of ~6500 cells and the *M.musculus* tree of ~10,000 cells
611 match for both metrics (**Fig 6C**). We conclude that when reconstructing a cell lineage tree, the
612 results obtained with an algorithm for a training set of trees with a number of leaves an order of

613     magnitude smaller than the test set are comparable, although the triplet distance is more unstable

614     than the RF distance.

615         Regarding the generalization of the results obtained with the intMEMOIR technology

616     which is difficult to compare at the molecular level to the sequence-based approaches for lineage

617     reconstruction as it also shows differences such as the absence of accidental deletions or

618     *dropouts*, we think that in conjunction with the results from the *in silico* approaches, the

619     generalizable conclusions are the necessity of having well calibrated mutation rates to avoid too

620     little mutations but also array degenerations, the utility of having a training set of smaller trees to

621     optimize lineage reconstruction methods including distances and clustering, and allowing for a

622     clear interpretation of the effect of the two different metrics with different tree sizes.

623         Overall, we think that the decisions taken while producing the datasets for the *in silico*

624     challenges were the correct ones. We were able to pose a problem that we think is close enough

625     to a biological situation and difficult enough so that the lessons learned and solutions generated

626     can be implemented in other contexts. Indeed, it has been estimated that under ideal conditions

627     of optimized mutation rates, uniform cell divisions and fully sequenced targets, 30 targets should

628     be sufficient to reach a high level of accuracy for the lineage reconstruction of a tree of about

629     65,000 cells (Salvador-Martínez et al., 2019). In this situation 100 targets would theoretically

630     yield almost perfect accuracy, far from the results obtained by the solutions submitted to both

631     challenges.

632         Finally, as new DNA-editing-based molecular tools promise the reconstruction of single-

633     cell lineages from complex model organisms, including the human cell lineage, an important

634     question is whether the access to smaller trees and the molecular data from their cell lineages

635     could help find solutions to be implemented for larger trees of the same origin. The *M. musculus*

636     lineage tree being the current experimental frontier for lineage reconstruction(Bowling et al.,

637     2020; Kalhor et al., 2018), our results show that indeed, in order to obtain an accurate full cell

638     lineage for mouse or human, it could be possible to train algorithms on smaller trees obtained

639     from organs (Bowling et al., 2020) or *in vitro* dividing cells and these can then be implemented

640     for building algorithms that can then be applied to the reconstruction of much larger trees.  This

641     DREAM challenge was a first attempt to rigorously examine the performance and robustness of

642    various algorithms under the same conditions. It took advantage of the unique opportunity to use

643    unpublished datasets of molecular and simulated character arrays. We hope that showing that

644    machine learning methods can indeed be successfully implemented will pave the way for other

645    benchmarking efforts based on emerging technologies for monitoring cell lineages and the

646    application of new algorithmic approaches, but also that the approaches described here will pave

647    the way for the solution of the mouse and human cell lineages.


648


649

## References

Alemany, A., Florescu, M., Baron, C.S., Peterson-Maduro, J., and Oudenaarden, A. van (2018). Whole-organism clone tracing using single-cell sequencing. Nature *556*, 108–112.

Becattini, S., Latorre, D., Mele, F., Foglierini, M., Gregorio, C.D., Cassotta, A., Fernandez, B., Kelderman, S., Schumacher, T.N., Corti, D., et al. (2015). Functional heterogeneity of human memory CD4$^+$ T cell clones primed by pathogens or vaccines. Science *347*, 400–406.

Behjati, S., Huch, M., Boxtel, R. van, Karthaus, W., Wedge, D.C., Tamuri, A.U., Martincorena, I., Petljak, M., Alexandrov, L.B., Gundem, G., et al. (2014). Genome sequencing of normal cells reveals developmental lineages and mutational processes. Nature *513*, 422–425.

Bowling, S., Sritharan, D., Osorio, F.G., Nguyen, M., Cheung, P., Rodriguez-Fraticelli, A., Patel, S., Yuan, W.-C., Fujiwara, Y., Li, B.E., et al. (2020). An Engineered CRISPR-Cas9 Mouse Line for Simultaneous Readout of Lineage Histories and Gene Expression Profiles in Single Cells. Cell *181*, 1410-1422.e27.

Chan, M.M., Smith, Z.D., Grosswendt, S., Kretzmer, H., Norman, T.M., Adamson, B., Jost, M., Quinn, J.J., Yang, D., Jones, M.G., et al. (2019). Molecular recording of mammalian embryogenesis. Nature *570*, 77–82.

Chow, K.-H.K., Budde, M.W., Granados, A.A., Cabrera, M., Yoon, S., Cho, S., Huang, T., Koulena, N., Frieda, K.L., Cai, L., et al. (2021). Imaging cell lineage with a synthetic digital recording system. *Science* 327, eabb3099.

Evrony, G.D., Lee, E., Mehta, B.K., Benjamini, Y., Johnson, R.M., Cai, X., Yang, L., Haseley, P., Lehmann, H.S., Park, P.J., et al. (2015). Cell Lineage Analysis in Human Brain Using Endogenous Retroelements. Neuron *85*, 49–59.

Felsenstein, J. (1985). CONFIDENCE LIMITS ON PHYLOGENIES: AN APPROACH USING THE BOOTSTRAP. Evolution *39*, 783–791.

Frieda, K.L., Linton, J.M., Hormoz, S., Choi, J., Chow, K.-H.K., Singer, Z.S., Budde, M.W., Elowitz, M.B., and Cai, L. (2017). Synthetic recording and in situ readout of lineage information in single cells. Nature *541*, 107–111.

Frumkin, D., Wasserstrom, A., Kaplan, S., Feige, U., and Shapiro, E. (2005). Genomic Variability within an Organism Exposes Its Cell Lineage Tree. Plos Comput Biol *1*, e50.

Garcia-Marques, J., Espinosa-Medina, I., Ku, K.-Y., Yang, C.-P., Koyama, M., Yu, H.-H., and Lee, T. (2020). A programmable sequence of reporters for lineage analysis. Nat Neurosci *23*, 1618–1628.

