

Supplemental Material – Modeling and predicting cancer clonal evolution with reinforcement learning

Stefan Ivanovic¹ and Mohammed El-Kebir^{1,2,†}

¹Department of Computer Science, University of Illinois at Urbana-Champaign, IL 61801, USA

²Cancer Center at Illinois, University of Illinois Urbana-Champaign, Urbana, IL 61801, USA

[†]Correspondence: melkebir@illinois.edu

Contents

A	Supplemental methods	3
A.1	CloMu’s model	3
A.1.1	Neural network	3
A.1.2	Basic introduction to policy gradients	4
A.1.3	Policy gradients for CloMu	6
A.1.4	Policy gradients with few trees per patient for CloMu	7
A.2	Additional CloMu prediction tasks	8
A.2.1	Determining evolutionary pathways	8
A.2.2	Determining multi-mutation causality	10
A.2.3	Identifying consensus trees	10
A.3	Bootstrapping	11
B	Supplemental results for simulations	12
B.1	Setup	12
B.2	Baseline methods and parameters	15
B.3	Results	16
B.3.1	Bootstrapping – Simulations (I-a)	17
B.3.2	Decreasing numbers of patients – adaptation of Simulations (I-a)	20
B.3.3	Decreasing numbers of mutations – adaptation of Simulations (I-a)	21
B.3.4	Causation strength based on real data – adaptation of Simulations (I-a)	22
B.3.5	Large numbers of passenger mutations – adaptation of Simulations (I-a)	23
B.3.6	Violations of independent clonal evolution – adaptation of Simulations (I-a)	24
B.3.7	Including highly incorrect possible trees – adaptation of Simulations (I-a)	25
B.3.8	Inhibition between interchangeable mutations – adaptation of Simulations (I-b,c)	28
B.3.9	CloMu with increased number L of hidden neurons – Simulations (I-c)	29
B.3.10	Mutations with shared properties	30
B.3.11	Isolating multi-mutation effects	34
B.3.12	Results on RECAP simulations – Simulations (III)	35
B.3.13	Results on TreeMHN simulations – Simulations (IV)	37

36	C Supplemental results for real data	38
37	C.1 Breast cancer	38
38	C.1.1 Bootstrapping	38
39	C.1.2 Breast cancer hormone receptor status and latent representations	42
40	C.1.3 Systematic comparison with TreeMHN	45
41	C.1.4 Predicting subsequent mutations on subtrees	46
42	C.2 AML	48
43	C.2.1 Additional analysis of high fitness mutations and prevalence	48
44	C.2.2 Additional analysis of latent representations	53
45	C.2.3 Systematic comparison with TreeMHN	55
46	C.2.4 Predicting subsequent mutations on subtrees	56

47 List of Figures

48	S1 Demonstration of tree generation	5
49	S2 A diagram of evolutionary pathways	10
50	S3 Number of trees in simulated data	13
51	S4 Bootstrapping demonstrates the robustness of CloMu's predictions on Simulations (I-a)	18
52	S5 Bootstrapping causal strength rankings show CloMu's robustness on Simulations (I-a)	19
53	S6 CloMu gives accurate predictions even on simulations with few patients	20
54	S7 The runtimes and memory usage of CloMu	21
55	S8 CloMu's performance remains high for simulations with fewer mutations per patient	22
56	S9 CloMu's performance is slightly reduced when using varying effect sizes in causality	23
57	S10 CloMu's performance is reduced when including a large number of passenger mutations	24
58	S11 CloMu is robust to violations of the independent clonal evolution assumption	25
59	S12 CloMu is robust to random trees in the set of possible trees	27
60	S13 Relative probabilities of trees in simulations with random trees included	28
61	S14 Adding inhibition between interchangeable mutations on Simulations (I-b) and (I-c)	29
62	S15 CloMu is robust to increasing the number L of hidden neurons	30
63	S16 A demonstration of mutations with shared properties	32
64	S17 CloMu accurately determines similar mutations and reconstructs mutation properties	33
65	S18 CloMu accurately predicts causal relationships on simulations with shared mutation properties	33
66	S19 CloMu detects non-linear multi-mutation causality caused by pairs of mutations	35
67	S20 CloMu's matches RECAP's performance on RECAP's simulated data	36
68	S21 CloMu performs slightly worse than TreeMHN on TreeMHN's simulated data	38
69	S22 Bootstrapping provides statistical bounds for breast cancer relative causality predictions	39
70	S23 Statistical bounds for ranking the strength of relative causality predictions	40
71	S24 Bootstrapping gives statistical bounds for fitness predictions on the breast cancer data set	41
72	S25 CloMu produces latent representations somewhat associated with hormone receptor status	43
73	S26 PCA of latent representations on the breast cancer dataset	44
74	S27 Many mutations have highly similar latent representations in the breast cancer data	44
75	S28 A diagram of next mutation prediction on subtrees	47
76	S29 CloMu accurately predicts future mutations given a partial tree on breast cancer data	48
77	S30 High fitness mutations tend to occur earlier in evolution	51
78	S31 Analyzing the prevalence of initial clones and terminal clones	52
79	S32 PCA of latent representations on the AML dataset	54
80	S33 Many mutations have highly similar latent representations in the AML dataset	54

81	S34 CloMu accurately predicts future mutations given a partial tree	57
----	---	----

82 **List of Tables**

83	S1 Demonstration of relative and absolute causality calculations	12
84	S2 Calculation of false positives/negatives using signed causality	16
85	S3 Comparison of TreeMHN and CloMu on breast cancer cohort	46
86	S4 Comparison of TreeMHN and CloMu on AML cohort	56

87 **A Supplemental methods**

88 Supplementary Material [A.1](#) provides details on CloMu’s model and its training. We describe additional
 89 prediction tasks in Supplementary Material [A.2](#) Finally, in Supplementary Material [A.3](#), we discuss boot-
 90 strapping.

91 **A.1 CloMu’s model**

92 We begin by providing details on CloMu’s neural network in Supplementary Material [A.1.1](#). We then
 93 discuss how we used reinforcement learning to identify model parameters f_θ for the neural network. Specif-
 94 ically, Supplementary Material [A.1.2](#) gives a brief introduction to policy gradients. In Supplementary Ma-
 95 terial [A.1.3](#) we discuss how policy gradients can be applied to CloMu’s model. Finally, in Supplementary
 96 Material [A.1.4](#) we discuss a speedup for the typical case of a small number of trees per patient.

97 **A.1.1 Neural network**

98 CloMu uses a two layer neural network with a small number L of hidden neurons for the function f_θ .
 99 Specifically, the input to the neural network is the vector \mathbf{c} and the output is a log rate $f_\theta(\mathbf{c}, s)$ for each
 100 mutation s . This model has $2Lm + m + L = O(Lm)$ total parameters. We use a tanh nonlinearity rather
 101 than a ReLU nonlinearity to reduce the risk of latent variables being stuck at zero during training. To
 102 optionally enforce the infinite sites assumption, we set $f_\theta(\mathbf{c}, s)$ to a very large negative number if $c_s = 1$ to
 103 ensure a mutation is not added twice to the same clone.