694    Gascuel, O., and Steel, M. (2006). Neighbor-Joining Revealed. Mol Biol Evol *23*, 1997–2000.

695    Jones, M.G., Khodaverdian, A., Quinn, J.J., Chan, M.M., Hussmann, J.A., Wang, R., Xu, C.,
696    Weissman, J.S., and Yosef, N. (2020). Inference of single-cell phylogenies from lineage tracing
697    data using Cassiopeia. Genome Biol *21*, 92.

698    Kalhor, R., Kalhor, K., Mejia, L., Leeper, K., Graveline, A., Mali, P., and Church, G.M. (2018).
699    Developmental barcoding of whole mouse via homing CRISPR. Science *361*, eaat9804.

700    Kebschull, J.M., and Zador, A.M. (2018). Cellular barcoding: lineage tracing, screening and
701    beyond. Nat Methods *15*, 871–879.

702    Kester, L., and Oudenaarden, A. van (2018). Single-Cell Transcriptomics Meets Lineage
703    Tracing. Cell Stem Cell *23*, 166–179.

704    Kretzschmar, K., and Watt, F.M. (2012). Lineage Tracing. Cell *148*, 33–45.

705    Lefort, V., Desper, R., and Gascuel, O. (2015). FastME 2.0: A Comprehensive, Accurate, and
706    Fast Distance-Based Phylogeny Inference Program. Mol Biol Evol *32*, 2798–2800.

707    Lemoine, F., Entfellner, J.-B.D., Wilkinson, E., Correia, D., Felipe, M.D., Oliveira, T.D., and
708    Gascuel, O. (2018). Renewing Felsenstein's phylogenetic bootstrap in the era of big data. Nature
709    *556*, 452–456.

710    Livet, J., Weissman, T.A., Kang, H., Draft, R.W., Lu, J., Bennis, R.A., Sanes, J.R., and
711    Lichtman, J.W. (2007). Transgenic strategies for combinatorial expression of fluorescent
712    proteins in the nervous system. Nature *450*, 56–62.

713    Lodato, M.A., Woodworth, M.B., Lee, S., Evrony, G.D., Mehta, B.K., Karger, A., Lee, S.,
714    Chittenden, T.W., D'Gama, A.M., Cai, X., et al. (2015). Somatic mutation in single human
715    neurons tracks developmental and transcriptional history. Science *350*, 94–98.

716    McKenna, A., and Gagnon, J.A. (2019). Recording development with single cell dynamic
717    lineage tracing. Development (Cambridge, England) *146*, dev169730.

718    McKenna, A., Findlay, G.M., Gagnon, J.A., Horwitz, M.S., Schier, A.F., and Shendure, J.
719    (2016a). Whole-organism lineage tracing by combinatorial and cumulative genome editing.
720    Science *353*, aaf7907.

721    McKenna, A., Findlay, G., Gagnon, J.A., Horwitz, M., Schier, A.F.F., and Shendure, J. (2016b).
722    Whole organism lineage tracing by combinatorial and cumulative genome editing.

723    Perli, S.D., Cui, C.H., and Lu, T.K. (2016). Continuous genetic recording with self-targeting
724    CRISPR-Cas in human cells. Science *353*, aag0511.

725  Price, M.N., Dehal, P.S., and Arkin, A.P. (2010). FastTree 2--approximately maximum-
726  likelihood trees for large alignments. Plos One *5*, e9490.

727  Qiu, P. (2020). Embracing the dropouts in single-cell RNA-seq analysis. Nat Commun *11*, 1169.

728  Raj, B., Wagner, D.E., McKenna, A., Pandey, S., Klein, A.M., Shendure, J., Gagnon, J.A., and
729  Schier, A.F. (2018). Simultaneous single-cell profiling of lineages and cell types in the vertebrate
730  brain. Nat Biotechnol *36*, 442–450.

731  Robinson, D.F., and Foulds, L.R. (1981). Comparison of phylogenetic trees. Math Biosci *53*,
732  131–147.

733  Saez-Rodriguez, J., Costello, J.C., Friend, S.H., Kellen, M.R., Mangravite, L., Meyer, P.,
734  Norman, T., and Stolovitzky, G. (2016). Crowdsourcing biomedical research: leveraging
735  communities as innovation engines. Nat Rev Genet *17*, 470–486.

736  Salcedo, A., Tarabichi, M., Espiritu, S.M.G., Deshwar, A.G., David, M., Wilson, N.M., Dentro,
737  S., Wintersinger, J.A., Liu, L.Y., Ko, M., et al. (2020). A community effort to create standards
738  for evaluating tumor subclonal reconstruction. Nat Biotechnol *38*, 97–107.

739  Salvador-Martínez, I., Grillo, M., Averof, M., and Telford, M.J. (2019). Is it possible to
740  reconstruct an accurate cell lineage using CRISPR recorders? Elife *8*, e40292.

741  Salvador-Martínez, I., Grillo, M., Averof, M., and Telford, M.J. (2020). CeLaVi: An Interactive
742  Cell Lineage Visualisation Tool. Biorxiv 2020.12.14.422765.

743  Sevillya, G., Frenkel, Z., and Snir, S. (2016). Triplet MaxCut: a new toolkit for rooted supertree.
744  Methods Ecol Evol *7*, 1359–1365.

745  Spanjaard, B., Hu, B., Mitic, N., and Junker, J.P. (2017). Massively parallel single cell lineage
746  tracing using CRISPR/Cas9 induced genetic scars. BioRxiv 205971.

747  Spanjaard, B., Hu, B., Mitic, N., Olivares-Chauvet, P., Janjuha, S., Ninov, N., and Junker, J.P.
748  (2018). Simultaneous lineage tracing and cell-type identification using CRISPR–Cas9-induced
749  genetic scars. Nat Biotechnol *36*, 469–473.

750  Spiro, A., and Shapiro, E. (2016). eSTGt: a programming and simulation environment for
751  population dynamics. Bmc Bioinformatics *17*, 187.

752  Sugino, K., Garcia-Marques, J., Espinosa-Medina, I., and Lee, T. (2019). Theoretical modeling
753  on CRISPR-coded cell lineages: efficient encoding and optimal reconstruction. Biorxiv 538488.

754  Sulston, J.E., and Horvitz, H.R. (1977). Post-embryonic cell lineages of the nematode,
755  Caenorhabditis elegans. Dev Biol *56*, 110–156.

756  Wagner, D.E., and Klein, A.M. (2020). Lineage tracing meets single-cell omics: opportunities
757  and challenges. Nat Rev Genet *21*, 410–427.

758  Weissman, T.A., and Pan, Y.A. (2015). Brainbow: New Resources and Emerging Biological
759  Applications for Multicolor Genetic Labeling and Analysis. Genetics *199*, 293–306.

760  (1987). The neighbor-joining method: a new method for reconstructing phylogenetic trees. Mol
761  Biol Evol.

762
763

**Figure Legends**

765

766 **Box 1 Training set**
767 One of the main goals of this challenge was to provide participants with a training set composed
768 of several trees, their cell's character arrays and the gold standard tree solution. This allowed
769 participants to train or optimize their methods.
770 **A.** In the *in vitro* experiments to obtain mouse stem cell lineages, mutations were induced for the
771 first 36 hrs of growth (approximately 3 cell divisions) and cells were then allowed to grow with
772 no further changes in the recording arrays for an additional 24 hrs. For all these cells the final
773 values (unmodified encoded as 1, inverted encoded as 2 or deleted encoded as 0) of the 10
774 character arrays were obtained by smFISH, while cell divisions were tracked by video-
775 microscopy (see **Table S1**). Two partitions were created from the original unpublished dataset
776 containing 106 lineages, which represent sufficient experimental data to extract a training set: the
777 training partition composed of 76 trees was provided for the teams along with the corresponding
778 ground truth lineages, for the test partition composed of 30 trees only the cells character arrays
779 were provided without ground truth. The partitions were defined to have similar tree size
780 distribution, given that the lineages were composed of a different number of cells depending on
781 the cell division and survival rates, shown in *middle* histogram panel. Also, a similar median RF
782 score distribution between the two data sets when using a maximum-likelihood method described
783 in *Chow et at* was used as partition criteria, see *bottom* panel. **B.** For the *in silico* challenges, both
784 character arrays for the training and test sets were simulated in a similar way. The type of Cas9-
785 induced mutations consisted of 32 characters 'A' to 'Z' and 'a' to 'e' and character deletion '-'.
786 The characters represent DNA targets for Cas9 but no specific relationship with actual DNA
787 sequences was established. The starting character was '0' and the probability of mutating to one
788 of the 30 characters or of being deleted (insertions were not considered) followed in alphabetical
789 order the Gamma probability distribution used to sample the mutations, shown in blue, and in red
790 a fit on the histogram of the actual results. Mutations are irreversible, once a target is mutated, it
791 can no longer change, either to revert to the unmutated state or to transit to a new state. **C.** Inter-
792 target deletions were simulated for both *in silico* challenges where *C. elegans* arrays were
793 composed of one hundred characters and *M. musculus* of one thousand characters. When a Cas9-
794 induced mutation occurs, the double strand of DNA is broken, which is eventually repaired by
795 the cell. However, in cases where 2 or more relatively close double strands break before the cell
796 repair machinery can act, the DNA between these breaks can be lost this is known as an "inter-
797 target deletion". We implemented these so that when two mutations occur in close targets (less
798 than 20 targets apart in the recording array) within a short interval of time during a given cell
799 division, all the targets between them are removed. In these simulations, 5-10% of targets are
800 missing due to inter-target deletions. **D.** Acquisition dropout distributions were implemented
801 only for the *M. musculus* challenge. In order to capture the variability of the signal quality in
802 both the individual samples and the different sites we modeled the 'sequencing dropout' of single
803 cell samples by assigning distinct coverage factors for each sample and for each locus. The
804 density of cell coverage factors P = ($pi$: i = 1 to M) is the probability of obtaining a signal in each
805 sample or and the density of site coverage factors Q = ($qj$: j = 1 to N) as the probability of
806 obtaining a signal in each locus. The probability of obtaining a signal in sample i and locus j thus
807 equals $pi.qj.r$. Those are multiplied to get the individual coverage factor of a specific site in a
808 specific cell, finally deriving the acquisition dropout status as a factor of a global coverage

809 parameter *r*. **E.** We provided 100 training cell lineage trees of 100 cells for *C. elegans* and of
810 ~1000 cells for *M. musculus.* As the *C. elegans* tree has been experimentally solved, its topology
811 was used to generate the training set. The *M. musculus* tree being completely synthetically
812 generated, the training set was obtained by simply running shorter simulations to obtain ~1000
813 cells trees instead of the ~10,000 cells tree for the test set. **F.** *Top*. We extracted the *C. elegans*
814 training set from its tree topology by cutting and pasting subsets of tree branches. We followed
815 the indicated schematic of cutting and pruning 100 times subsets of the whole tree. Note only
816 one prune and regraft event is shown in red in the diagram. From the obtained topology, the
817 mutation arrays were generated from the Gamma distribution and then 100 cells were sampled.
818 This process was repeated 100 times to obtain a full training set. *Bottom* The boxplots show the
819 performance of each submitted method for inferring the lineage trees from 100 training lineages
820 used in the *C. elegans in vitro* challenge. The similarity between the inferred trees and the
821 ground truth trees was measured by Robinson-Foulds distance *left* and Triplet distance *right*. Red
822 stars indicate the score for the *C. elegans* 1000 cell tree. The values for the *M. musculus* training
823 set, were not established due to excessive computational time required.
824