104 To enable easier regularization, which is used primarily for the sake of interpretability but also used to
 105 avoid overfitting on some extremely small datasets, we made an two additional modifications to the neural
 106 network. First, the neural network also contains a single scalar value for each input mutation, which is added

107 to all variables in the output equally. This allows one variable per mutation to influence the general fitness
108 of the clone, allowing for a mutation to affect clonal fitness while only modifying one parameter (rather
109 than modifying one parameter per output mutation). Specifically, this is equivalent to adding a function
110 independent of s , $h_\theta(\mathbf{c})$, to $f_\theta(\mathbf{c}, s)$ like the below equation.

$$111 \quad f_\theta(\mathbf{c}, s) := f_\theta(\mathbf{c}, s) + h_\theta(\mathbf{c}). \quad (1)$$

112 Second, a one-hot encoded variable representing the number of mutations in the clone is included. This
113 allows the model to take into account the number of mutations on a clone without modifying the parameter
114 for every individual input mutation. This therefore allows the model to easily know if the clone is the
115 non-mutated root clone with only one parameter.

116 These two modifications are essentially useless and give no empirical advantage when no regularization
117 is used, since these variables can automatically be computed and used from the already existing data. How-
118 ever, once L1 regularization is included during training, these variables help the model achieve sparsity. On
119 very small datasets sparsity can be important to avoid overfitting; however, on larger datasets, sparsity is
120 simply useful because it helps with interpretability. Having all variables that do not need to be nonzero set
121 to zero helps with creating good plots and thus interpretability. This also likely helps with any downstream
122 tasks that would use the latent representation of mutations or clones generated by the model.

123 **A.1.2 Basic introduction to policy gradients**

124 The type of reinforcement learning used in this paper is policy gradients. Define A as a sequence of actions
125 the model may take, $R(A)$ as the reward for those actions, and $\Pr(A | f_\theta)$ as the probability our model takes
126 those actions. As an example, if the environment was a video game, A could be a sequence of actions taken
127 in the game, and $R(A)$ could be the score achieved by taking those actions. In our application, A will be the
128 phylogenetic generation processes sampled from our model, which is demonstrated in Fig. S1. Define S_A as
129 the set of all possible sequences of actions our model could take. In our application, S_A would correspond
130 to the set of all tree generating processes. To optimize the reward that the model gets, one wants to increase
131 the probability of high reward actions. Therefore, one wants to identify model parameters θ that maximize

Example:: There exists two mutations, **Red** and **Blue**. Both mutations naturally occur on clones at a rate of 1, but if a clone has mutation **Red**, then mutation **Blue** occurs with rate 2.

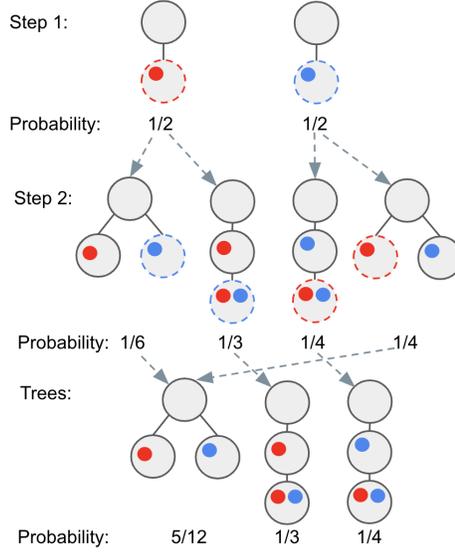


Figure S1: **Demonstration of tree generation.** This is a toy example to demonstrate how causal relationships affect the probability of phylogeny trees in a simple system. It also demonstrates how there are multiple ways of generating the same tree; specifically, trees 1 and 4 in step 2 are the same tree.

132 the below equation.

$$133 \quad \sum_{A \in S_A} R(A) \Pr(A | f_\theta). \quad (2)$$

134 The ideal gradient for optimizing this would theoretically be the below equation.

$$135 \quad \frac{d}{d\theta} \sum_{A \in S_A} R(A) \Pr(A | f_\theta) = \sum_{A \in S_A} R(A) \frac{d}{d\theta} \Pr(A | f_\theta). \quad (3)$$

136 In practice, it is infeasible to sum over all $A \in S_A$ since S_A may be a very large set. However, one can
137 estimate (3) by sampling from $\Pr(A | f_\theta)$:

$$138 \quad \sum_{A \in S_A} R(A) \frac{d}{d\theta} \Pr(A | f_\theta) = \sum_{A \in S_A} R(A) \Pr(A | f_\theta) \frac{d}{d\theta} \log(\Pr(A | f_\theta)) = \mathbb{E}_A[R(A) \frac{d}{d\theta} \log(\Pr(A | f_\theta))] \quad (4)$$

139 Therefore one uses the standard REINFORCE (Williams, 1992) gradient estimate

$$140 \quad \mathbb{E}_A[R(A) \frac{d}{d\theta} \log(\Pr(A | f_\theta))]. \quad (5)$$

141 Note, this section is not intended to give a complete description of policy gradients. Rather, it explains
 142 the necessary content for an understanding of CloMu.

143 A.1.3 Policy gradients for CloMu

144 In our application, a sequence A of actions will be the phylogenetic generation processes G sampled from
 145 our model. Moreover, the set S_G of all possible sequences of actions corresponds to the set of all tree
 146 generating processes on m mutations. To use reinforcement learning, we need to meet two criteria. First,
 147 we need to be able to sample from our distribution $\Pr(G | f_\theta)$. For our model, this is very easy since G
 148 is generated by a simple probabilistic process (Section Independent clonal evolution problem). Second, we
 149 need to be able to formulate our objective as the below equation, where $R(G)$ is the reward function.

$$150 \quad \sum_{G \in S_G} R(G) \Pr(G | f_\theta). \quad (6)$$

151 An important aspect of $R(G)$ is that the gradient with respect to θ that is used for optimization does not
 152 affect $R(G)$. To make this explicit and clear, instead of formulating our objective function as that objective
 153 function, we formulate the gradient of our objective function to fit the reinforcement learning objective
 154 function gradient. Specifically, to apply policy gradients, the gradient of our objective function must be
 155 formulated as below.

$$156 \quad \sum_{G \in S_G} R(G) \frac{d}{d\theta} \Pr(G | f_\theta). \quad (7)$$

157 We can do so through the below manipulations of the gradient of the log likelihood.

$$158 \quad \frac{d}{d\theta} \log \Pr(\mathcal{T}_1, \dots, \mathcal{T}_n | f_\theta) = \frac{d}{d\theta} \sum_{p=1}^n \log \left(\sum_{T \in \mathcal{T}_p} \sum_{G \in \mathcal{G}(T)} \Pr(G | f_\theta) \right) \quad (8)$$

$$159 \quad = \sum_{p=1}^n \left(\sum_{T \in \mathcal{T}_p} \sum_{G \in \mathcal{G}(T)} \Pr(G | f_\theta) \right)^{-1} \frac{d}{d\theta} \sum_{T \in \mathcal{T}_p} \sum_{G \in \mathcal{G}(T)} \Pr(G | f_\theta) \quad (9)$$

$$160 \quad = \sum_{p=1}^n \left(\sum_{T \in \mathcal{T}_p} \sum_{G \in \mathcal{G}(T)} \Pr(G | f_\theta) \right)^{-1} \sum_{T \in \mathcal{T}_p} \sum_{G \in \mathcal{G}(T)} \frac{d}{d\theta} \Pr(G | f_\theta) \quad (10)$$

$$161 \quad = \sum_{p=1}^n \sum_{T \in \mathcal{T}_p} \sum_{G \in \mathcal{G}(T)} \left(\sum_{T' \in \mathcal{T}_p} \sum_{G' \in \mathcal{G}(T')} \Pr(G' | f_\theta) \right)^{-1} \frac{d}{d\theta} \Pr(G | f_\theta) \quad (11)$$

$$= \sum_{G \in \mathcal{S}_G} \sum_{p=1}^n \frac{\delta(p, G)}{\sum_{T \in \mathcal{T}_p} \sum_{G' \in \mathcal{G}(T)} \Pr(G' | f_\theta)} \cdot \frac{d}{d\theta} \Pr(G | f_\theta) \quad (12)$$

In the above derivation, we define $\delta(p, G)$ as 1 if $G \in \mathcal{G}(T)$ for some $T \in \mathcal{T}_p$, and $\delta(p, G) = 0$ otherwise. In other words $\delta(p, G)$ indicates whether G corresponds to an observed tree of patient p . Equation (12) corresponds to the reinforcement learning gradient if we set the reward function as below.

$$R(G) = \sum_{p=1}^n \frac{\delta(p, G)}{\sum_{T \in \mathcal{T}_p} \sum_{G' \in \mathcal{G}(T)} \Pr(G' | f_\theta)}. \quad (13)$$

In reinforcement learning, one is sampling from $\Pr(G | f_\theta)$, so we can easily estimate $\sum_{T \in \mathcal{T}_p} \sum_{G \in \mathcal{G}(T)} \Pr(G | f_\theta)$ by simply observing what fraction of generated G are in $\mathcal{G}(T)$ with $T \in \mathcal{T}_p$. In other words, we can observe what fraction of generated trees are possible for patient p . This fraction changes as we train the model, and therefore must constantly be updated. However, assuming the model does not change much in a small number of iterations, it does not need to be calculated from scratch in every gradient update. The reinforcement learning gradient is now the below expression

$$\mathbb{E}_G [R(G) \frac{d}{d\theta} \log(\Pr(G | f_\theta))]. \quad (14)$$

A.1.4 Policy gradients with few trees per patient for CloMu

When the set of possible trees per patient is small enough to be reasonably explicitly enumerated, it can be beneficial to directly model the probability of each individual tree, rather than indirectly modeling the probability of each tree through reinforcement learning. In particular, when there are very few possible trees per patient, the reinforcement learning method introduced in the previous section may struggle to find those trees and would perform much better if the probability of each tree was explicitly modeled. However, one cannot just avoid using reinforcement learning in those cases because there could still be a very large number of generation processes per tree. To solve both of these problems, we use a modified reinforcement learning system where the probability of each possible tree is explicitly modeled, but the probability of each generation process is not explicitly modeled and only learned through reinforcement learning. This version of CloMu was implemented and is used for all experiments as well as made publicly available.

More specifically, for each possible tree T in the dataset (the union of \mathcal{T}_p for all patients $p \in [n]$), we only sample generation processes $G \in \mathcal{G}(T)$. Of course, this requires a modification of the objective

187 function to compensate for the change in the sampling procedure. Define $P_T(G | f_\theta)$ as the probability
 188 of generating G with our procedure restricted such that it always generates processes which give the tree
 189 T . Define $R'(G) = R(G) \Pr(G | f_\theta) / P_T(G | f_\theta)$. $R'(G)$ is the new reward function for this adjusted
 190 sampling procedure as shown in the below calculation.

$$191 \quad \mathbb{E}_{G \sim P_T}[R(G) \log(\Pr(G | f_\theta))] = \mathbb{E}_{G \sim P_T}[\Pr(G | f_\theta) / P_T(G | f_\theta) R(G) \log(\Pr(G | f_\theta))] \quad (15)$$

$$192 \quad = \mathbb{E}_{G \sim P_T}[R'(G) \log(\Pr(G | f_\theta))]. \quad (16)$$

193 A.2 Additional CloMu prediction tasks

194 This section provides details on determining evolutionary pathways (Supplementary Material A.2.1), multi-
 195 mutation causality (Supplementary Material A.2.2) and consensus trees (Supplementary Material A.2.3).

196 A.2.1 Determining evolutionary pathways

197 For the sake of interpretability, it is useful to represent patterns in cancer evolution as evolutionary pathways
 198 that frequently occur. We define an evolutionary pathway of interchangeable mutations as a list of sets of
 199 mutations such that the cancer obtained one mutation from each set in order. More precisely, we define
 200 an *evolutionary pathway* of length ℓ as sets of mutations S_1, \dots, S_ℓ . A clone with mutations s_1, \dots, s_ℓ in
 201 temporal order fits the evolutionary pathway S_1, \dots, S_ℓ if $s_k \in S_k$ for all $k \in [\ell]$.

202 The probability of an evolutionary pathway S_1, \dots, S_ℓ is defined as the probability that a clone with ℓ
 203 mutations fits the evolutionary pathway. Define $\mathbf{c}(s_1, \dots, s_k)$ as the clone with mutations s_1, \dots, s_k . The
 204 probability of the evolutionary pathway S_1, \dots, S_ℓ is given by the below expression.

$$205 \quad \Pr(S_1, \dots, S_\ell | f_\theta) = \sum_{(s_1, \dots, s_\ell) \in S_1 \times \dots \times S_\ell} \Pr((\mathbf{c}_0, s_1), (\mathbf{c}(s_1), s_2), \dots, (\mathbf{c}(s_1, \dots, s_{\ell-1}), s_\ell) | f_\theta). \quad (17)$$

206 Fig. S2 gives an example pathway and set of trees for this calculation. As a null model, we consider the
 207 probability $\Pr_0(s | f_\theta, \ell)$ of a mutation s occurring in a pathway of length ℓ .

$$208 \quad \Pr_0(s | f_\theta, \ell) = \sum_{\substack{\{s_1, \dots, s_\ell\} \subseteq [m] \\ \text{with } s \in \{s_1, \dots, s_\ell\}}} \Pr(\{s_1\}, \dots, \{s_\ell\} | f_\theta), \quad (18)$$

209 where $\Pr(\{s_1\}, \dots, \{s_\ell\} | f_\theta)$ is the probability of the pathway $\{s_1\}, \dots, \{s_\ell\}$.

210 The null probability, ignoring causal relationship, of a mutation s occurring at a particular step for a
 211 particular clone would be $\Pr_0(s | f_\theta, \ell)/\ell$. The null probability for a pathway would then be the null
 212 probability of generating every clone that fits the pathway based on the default probability in every step
 213 of generating the clone. Specifically, the null probability of a pathway S_1, \dots, S_ℓ is defined by the below
 214 equation.

$$215 \quad \Pr_0(S_1, \dots, S_\ell | f_\theta) = \sum_{(s_1, \dots, s_\ell) \in S_1 \times \dots \times S_\ell} \prod_{k=1}^{\ell} \Pr_0(s_k | f_\theta, \ell)/\ell = \prod_{k=1}^{\ell} \sum_{s_k \in S_k} P_0(s_k | f_\theta, \ell)/\ell. \quad (19)$$

216 The expression $\Pr(S_1, \dots, S_\ell | f_\theta) / \Pr_0(S_1, \dots, S_\ell | f_\theta)$ defines how much more likely a pathway is to
 217 occur than what one would expect if there were no causal relationships. For a pathway to be meaningful,
 218 one wants $\Pr(S_1, \dots, S_\ell | f_\theta) / \Pr_0(S_1, \dots, S_\ell | f_\theta)$ to be high in order to indicate that the pathway
 219 demonstrates meaningful causal relationships. Additionally, one wants $\Pr(S_1, \dots, S_\ell | f_\theta)$ itself to be high
 220 so that the pathway occurs frequently enough to be meaningful. Define the score of a pathway with the
 221 below equation, where $a > 0$ is a nonnegative constant.

$$222 \quad \text{score}(S_1, \dots, S_L, a) = \Pr(S_1, \dots, S_L | f_\theta)^a \frac{\Pr(S_1, \dots, S_L | f_\theta)}{\Pr_0(S_1, \dots, S_L | f_\theta)}. \quad (20)$$

223 The value of a can be adjusted based on the specific application. When implementing this we operate in
 224 logarithmic space.