825

826

**Box 2 Scoring approach**
827
828  We applied two widely used metrics for tree comparison: the Robinson-Foulds distance and the
829  triplets distance. While both metrics are applied to assess tree similarities there is no clear
830  agreement as to which one is more relevant for lineage trees. We decided to use both metrics as a
831  way of evaluating their correlation and the insight they provide about the lineage relationships.
832  The Robinson-Foulds distance is commonly defined as the number of partitions shared by a pair
833  of trees across all possible partitions. A partition refers to any cut in the internal branches of a
834  tree that would generate two sub-trees containing complementary leaves. Since the ground-truth
835  and the inferred lineage contain in total the same set of leaves, we can define a shared partition if
836  there is a way to cut both the inferred and ground-truth trees such that the resulting sub-trees
837  share the same sets of leaves. We obtain the RF distance by normalizing to the maximum
838  possible distance of 1, when there are no shared partitions by the trees (Robinson and Foulds,
839  1981). On the other hand, the triplet distance enumerates all possible combinations of three
840  leaves and their corresponding lineage relationship in both the ground truth and the inferred
841  trees. One then counts the number of shared triplets and normalizes by the total possible number
842  of triplets to obtain the triplet distance. For both metrics, a distance value of 0 means that the
843  ground truth and inference trees are identical under the specific criteria while a distance value of
844  1 means that the inference is comparable to a random guess on the tree structure. Overall, the
845  Robinson-Foulds metric detects main branching events, while the triplet metric is a better
846  measure of local branching events.
847  We here present an illustrative example with left the ground truth and *right* the predicted tree. In
848  this case, the tree has three possible partitions *top right* and ten possible triplets *bottom left*. Since
849  1 out of three partitions was incorrect the RF distance is 1/3 or 0.66. Similarly, 4 out of 10
850  triplets were incorrect for a triplet distance of 4/10 or 0.4. Higher distance implies more
851  differences between the ground truth and the inference and therefore a lower score. As observed
852  in the results of this challenge, the relationship between the two metrics will depend on the tree
853  topology but also on the tree size. Indeed the number of triplets will size as the cube of the
854  number of nodes, while the RF partitions will scale linearly with the number of nodes.
855

856

857

858  **Figure 1. Three challenges for lineage reconstruction from experimental and in silico**
859  **generated character arrays.**   **A.** Challenge consisting of reconstructing *in vitro* growing cell
860  lineages. The lineage tracing intMEMOIR system consists of a character array of editable DNA
861  elements –or barcodes- and the integrase enzyme Bxb1. A mouse stem cell line was engineered
862  with both components. A recording event happens when the integrase stochastically edits one of
863  the 10 elements in the array, resulting in two possible outcomes, deletion and inversion (blue and
864  red squares). As cells divide, each individual daughter cell acquires unique edit patterns (right
865  panel). Finally, *in situ* readouts by smFISH enables the extraction of recorded data for individual
866  cells. Since the whole experiment is done under a microscope, a ground truth lineage tree is also
867  generated which we use as our ground truth. **B.** Diagram showing the simulations performed to
868  generate the character arrays for the two *in silico* datasets, an initial cell with N multiple targets
869  (200 or 1000 for the *C. elegans* or *M. musculus* challenge respectively) accumulates one of the
870  30 independent mutations with a given probability, which are inherited in subsequent cell
871  divisions. The pattern of mutations accumulated in each cell is used to infer the lineage tree. **C.**
872  *In silico* challenge consisting of reconstructing the ~1000 cells *C. elegans* cell lineage from the
873  simulated cell character arrays. For visualization purposes the ground truth cell lineage shows
874  only the first 9 cell divisions. **D.** Challenge consisting of reconstructing ~10,000 cells from a
875  simulated *M. musculus* cell lineage developmental tree generated using Stochastic Tree Grammar
876  (STG). The tree simulation describes the early stages of mouse development up to the three germ
877  layers (Mesoderm, Ectoderm and Endoderm are highlighted with colors in the equations and
878  resulting tree), those in turn continue to differentiate to the final populations of about $10^{12}$ cells
879  and 11 cell types simulated in the challenge. **E.** Displayed is a simulation example of the ground
880  truth tree for a subset of cells from the Mesoderm and Ectoderm, highlighted with the respective
881  colors, throughout 1 year of development. The edges width and color reflect the hypergeometric
882  score of its descending leaves.

883    **Figure 2**.  **Analysis of challenge results.  A.** Average performance across 30 lineages of all
884    teams by both triplets and RF metrics for the *in vitro* challenge. **B.** Average bootstrapped
885    performance of all teams by both triplets and RF metrics for the *C. elegans in silico* challenge. **C.**
886    Average bootstrapped performance of all teams by both triplets and RF metrics for the *in silico*
887    *M. musculus* challenge **D.** Number of perfectly reconstructed lineages for each team in the *in*
888    *vitro* challenge. **E.** We partitioned the *in vitro* challenge test data into large (more than 8 cells)
889    and small (less or equal to 8 cells) trees, to assess performance by tree size. **F.** The scores for the
890    two metrics of all 30 trees for all 9 teams for the *in vitro* challenge are plotted against each other
891    and color coded depending on the size of the tree. Deep blue dots, small trees #cells<10, gray
892    blue dots, trees with 10<cells<20, light blue dots, trees with #cells>20. Scores show a general
893    correlation r=0.77 between the two metrics, but also significant dispersion especially for larger
894    trees.
895