225 In order to determine the most meaningful pathways, one can optimize for the score of the pathway.
 226 Specifically, given a length $\ell \leq m$ the procedure used in this paper starts with an evolutionary pathway
 227 S_1, \dots, S_ℓ where each $S_k = [m]$. Then, in a greedy fashion, we add or remove a single mutation that
 228 maximizes $\text{score}(S_1, \dots, S_\ell, a)$. We repeat this until a locally optimal pathway is found. Once an individual
 229 optimal pathway is found, we set the probability $\Pr((\mathbf{c}_0, s_1), \dots, (\mathbf{c}(s_1, \dots, s_{\ell-1}), s_\ell) | f_\theta)$ used in (17)
 230 to 0 in a lookup table so as to discourage repeatedly finding the same or similar pathways. Then, the
 231 optimization procedure is run again to generate a new pathway. If the score of a pathway reaches below
 232 $\Pr([m], \dots, [m] | f_\theta)^a$, it is not at all meaningful, and the process of generating optimal pathways ends.
 233 This cut off is used since the pathway consisting of only the set of all mutations at every step has that score.

$$\Pr \left(\begin{array}{l} S_1 = \{ \bullet, \bullet \} \\ S_2 = \{ \bullet \} \\ S_3 = \{ \bullet, \bullet \} \end{array} \right) = \Pr \left(\begin{array}{c} \text{Tree 1} \\ \text{Tree 2} \\ \text{Tree 3} \\ \text{Tree 4} \end{array} \right) + \Pr \left(\begin{array}{c} \text{Tree 5} \\ \text{Tree 6} \\ \text{Tree 7} \\ \text{Tree 8} \end{array} \right) + \Pr \left(\begin{array}{c} \text{Tree 9} \\ \text{Tree 10} \\ \text{Tree 11} \\ \text{Tree 12} \end{array} \right) + \Pr \left(\begin{array}{c} \text{Tree 13} \\ \text{Tree 14} \\ \text{Tree 15} \\ \text{Tree 16} \end{array} \right)$$

Figure S2: **A diagram of evolutionary pathways.** Eq. (17) states that the probability of an evolutionary pathway is equal to the sum of the probability of all trees that fit the evolutionary pathway (sequentially acquiring one mutation from each step in the pathway). In this example, S_1 has two mutations, S_2 has a single mutation and S_3 has two mutations, leading to $2 \cdot 1 \cdot 2 = 4$ possible trees. The probability of this evolutionary pathway equals the sum of the probabilities of the four resulting trees.

234 A.2.2 Determining multi-mutation causality

235 One advantage of CloMu is that it allows for complex interactions between mutations while utilizing a
 236 small number of parameters. The ability of CloMu to detect the effects of combinations of mutations was
 237 demonstrated in Simulations (II) in the Main Text. However, here we introduce *multi-mutation causality*
 238 to more explicitly describe this phenomenon. While in the Main Text we defined absolute and relative
 239 causality between a mutation s and another mutation t , here we extend the definition to capture the causal
 240 relationship between a combination of mutations and some other mutation. Specifically, we define absolute
 241 *multi-mutation causality* from a set of mutations in a clone \mathbf{c} to a mutation t as $f_\theta(\mathbf{c}, t) - f_\theta(\mathbf{c}_0, t)$. Therefore,
 242 for multi-mutation causality, we say the mutations in \mathbf{c} *cause* t if $f_\theta(\mathbf{c}, t) - f_\theta(\mathbf{c}_0, t) > \tau$, the mutations
 243 in \mathbf{c} *inhibit* t if $f_\theta(\mathbf{c}, t) - f_\theta(\mathbf{c}_0, t) < -\tau$, and there is *no causal relationship* from the mutations in \mathbf{c} to t
 244 otherwise. The simplest non-trivial case of multi-mutation causality is a pair of mutations causing some third
 245 mutation. Let \mathbf{c} be the clone with exactly mutations s and t . We say s and t *cause* r if $f_\theta(\mathbf{c}, r) - f_\theta(\mathbf{c}_0, r) >$
 246 τ , s and t *inhibit* r if $f_\theta(\mathbf{c}, r) - f_\theta(\mathbf{c}_0, r) < -\tau$, and there is no causal relationship from s and t to r
 247 otherwise. Our ability to detect multi-mutation causality is tested in Supplementary Material B.3.11.

248 A.2.3 Identifying consensus trees

249 One prediction task introduced in the RECAP (Christensen et al., 2020) and REVOLVER (Caravagna et al.,
 250 2018) papers is the identification of a small number of consensus trees. Accomplishing this task requires
 251 the knowledge that there are a small number of true trees in each dataset, which is not an intrinsic property

252 of our model. To adjust for this fact, a post processing step to our model must be added. Specifically, tree
 253 predictions are generated for each patient according to the following procedure. As previously defined, let
 254 $\Pr(T \mid f_\theta)$ be the probability that tree T is generated according to the model, and \mathcal{T}_p be the set of possible
 255 trees for patient p . For a set \mathcal{S} of trees define $\Pr(\mathcal{S} \mid f_\theta)$ as the below equation.

$$256 \quad \Pr(\mathcal{S} \mid f_\theta) = \prod_{p=1}^n \max_{T \in \mathcal{S} \cap \mathcal{T}_p} \Pr(T \mid f_\theta). \quad (21)$$

257 In other words, $\Pr(\mathcal{S} \mid f_\theta)$ is the probability of the maximum likelihood possible trees $(T_1, \dots, T_n) \in$
 258 $\mathcal{T}_1 \times \dots \times \mathcal{T}_n$ when we restrict the set of trees to \mathcal{S} . Thus, $\Pr(\mathcal{S} \mid f_\theta)$ represents how well \mathcal{S} works as the
 259 set of true trees.

260 Define $\mathcal{S}_0 = \bigcup_{p=1}^n \mathcal{T}_p$ as the set of all input trees. Define \mathcal{S}_{t+1} given \mathcal{S}_t as follows.

$$261 \quad \mathcal{S}_{t+1} = \mathcal{S}_t \setminus \{\operatorname{argmax}_{T \in \mathcal{S}_t} \Pr(\mathcal{S}_t \setminus \{T\} \mid f_\theta)\}. \quad (22)$$

262 In other words, \mathcal{S}_{t+1} is \mathcal{S}_t with the element removed which decreases $\Pr(\mathcal{S} \mid f_\theta)$ by the smallest quantity.
 263 Define t^* as the last t with $\mathcal{T}_p \cap \mathcal{S}_t$ non-empty for all p . Our final predicted tree T_p^* for each patient p is given
 264 by the below equation.

$$265 \quad T_p^* = \operatorname{argmax}_{T \in \mathcal{S}_{t^*} \cap \mathcal{T}_p} \Pr(T \mid f_\theta) \quad (23)$$

266 In simple terms, this procedure works by iteratively removing trees from the set of possible trees while
 267 predicting the maximum likelihood tree for each patient until no more trees can be removed from the set of
 268 possible trees. This then results in predictions T_1^*, \dots, T_n^* .

269 **A.3 Bootstrapping**

270 One way to verify the stability of predictions is to perform bootstrapping on the patients and running CloMu
 271 on the bootstrapped data. Specifically, one can randomly sample patients with replacement to produce a
 272 new datasets of the same size as the original data set. After running CloMu on bootstrapped instances, one
 273 can produce predictions on each instance such as causality or fitness. The distribution of predictions across
 274 bootstrapped instances gives information about the confidence of these predictions.

Row	Measurement	Mutation to be acquired by clone		
		Mut. $t = 2$	Mut. $t = 3$	Mut. $t = 4$
1	Mutation rates $\lambda(\mathbf{c}_0, t)$	1	1	1
2	Mutation rates $\lambda(\mathbf{c}_1, t)$	5	10	15
3	Absolute causality $A(s, t)$	$\log(5/1)$	$\log(10/1)$	$\log(15/1)$
4	Mutation probabilities for \mathbf{c}_1	1/6	1/3	1/2
5	Mutation rates $\lambda(\mathbf{c}_2, t)$	0	1	1
6	Mutation probabilities for \mathbf{c}_2	0	1/2	1/2
7	Mutation rates $\lambda(\mathbf{c}_3, t)$	1	0	1
8	Mutation probabilities for \mathbf{c}_3	1/2	0	1/2
9	Mutation rates $\lambda(\mathbf{c}_4, t)$	1	1	0
10	Mutation probabilities for \mathbf{c}_4	1/2	1/2	0
11	Averaged mutation probabilities	7/24	8/24	9/24
12	Relative causality $R(s, t)$	$\log(4/7) < 0$	$\log(1/3) = 0$	$\log(4/3) > 0$

Table S1: **An example demonstrating absolute and relative causality in the presence of a highly fit mutation.** We consider a theoretical dataset with $m = 4$ mutations including a highly fit mutation $s = 1$. For the sake of simplicity, we assume mutation 1 occurs with a rate $\lambda(\mathbf{c}, 1)$ of 0 on all clones \mathbf{c} . Row 3 is calculated as the log of row 2 divided by row 1. Row 4 is calculated from row 2 by dividing each column by the sum of the three columns. This is calculating $\lambda(s, t) / \sum_r \lambda(s, r)$. The same calculation gives row 6 from row 5, row 8 from row 7 and row 10 from row 9. Row 11 is the average of rows 4, 6, 8, and 10. Row 12 is the log of row 4 divided by row 11. Mutation 1 is highly fit as reflected by the high rates in row 2 (as compared to rows 1, 5, 7, and 9). Consequently, it causes all other mutations in terms of absolute causality (row 3). However, in terms of relative causality, it causes mutation 4, inhibits mutation 2, and has no effect on mutation 3 (row 12).

275 B Supplemental results for simulations

276 We discuss the setup of our simulations in Supplemental Material B.1. Next, we discuss the baseline meth-
277 ods and their parameters in Supplemental Material B.2. Finally, we provide detailed simulation results in
278 Supplemental Material B.3.

279 B.1 Setup

280 **Simulations (I-a).** In the Main Text, we discussed simulations with $n = 500$ patients, $m = 10$ muta-
281 tions including 5 driver mutations, and random positive causal relationships between driver mutations. To
282 generate the set \mathcal{T}_p of trees for each patient p , we simulated bulk sequencing using five sequencing sam-
283 ples. Specifically, we determined mixture proportions for each clone in each sample uniformly at random,
284 obtained the corresponding mutation frequencies, and then enumerated all possible trees that could have
285 possibly generated these mutation frequencies (using Prüfer sequences (Prüfer, 1918)). These simulation
286 instances are referred to as Simulations (I-a). The rate multipliers were chosen to be 11. That is, log rates
287 were increased by $\log(11)$ so that rates are increased by a factor of 11. The rationale is that increasing the

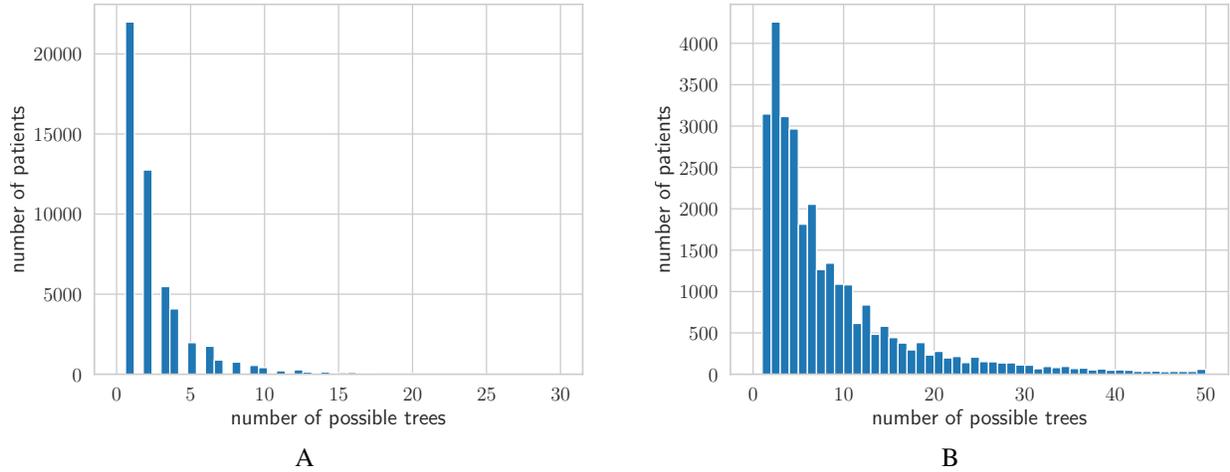


Figure S3: **Number of trees in simulated data.** Histogram of the number of trees for each patient in (A) Simulations (I-a,b,c) used to assess tree selection and causality. (B) Simulations (II) used to assess mutation interchangeability and evolutionary pathway identification.

288 rate of a single mutation by a factor of 11 results in the overall rate of mutations on a clone doubling if there
 289 exists a total of 10 mutations. The initial rate of mutations on the clone is $10 = 1 + 9$, and the new rate is
 290 $20 = 11 + 9$. Thus, a driver mutation that causes N other mutations has a fitness $(1 + N)$ times that of a
 291 passenger (which is very realistic as shown in our analysis fitness in breast cancer and AML datasets).

292 **Simulations (I-b) and (I-c).** To simulate causality with interchangeable mutations, we considered each
 293 mutation as a driver mutation. We partitioned the set of m mutations into 5 pairwise disjoint sets S_1, \dots, S_5
 294 of equal size. For each pair (S_i, S_j) of mutation sets (where $i, j \in \{1, \dots, 5\}$), we determined S_i as causing
 295 S_j , S_i as inhibiting S_j and S_i and S_j having no causal effect with uniform probabilities of $1/3$ for each
 296 case. With a slight abuse of notation, we set the rate multiplier $g(S_i, S_j)$ to 11 if S_i causes S_j , $1/11$ if S_i
 297 inhibits S_j , and 1 if there is no causal relationship from S_i to S_j . The resulting rate multiplier $g(s, t)$ of
 298 mutations $s, t \in [m]$ equals the rate multiplier $g(S_i, S_j)$ of their comprising sets (i.e., $s \in S_i$ and $t \in S_j$).
 299 Similarly, we computed rates $\lambda(c, t)$ as in the Main Text and generated a ground-truth T_p^* for each patient
 300 $p \in [n]$ by first drawing the number ℓ of mutations of tree T_p^* uniformly from $\{5, 6, 7\}$. Subsequent bulk
 301 sequencing simulation with five sequencing samples, as discussed above, resulted in sets $\mathcal{T}_1, \dots, \mathcal{T}_n$ of trees
 302 per patient (Fig. S3A). We varied the number $m \in \{10, 15\}$ of mutations. Consequently, the size of each
 303 set S_i will either be 2 (for $m = 10$) or 3 (for $m = 15$). We simulated $n = 400$ patients for $m = 10$ and
 304 $n = 600$ patients for $m = 20$. For each setting of m , we generated 20 simulation instances. These instances

305 are denoted as Simulations (I-b) and (I-c) (see Table 2).