896 **Figure 3. Different approaches for solving lineage trees and using the training data**.
897 **A.** In the *in vitro* challenge, the transition rates from the unedited state (*1*) to either of the two
898 edited states (*0, 2*) can be learned directly from the training data, the probabilities for all possible
899 transitions at each of the ten array positions are shown as extracted from the training set. **B.** The
900 schematic shows that when computing the sequence distances, instead of assigning equal weight
901 to different character replacement as in Hamming distance, the weighted Hamming distance
902 assigns different weights to different character replacements. **C.** Description of *Liu lab*'s method
903 in all 3 challenges. First, for the *in vitro* challenge, the transition probability is calculated by
904 counting the frequency of every state transition from parent node to child node. For the *in silico*
905 challenges the transition probability for all character arrays is extracted. Next, the pairwise cell
906 distance is defined as the likelihood of two cells' states arising from two independent events.
907 Finally, the cell lineage is reconstructed from the distance matrix using the UPGMA method. **D.**
908 This schematic shows the *Guan Lab*'s method used to reconstruct the *C. elegans* tree.  First, all
909 gap mutations are remarked based on mutation types at both ends, since gaps, even at the same
910 sites, could be the results of different mutation incidents from simultaneous mutations at both
911 ends. Then the mutation weights are generated for each mutation state at each of the 200 sites in
912 the array and are given by 1-log10(p), where p is the observed probability of the mutation on that
913 site. The weights define how important characters should be considered when comparing the
914 mutation states between cells. Then bifurcate clustering of nearest cells was carried out based on
915 matrix calibration. In the training set, the characters of all cells at all sites will be presented as *n*
916 200 by 100 matrices, where *n*=30 is the number of array characters (0, A, B, ... ). The inner
917 product of the matrices, which is *n* 100 by 100 matrices, reveals the relationship between the 100
918 cells in each tree of the training set themselves according to the 200 states, and the sum of *n*
919 product matrices gives the overall pairwise similarity relationship of the 100 cells, where we can
920 extract the most similar cell pair by the maximum value in that matrix (denoted as dark red, and
921 the indices of the cells are denoted as *i* and *j*). Then a parent cell, generated based on the shared
922 mutations of the two cells, replaces the two cells and is sent back to next iterations of bifurcate
923 clustering, until only one pair of cells is left and their parent cell will become the tree root. **E.**
924 *Top.* For the *in vitro* challenge Cassiopeia-*ILP* (*Yosef Lab*) takes as input a "character matrix,"
925 summarizing the mutations seen at heritable target sites across cells and infers a Steiner Tree,
926 finding the tree of minimum weight connecting all observed cell states across all possible
927 evolutionary histories using integer linear programming (ILP). Importantly, the edges connecting
928 cell states can be weighted by the number of mutations along that edge or the log-likelihood of
929 these mutations.  *Bottom.* For both *in silico* challenges *Cassiopeia-Greedy* infers a phylogeny
930 from the observed character-states across all cells, which can be summarized in a cell's x cut-site
931 "character-matrix". To do so, the algorithm recursively applies a heuristic to split cells into two
932 groups based on the frequency of a given state at a character, *n(i, s)*, and the likelihood of that
933 state arising, *p(s)*. This procedure is applied until a full phylogeny is resolved. **F.** Using the 76
934 trees in the training set of the *in vitro* challenge to compare the relationships between cells that
935 share a particular state, *Liu lab* quantified how rarer states are more predictive of the true

936    relationship between pairs of cells. As observed in the plot, these relative rates can vary by both
937    identity and for each of the ten positions in the target array. **G.** *Cassiopeia-ILP* (*Yosef Lab*) is
938    able to incorporate learned state priors by weighting evolutionary transitions by their log-
939    likelihoods and find a Weighted Parsimony solution. Performance on the training data can
940    inform whether Weighted or Unweighted Parsimony is better suited.

**Figure 4. *DCLEAR* Learning k-mer replacement distances by simulation.** **A.** The input sequences were first used to estimate the summary statistics such as mutation rate ($\mu$), outcome probability of each character, number of targets and number of tips. These estimated parameters, combined with the pre-defined parameters such as cell divisions, were used to simulate multiple lineage trees from the root node. The k-mer nodal distances were estimated from these simulated lineage trees and then used to compute the distances between input sequences. **B.** The schematic shows a simulated lineage tree with one root, two internal nodes and three tips. The nodal distance is defined as the distance between any two nodes on the lineage tree. The expected nodal distance can be estimated from the replacement of individual characters (e.g. between A and C), the replacement of k-mers (e.g. between 0A and 0C), or sequences (e.g. between A000A and E00C). **C.** The heatmap shows the expected nodal distance of the replacement of the most frequent individual characters. **D.** The heatmap shows the probability of replacement of the most frequent individual characters at a nodal distance of 15. **E.** The histogram shows the posterior distribution of nodal distance of two sequences when having the same characters A or C at any specific position. **F-G.** The histograms show the observed distribution (red bars) and estimated posterior distribution of nodal distance of two sequences **F** with the replacement of C- by CC, or **G** with BBBB at the same position. The posterior distributions were estimated by using an independent model (blue bars) and a conditional model (green bars). In both cases, the posterior distribution estimated by the conditional model is more consistent with the observed distribution. **H.** The simulated trees were used to compare the performance of lineage reconstruction by using Hamming distance and k-mer replacement distances with different k's. We simulated 1,000 lineage trees with cell division of 16, mutation probability of 0.1, 200 targets and 200 tips. The outcome probability was sampled from a Gamma distribution with shape of 0.1 and rate of 2. For both k-mer replacement distances and Hamming distance, we used a balanced minimum evolution (ME) algorithm with tree rearrangement (nearest neighbor interchange, subtree pruning and regrafting, and tree bisection and reconnection) to infer the tree topology. The similarity between the inferred tree and the simulated tree was measured by the Robinson-Foulds (RF) distance.

971    **Figure 5. *AMbeRland* A decision tree based approach for reconstruct cell lineages**
972    **A.** After selecting manually different model features for *left* the *in vitro* challenge (F1 to F4) and
973    *right* the *C. elegans* challenge, *AMbeRland* learns the features importance represented by
974    histograms of the weights, for predicting phylogenetic relationships directly from the training
975    data using a Gradient Boosting Machine (GBM) *middle*. These learned weights are then used to
976    predict the probability of sister-cell relationships on the hold out test data creating a probability
977    matrix used for hierarchical reconstruction *bottom*. **B.** *Left* Trees are reconstructed from
978    probability matrices by performing a grid search to obtain the clustering thresholds at each tree
979    level while maximizing the RF and triplets metrics. *Right* Example of differences when
980    establishing thresholds for Tree 29, the largest correctly reconstructed tree in the *in vitro*
981    challenge. See also detailed examples in **Fig S7 & S8**.