306 **Simulations (II).** For evolutionary pathway assessment, we generated simulation instances that consist of
307 one pathway (25% probability) or two pathways (75% probability). We considered $m = 20$ mutations, some
308 of which will occur in pathways and are drivers and the remaining mutations are passengers. Each pathway
309 consists of three pairwise disjoint sets S_1, S_2, S_3 of mutations. In the case of two pathways all mutation sets
310 are pairwise disjoint. The number of mutations in a set S_i is drawn uniformly from $\{1, 2, 3\}$. As such, the
311 number of drivers ranges from 3 mutations (in the case of a single pathway with one mutation per set) to 18
312 mutations (in the case of two pathways with three mutations per set). To generate trees according to these
313 evolutionary pathways we defined rates $\lambda(\mathbf{c}, t)$ according to the stage and the pathway a particular clone is
314 at. Specifically, clone \mathbf{c} is at stage $i \in \{0, \dots, 3\}$ of the pathway S_1, S_2, S_3 if clone \mathbf{c} contains at least one
315 mutation from each set S_1, \dots, S_i and does not contain a mutation from set S_{i+1} (here if $i = 3$ then $S_4 = \emptyset$).
316 We set the rate $\lambda(\mathbf{c}, s)$ to 5 if \mathbf{c} is at stage 0 for both pathways and $s \in S_1$ for some pathway. We set the
317 rate $\lambda(\mathbf{c}, s)$ to 20 if \mathbf{c} is at stage $i > 0$ for some pathway and $s \in S_{i+1}$ for that pathway. If neither condition
318 is met, we set rate $\lambda(\mathbf{c}, s) = 1$. Note, by definition, the normal clone \mathbf{c} is at stage 0 for both pathways.
319 Additionally, if a clone is at stage 3 then all mutations have rate 1. Similarly to the other simulations, we use
320 these rates to generate a ground-truth T_p^* for each patient $p \in [n]$ by first drawing the number ℓ of mutations
321 of tree T_p^* uniformly from $\{5, 6, 7\}$. Subsequent bulk sequencing simulation with five sequencing samples,
322 as discussed above, resulted in sets $\mathcal{T}_1, \dots, \mathcal{T}_n$ of trees per patient (Fig. S3B). We set $n = 500$ patients, and
323 generated a total of 30 simulation instances. These instances are denoted as Simulations (II) (see Table 2).

324 **Simulations (III).** In addition to newly generated simulation instances, we also considered the previously
325 published RECAP simulated dataset. These simulations have a small number of patient subtypes each of
326 which has a single phylogeny tree. Specifically, within the RECAP simulated data, two groups of simulated
327 datasets are used. One group has 5 mutations per patient and 5 mutations total. The second group has 7
328 mutations per patient and 12 mutations total. For each dataset (in each group), there are $k \in \{1, \dots, 5\}$
329 clusters among $n \in \{50, 100\}$ patients. Each patient cluster corresponds to one true tree. In the published
330 RECAP simulations, each patient was randomly assigned to a true tree, then the set of possible trees is
331 generated by simulated bulk frequency measurements. The goal in this simulated dataset is to assign the
332 correct tree to each patient. We denote these instances as Simulations (III) as shown in Table 2.

333 **Simulations (IV).** Finally, we considered previously published TreeMHN simulations. Specifically, we
334 restricted our comparison to instances with $n = 300$ patients, and randomly selected 20 instances for each
335 number $m \in \{10, 15, 20\}$ of mutations. These simulations instances are referred to as Simulations (IV).

336 **B.2 Baseline methods and parameters**

337 We compared CloMu to TreeMHN (Luo et al., 2023), RECAP (Christensen et al., 2020), REVOLVER (Car-
338 avagna et al., 2018), GeneAccord (Kuipers et al., 2021), and HINTRA (Khakabimamaghani et al., 2019).
339 In the following, we will discuss which methods were included for each class of simulation instances, list
340 method parameters and provide details on any post-processing that was conducted.

341 **Methods per simulation class.** We used Simulations (I-a) to assess tree selection and causality inference.
342 For the former, we compared CloMu to RECAP and REVOLVER. For the latter, we included all methods in
343 our benchmarking. We used Simulations (I-b) and (I-c) to further asses the differences between CloMu and
344 TreeMHN in causality prediction in the presence of interchangeable mutations. Simulations (II) were used
345 to demonstrate the detection of evolutionary pathways of interchangeable mutations. The methods RECAP
346 and REVOLVER had their predictions post-processed in order to form a baseline to compare with CloMu.
347 Simulations (III) are from the RECAP paper and further asses tree selection. On these simulations, CloMu is
348 compared against results form the RECAP paper including results for RECAP, REVOLVER, and HINTRA.
349 Simulations (IV) are from the TreeMHN paper and give an additional comparison between TreeMHN and
350 CloMu for the task of predicting causal relationships.

351 **Method parameters.** On Simulations (IV), TreeMHN was ran with default parameters and the use of
352 stability selection. For comparison on Simulations (I-a,b,c), we ran TreeMHN with the penalty parameter
353 set to maximize precision and recall. Additionally, we did not use stability selection, since improving
354 precision at the cost of recall would not improve their results. We ran RECAP and REVOLVER with their
355 default parameters. We ran GeneAccord with default parameters. We did not run HINTRA due to its lack
356 of scalability. However, we did compare with HINTRA on RECAP’s data (Simulation (III)), which was ran
357 with default parameters.

Predicted relationship	True relationship		
	No relationship	Positive relationship	Negative relationship
No relationship	True Negative	False Negative	False Negative
Positive relationship	False Positive	True Positive	False Positive
Negative relationship	False Positive	False Positive	True Positive

Table S2: This table demonstrates how true positives, true negatives, false positives and false negatives are determined based on predicted and true causal relationships in the case of signed and bidirectional relationships.

358 **Post-processing.** For Simulations (I-a) and (II), which assess causality and evolutionary pathway infer-
359 ence, respectively, we included RECAP and REVOLVER despite their implementation not directly support-
360 ing these tasks. RECAP and REVOLVER are consensus tree methods (Main Text Table 1), which select one
361 tree per patient. To support causality inference for RECAP and REVOLVER, we applied a post-processing
362 step. Specifically, we processed RECAP’s and REVOLVER’s selected trees for each patient from the set
363 of possible trees. We then extracted all edges from these trees and determined the frequency of each edge.
364 For Simulations (I-a), we considered there to be a positive causal relationship from mutation s to mutation
365 t if the edge (s, t) occurred with a frequency above a fixed threshold. We used a threshold of 40 based on
366 the number $n = 500$ patients for Simulations (I-a). For Simulations (II), the goal is to predict the edges
367 which occur in the simulated pathways. Specifically, if a ground-truth pathway contains sets S_i and S_{i+1} ,
368 we say that pathway contains all edges (s, t) with $s \in S_i$ and $t \in S_{i+1}$. Given a threshold frequency, RE-
369 CAP/REVOLVER predicts an edge to exist for a pathway in a simulation if that edge occurs in the selected
370 trees with a frequency above the threshold. We used a threshold of 20 based on the number $n = 500$ of
371 patients for Simulations (II).