982

983

984

985 **Figure 6. Consensus methods and agreement in tree reconstruction. A.** Depth-dependent
986 agreement between reconstructed trees calculated by Felsenstein Bootstrap Proportion and
987 Transfer Bootstrap Expectation. Both metrics assess the degree of agreement that different trees
988 have on specific splits (or cell divisions). High agreement indicates that most teams resolved
989 splits correctly at that depth. The distribution is computed across all 30 trees in the *in vitro* test
990 sets. **B.** We computed the consensus trees by majority rule using the *consensus* function from the
991 *R package* **ape** v5.3. The consensus performance in the *in vitro* challenge is higher than any
992 individual team by RF distance but not by triplets (red dotted line indicates the best performed by
993 each metric). **C.** Scores summarizing all participating methods for the *in silico* challenges,
994 including the PHYLIP consensus and for reference *FastTree2*. **D.** Annotated subtree of *C.*
995 *elegans* challenge, edges are marked with tables listing the agreement of each of the 5 individual
996 submissions and the consensus in Transfer Bootstrap Distance where 1 is high agreement. Colors
997 refer to the table in **C**.
998
999

**Supplementary figures and tables**


**Figure S1.** *Mus musculus in silico* **challenge A.** Simulation of the Mouse lineage, "token" cells whose lineage are stochastically chosen to be followed as the lineage tree is formed, are shown in blue, in white are represented cells whose lineage is not followed. At the end of the simulation for the mouse lineage information for about 10,000 blue cells is stored, but it is estimated that the size of the tree is about $10^{12}$ or a trillion cells. **B.** Visualization of the 10,000 cell Mouse tree with 11 types of cells encoded by different colors.


**Figure S2.** *In vitro* **challenge rankings for all teams according to multiple metrics.**
 The ranks for each team were evaluated by calculating the ranksum values (left boxplots) for the Robison-Foulds (middle boxplots) and the triplet metric (right boxplots) sampled 1000 times with replacement from the scores for the 30 individual trees. The 9 teams were ordered by average ranksum and the Bayes Factor (BF) was calculated, yellow boxes show teams that are considered to be tied as they have a $1/3 < BF < 3$ and a $BF > 3$ against all the other teams in grey. Implementation of a third metric calculating quartets could not differentiate the top 3 teams: Yosef Lab (*Cassiopeia*) 0.4200, Guan Lab  0.4232, Jasper06  0.4243.


**Figure S3.** *In vitro* **challenge results with Robinson-Foulds and triplets metrics.**
The participant teams' distribution of scores across 30 reconstructed lineage trees is shown for **A.** triplets metric **B.** Robinson-Foulds metric **C.** Histogram showing the difference between the Robinson-Foulds and triplets metrics for all 30 trees across all teams. Median of zero indicates that overall the metrics agree but dispersion suggests a small bias for higher distance values in triplets. **D.** The histogram of scores of all 30 trees for all 9 teams are for *left* Robinson-Foulds and *right* triplets metrics, color coded depending on the percentage of unique barcode arrays in the tree. Deep blue dots trees with 25-50% unique arrays, gray blue dots trees with 50-75% unique arrays, light blue dots, trees with 75-100% unique arrays. **E.** Comparison of team performance depending on whether cells with degenerate barcodes are merged (gold boxes) or not (blue boxes). *Left* Boxplots represent the triplet distances, *Right* RF distances, of trees where for both predictions and ground truth, cells with the same barcodes were merged into a single

1030     leaf. The procedure followed for each tree a 100x bootstrap choosing each time a different cell

1031     with the same barcode as distances were recalculated for each fold.

1032

1033     **Figure S4. *In vitro* challenge list of trees that were reconstructed perfectly by at least one**

1034     **team.** Ground truth lineages are shown along with the array state for each cell.

1035

1036     **Figure S5. *In vitro* challenge largest trees with high reconstruction scores.**

1037     Two examples of large trees with 29 and 23 cells respectively and their RF and Triplets distance.

1038     These large but accurate trees were reconstructed by **A**) *AMberLand* and **B**) *Yosef Lab*

1039     (*Cassiopeia*).

1040

1041     **Figure S6. *Guan Lab* approach for *in vitro* challenge A.** Probability of mutations for the array

1042     sites and their corresponding weights for the Hamming distance. When calculating the weights

1043     for the Hamming distance, the mutation direction preference is set as reciprocal of the mutation

1044     frequency so that the rarer the mutation type, the more weight it is given to the distance

1045     between cells. **B.** A rule-based Hierarchical clustering approach was used to generate the trees.

1046     The cells character arrays final states were transformed by weights according to the observed

1047     probability of mutations, and the transformed states were used to calculate the distance

1048     between cells. The hierarchical clustering was done using a rule-based method to reconstruct

1049     parent cells, based on the fact that the editions from initial states (1) to edited states (0 and 2)

1050     are irreversible. **C.** Comparison of different clustering methods for the distance matrices

1051     including Rule-based hierarchical clustering, UPGMA and Neighbor Joining. The performance is

1052     shown for both triplets and RF distances. The Distribution across the 30 lineages in the test set

1053     and the average of the two tree measurements is shown by the violin plots. The rule-based

1054     hierarchical clustering method and UPGMA have similar performance on reconstructing cell

1055     lineage trees.

1056

1057     **Figure S7. Representation of the decision tree and weights** obtained by *Amberland* using

1058     GBM for the training set in the *in vitro* challenge. For each decision tree leaf are indicated: on

1059     top the feature's weight, the number of cells *n* and the percentage of the training set cells they

1060     represent, and in bold is the criteria of the feature used for selecting the next leaf *i.e* number of

1061     times the feature is present when comparing the 2 cells character arrays. Features in this case are:

1062 F1-both not mutated, F2-both same mutation F3-one mutation F4-different mutations. This

1063 figure was made using the *R package* "rattle".