372 B.3 Results

373 **Evaluating bidirectional and signed causal relationships.** To evaluate bidirectional causality for signed
374 causal relationships, one must calculate precision and recall from predicted and true causal relationships.
375 First, every ordered pair of mutations is considered. Then, if some causal relationship (positive or negative)
376 from s to t is predicted correctly, this is considered a true positive. Similarly, if no causal relationship is
377 predicted from s to t and there is no true causal relationship from s to t , then this is considered a true negative.
378 If there is a predicted causal relationship (positive or negative) from s to t , and this predicted relationship
379 is incorrect for any reason, then it is considered a false positive. Note, this implies that if the predicted

380 causal relationship is positive (or negative) and the true causal relationship is negative (or positive), then it is
381 considered a false positive, despite there being a true causal relationship. This information is summarized in
382 Table S2. Using this, each ordered pair of mutations can be classified as either a true positive, true negative,
383 false positive or false negative. Summing these numbers across all ordered pairs of mutations gives total true
384 positives, true negative, false positives and false negatives. Then, the precision is defined as the true positives
385 divided by the sum of true positives and false positives. The recall is defined as the true positives divided by
386 the sum of true positives and false negatives. These definitions follow those used in TreeMHN (Luo et al.,
387 2023).

388 **B.3.1 Bootstrapping – Simulations (I-a)**

389 We applied bootstrapping to the last simulation instance from Simulations (I-a), for a total of 100 times.
390 We recorded the median as well as the 90% confidence interval of the predicted absolute causality values.
391 Fig. S4 shows predictions on the instance from Simulations (I-a), showing that the lower bound for pairs of
392 mutations with a true causal relationship (Fig. S4C) vastly exceeds the upper bound for pairs of mutations
393 without a true causal relationship (Fig. S4D). Additionally, we ranked the strength of the predicted causal
394 relationship for the 90 ordered pairs of mutations. There were 10 actual causal relationships, and thus having
395 a ranking of at least 81 is consistent with there being a causal relationship and having a ranking of at most
396 80 is consistent with there being a non-causal relationship. Fig. S5 shows the 90% confidence interval on
397 the predicted causal relationship strength rankings. In all cases on all 100 simulation instances, actual causal
398 relationships have a ranking at least 81 and non-causal relationships have a ranking at most 80, indicating
399 perfectly correct rankings on all instances.

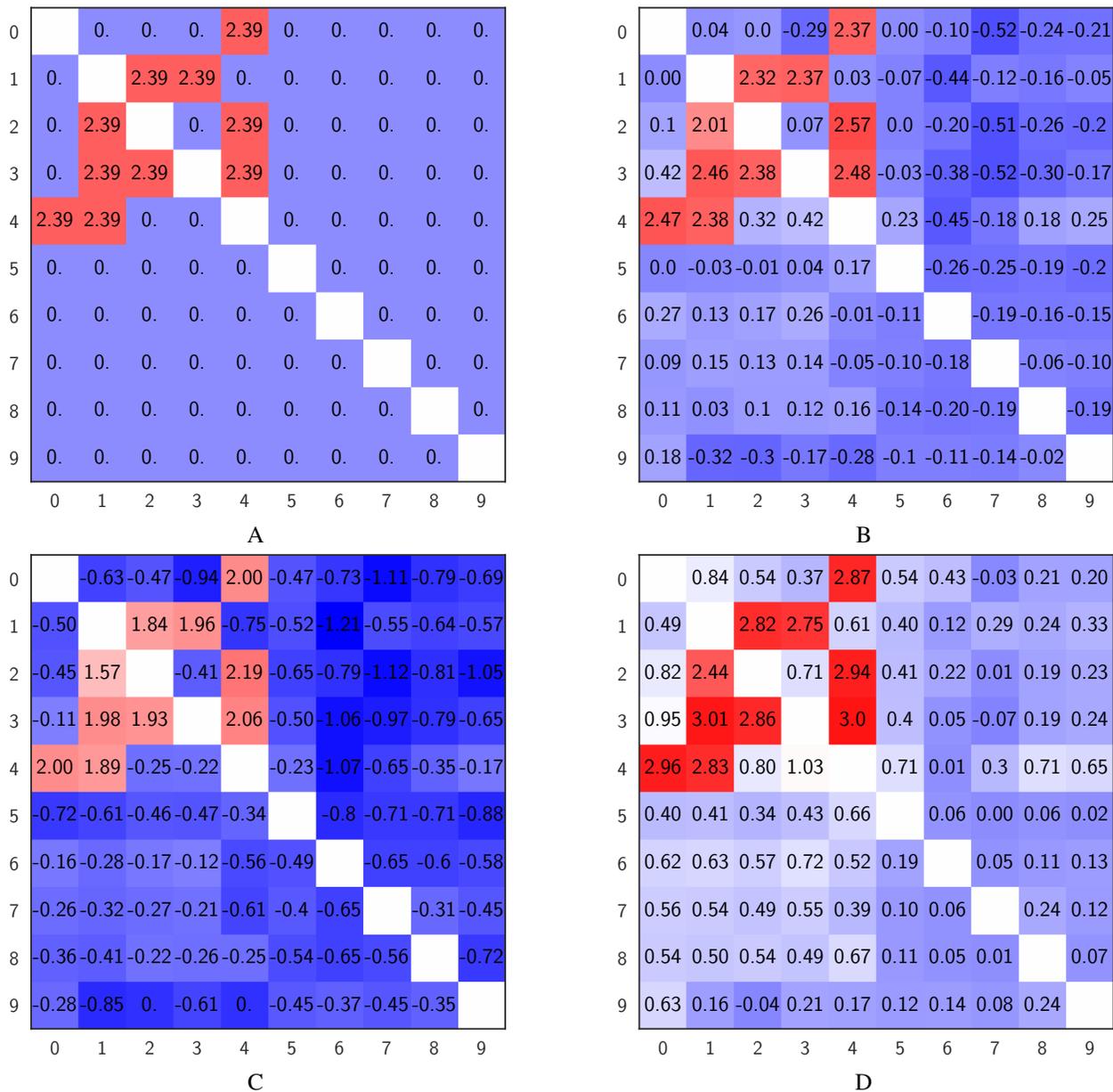


Figure S4: **Bootstrapping with 100 replicates demonstrates CloMu’s robust absolute causality predictions on a Simulations (I-a) instance.** True and predicted causal relationships are shown. Mutations are sorted such that the first 5 mutations are driver mutations and the remaining 5 mutations are passenger mutations. The 90% confidence interval for each mutation contains 90 values from 90 simulation instances and excludes the highest and lowest 5 values. (A) Ground-truth absolute causality values. (B) The median absolute causality predictions on 100 bootstrapped instances. (C) The 6th lowest absolute causality predictions. (D) The 6th highest absolute causality predictions. The lower bound for pairs of mutations with a true causal relationship vastly exceeds the upper bound for pairs of mutations without a true causal relationship.