1064

1065 **Figure S8. Reconstructing trees by clustering probability matrices as implemented by**

1066 **AMbeRland for the training set of the *in vitro* challenge.** Seventy six trees of different number

1067 of cells were used to optimize the tree reconstruction thresholds from the probability matrix of

1068 cells being sisters obtained from training a GBM algorithm **A.** Performance of the algorithm for

1069 four sets of thresholds: **set_A**=(0,0,0,0,0) results in mean RF=0.512 and triplets=0.389;

1070 **set_B**=(0.5,0,0,0,0) results in mean RF=0.519 and triplets=0.380; **set_C**=(0.8,0.4,0.2,0.1,0.05)

1071 results in mean RF=0.512 and triplets=0.433; and **set_D**=(0.3,0.1,0.05,0.01,0.005) results in

1072 mean RF=0.502 and triplets=0.375. The numbers shown in the scatter plots represent the tree ID

1073 and the color represents the number of cells in the tree. Threshold **set_D** was used to reconstruct

1074 the test dataset for submission. **B.** A perfectly reconstructed tree with 3 thresholds (tree ID 70

1075 from the training set, RF=0 and triplets=0) has 7 pairs joined into clusters at level 1, 4 pairs

1076 joined at level 2 and 2 pairs joined at level 3. **C.** Probability matrices for tree 70 are plotted for

1077 each level. From here it can be seen that cells 7 and 8 have the highest probability so they are

1078 first joined into cluster C1, the next pair with highest probability comprises cells 12 and 13

1079 which joined into cluster C2 and so on. Once all pairs are defined, the algorithm moves to Level

1080 2, where clusters C2 and C3 have the highest pairwise probability (cells on these two clusters can

1081 be seen on top right corner of level 1 probability matrix) so they are joined into a new cluster C1.

1082 The algorithm proceeds until all cells are joined into a single lineage.

1083

1084 **Figure S9. Clustering of cells into trees performed by *AMbeRland* for the training set in the**

1085 ***C. elegans in silico* challenge.** One hundred trees of a hundred cells each were used to optimize

1086 the tree reconstruction thresholds from the probability matrix of cells being sisters obtained from

1087 training a GBM algorithm **A.** Comparing performance of the algorithm for two sets of

1088 thresholds: **set_A**={0} gives mean RF distance=0.78 and triplets=0.59; **set_B**=(0.07, 0.04, 0.01,

1089 0.05, 0, 0, 0,0) gives mean RF distance=0.71 and triplets=0.49. Threshold **set_B** was used to

1090 reconstruct the test sample. **B.** Ground truth and reconstructed tree for training sample 100, with

1091 RF distance = 0.48 and triplets=0.44. **C.** Probability matrices for training sample 100 are plotted

1092 for each level. Clusters identified letters C. by Four clusters for level 7 (C1-C4) are indicated on

1093 the reconstructed tree in **B**.

1094

1095 **Figure S10. Agreement distribution across all reconstructed trees at different normalized**

1096 **tree depths for the *in vitro* challenge.** A depth of 0 represents the root of the tree whereas a

1097 depth of 1 corresponds to the leaves and therefore the depth of cell divisions within the lineage

1098 fall between [0,1]. *Top* For a given ground truth lineage, The Felsenstein Bootstrap Support is

1099 calculated across all reconstructed trees submitted by the teams corresponding to that lineage.

1100 We obtain a distribution by computing the FBS score for all 30 ground truth lineages. *Bottom*

1101 The Transfer Bootstrap Expectation is calculated in an analogous way.

1102

1103 **Table S1.**  Training and test datasets for the *in vitro* challenge.

1104

1105 **Table S2.** Comparing machine learning approaches for reconstruction of the *in vitro* cell lineage

1106 trees.

1107

1108 **Table S3.** Comparing machine learning approaches for reconstruction of the *in silico* large cell

1109 lineage trees (for comparison all methods were implemented in a two Intel(R) Xeon(R) CPUs @

1110 2.20GHz).

1111

1112

Box 1

**A** *in vitro* 106 lineages

Ground truth trees from videomicroscopy

Character array

**Train set:** 76 lineages w ground truth + character arrays
**Test set:** 30 lineages with only character arrays

Divide acording to: - Similar lineage size
- Similar reconstruction accuracy

106 lineages

Number of trees — Number of cells

RF distance

Partition
Test
Train

Test    Train
30%     70%     % of data

**B** **Distribution of simulated mutations**

Relative frequency — Mutation types

GESTALT v7
(replicate/target mean)
Gamma dist ($\kappa= 0.1$, $\theta= 2$)

**C** **Deletions** N targets (N=1000)

(N=200)

Two simultaneous CRISPR edits

10 targets are lost

N-10 targets

**D** **Dropouts**

Density of cell coverage factors

Information retained
Acquisition dropout

Density of site coverage factors

**E** *in silico*

**Train set:** 100 lineages with ground truth + character arrays
Trees of 100 cells/200 characters | Trees of 1000 cells/1000 characters
**Test set:** *C. elegans* lineage | simulated *M.musculus* lineage

Training set construction

Prune and regraft (100 iterations)

Simulate mutations in CRISPR recorder

Subsample 100 cells & change cells ID

**F**

Training set performance

Amberland
Cassiopeia
DCLEAR(WHD)
DCLEAR(KRD)
Liu
Guan

★ Test set    0.4  0.6  0.8  1.0    0.25 0.50 0.75 1.0
Robinson-Foulds distance    Triplet distance

Box 2



Ground truth    Reconstructed    RF = 0.66

Triplets = 0.4

Figure 1



**A**

Experiment: intMEMOIR in stem cells

1) MEMOIR system

mESC

DNA recording array    Int
                       Integrase

2) Recording events

Int → 2
1 → 0

3) FISH readout

Data: Recorded mutations          Inference

1
2
3
4
5
6
9
10
11
13
14
15

Time-lapse movie groud-truth

intMEMOIR stem cell line

Time (hr)
0    36   54

Tree comparison and scoring

**B**

Simulate
CRISPR-Cas9 Recorder

N Targets
(200 / 1,000)

Zygote

Mutational outcomes

Unmutated -> "0"
Outcome #1 -> "A"
Outcome #2 -> "B"
Outcome #3 -> "C"
⋮
Outcome #30-> "d"

c_0001    c_0002    c_0003    c_0004    **Cell**

0B00A000C0  0B00A00000  0000A0C000  0d00A0C000  **Barcode**

**C**

P1        P0        P0A

~1mm

**D**

Zygote
~E1.0

Blastocyst
~E4.0
ICM

~E6.5

$Zygote \xrightarrow{\tau} \{Zygote, Zygote\}_{p1} | \{Morula, Morula\}_{p2}$

$Morula \xrightarrow{\tau} \{Morula, Morula\}_{p1} | \{Morula, ICM\}_{p2} | \{Morula, Trophectoderm\}_{p3}$

$ICM \xrightarrow{\tau} \{ICM, ICM\}_{p1} | \{ICM, Epiblast\}_{p2} | \{ICM, PE\}_{p3}$

$Trophectoderm \xrightarrow{\tau} \{Trophectoderm, Trophectoderm\}_{p1} | \{Trophectoderm, Placenta\}_{p2}$

$Epiblast \xrightarrow{\tau} \{Epiblast, Epiblast\}_{p1} | \{Epiblast, Mesoderm\}_{p2} | \{Epiblast, Ectoderm\}_{p3}$

$PE \xrightarrow{\tau} \{PE, PE\}_{p1} | \{PE, Endoderm\}_{p2} | \{PE, YolkSac\}_{p3}$

Endoderm
(internal layer)

Lung cells (alveolar cell)    Thyroid cells    Digestive cells (pancreatic cell)

Mesoderm
(middle layer)

Cardiac muscle cells    Skeletal muscle cells    Tubule cells of the kidney    Red blood cells    Smooth muscle cells (in gut)

Ectoderm
(external layer)

Skin cells of epidermis    Neuron on brain    Pigment cells

Simulate 1y of mouse development

100 training trees, 1K cells
1 leaderboard tree 6K cells
1 final tree 10K cells

**E**

1 month

6 months

1 year

Figure 2



**A** *in vitro challenge*

**B** *in silico C. elegans challenge*

**C** *in silico M. musculus challenge*

**D** Number of perfect reconstructions *in vitro* challenge

**E** Large vs Small trees *in vitro challenge*

**F** All teams predictions for all trees *in vitro* challenge

Figure 3

Figure 4

A. [Flowchart]
Input sequences → • Mutation rate μ / • Outcome probability p(X) / • Number of targets / • Number of tips / • Cell division (16) → Simulate 1,000 lineage trees → Summarize k-mer nodal distance: **E[d|X,Y]**

Compute sequence distance
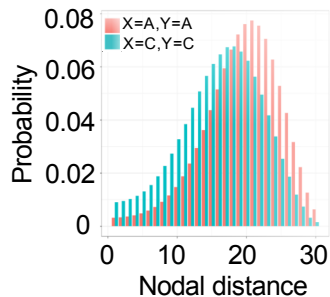
B. [Tree diagram]
root
0000
A00A    000C
A0BA    E00C    0---
tip 1   tip 2   tip 3

Nodal distance
Character distance (e.g. E[d|A,C])
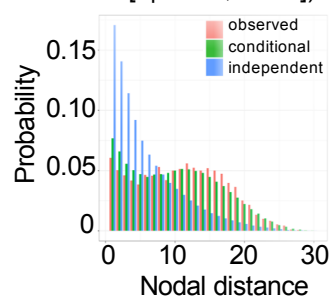k-mer distance (e.g. E[d|0A,0C])
Sequence distance (e.g. E[d|A00A,E00C])

C. Expected nodal distance (**E[d|X,Y]**)
distance
16
17
18
19
20
21

D. Character replacement (p(X,Y|d))
Prob
0.25
0.20
0.15
0.10
0.05

E. Distance distribution
X=A,Y=A
X=C,Y=C
Probability
Nodal distance

F. E[d|C-,CC])
observed
conditional
independent
Probability
Nodal distance

G. E[d|BBBB,BBBB])
observed
conditional
independent
Probability
Nodal distance

H.
k=4
k=3
k=2
k=1
Hamming
RF distance

# Figure 5

## A

### AMbeRland



## B

Figure 6

**A**



Normalized tree depth

0                                                    1

depth across all trees.

**B**



**C**

| Team name | | *C elegans* | | | leaderboard | | | final | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | #cells | RF | Triplets | #cells | RF | Triplets | #cells | RF | Triplets |
| DCLEAR | | 1000 | 0.5567 | 0.5062 | 6142 | 0.6631 | 0.2591 | 9745 | 0.5575 | 0.2588 |
| Liu | | 1000 | 0.6209 | 0.9647 | 6142 | 0.7276 | 0.4039 | 9745 | 0.7136 | 0.4993 |
| Cassiopeia | | 1000 | 0.9238 | 0.4745 | 6142 | 0.7487 | 0.5967 | 9745 | 0.8768 | 0.6197 |
| Guan | | 1000 | 0.676 | 0.9615 | 6142 | 0.5729 | 1 | | | |
| AmBerLand | | 1000 | 0.8215 | 0.9836 | | | | | | |
| SanGuo | | 1000 | 0.998 | 1 | | | | | | |
| Consensus | | 1000 | 0.4804 | 1 | 6142 | 0.6219 | 0.6854 | 9745 | 0.5909 | 1 |
| FastTree2 | | 1000 | 0.7202 | 1 | 6142 | 0.4058 | 0.2138 | 9745 | 0.402 | 0.1495 |

**D**