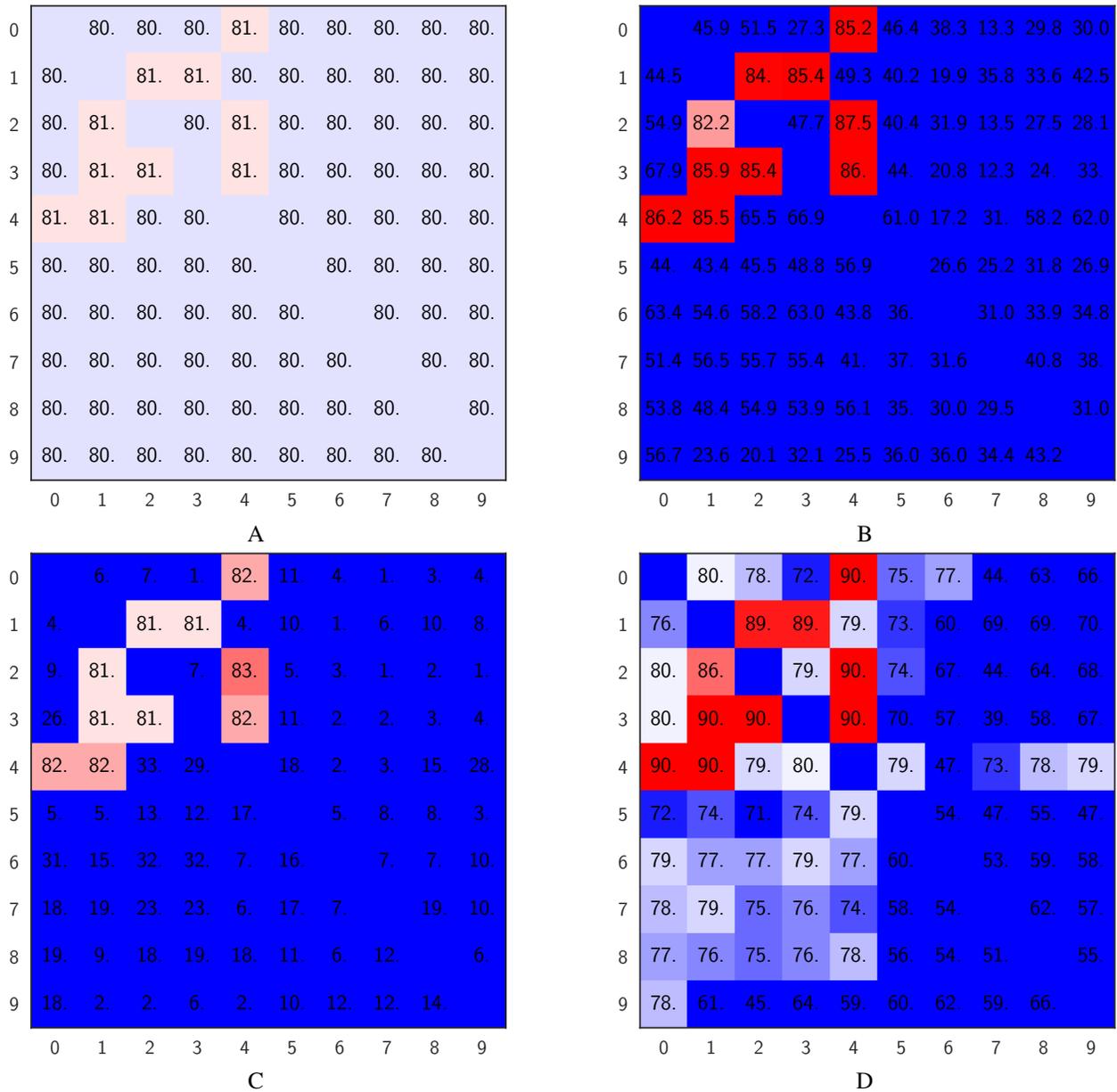


Figure S5: **The absolute causality strength rankings are shown to be robust via bootstrapping with 100 replicates on a Simulations (I-a) instance.** Causal relationships are ranked from 1 to 90 based on the absolute causality value. There are 10 ground truth causal relationships, and so a ranking of at least 81 indicates a causal relationship and a ranking of at most 80 indicates no causal relationship. To emphasize rankings near these values, the colormap ranges from $76 = 81 - 4$ to $85 = 81 + 4$. True and predicted causal relationships are shown. Mutations are sorted such that the first 5 mutations are driver mutations and the remaining 5 mutations are passenger mutations. The 90% confidence interval for each mutation contains 90 rankings from 90 simulation instances and excludes the highest and lowest 5 rankings. (A) Ground-truth rankings are shown. Specifically, causal relationships are given the lowest valid ranking for a causal relationship (81), and non-causal relationships are given the highest valid ranking for non-causal relationships (80). There is no ground-truth difference in ranking among causal relationships or among non-causal relationships. (B) The median absolute causality rankings on 100 bootstrapped instances. (C) The 6th lowest absolute causality rankings. (D) The 6th highest absolute causality rankings.

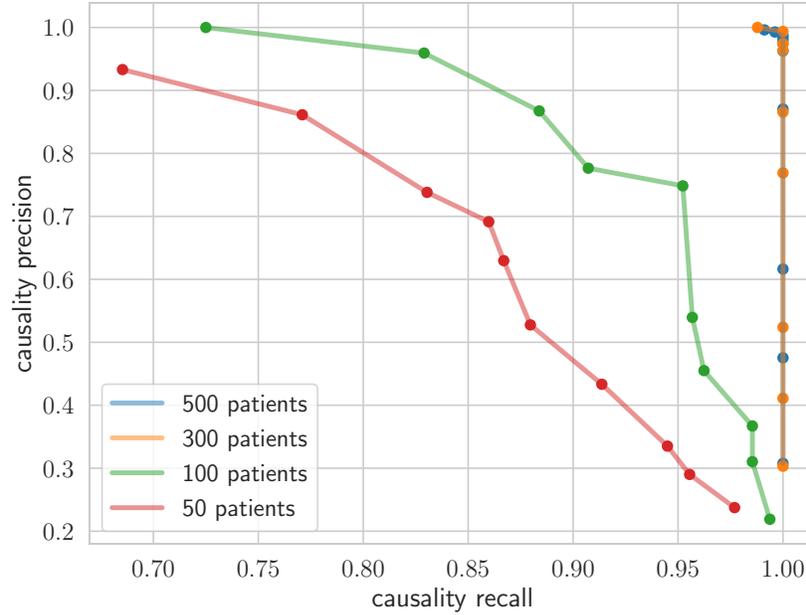


Figure S6: **CloMu gives accurate predictions even on simulations with few patients.** We consider a modification of Simulations (I-a), reducing the number n of patients from 500 (blue) to 300 (orange), 100 (green) and 50 (red). The resulting precision and recall values are somewhat reduced but remain high. The curve for 500 patients and 300 patients overlap since they are both nearly perfect.

400 B.3.2 Decreasing numbers of patients – adaptation of Simulations (I-a)

401 While the breast data had $n = 1756$ patients (Razavi et al., 2018), typically the number of patients in current
 402 cohort-level cancer phylogeny data consist of $n \approx 100$ patients (Morita et al., 2020; Jamal-Hanjani et al.,
 403 2017; Turajlic et al., 2018). However, Simulations (I-a) contain $n = 500$ patients per simulation instance.
 404 To assess CloMu performance with fewer patients, we generated new simulations with the same exact setup
 405 as Simulations (I-a), but with $n \in \{50, 100, 300\}$ patients per simulation instance. More specifically, for
 406 each number n of patients we generated 20 simulation instances. Fig. S6 demonstrates that although the
 407 performance of CloMu in terms of causality precision and recall decreased with decreasing number of
 408 patients, it remained acceptable.

409 The CloMu software has very reasonable runtimes and memory requirements. Consequently, for all
 410 experiments we trained CloMu on a laptop with a 2.4 GHz CPU and 64 GB of RAM without using a
 411 GPU. To demonstrate the effect of the number of patients on computational requirements, we plotted the
 412 memory usage and runtime of CloMu on simulations based on Simulations (I-a) but with differing numbers
 413 $n \in \{50, 100, 300, 500\}$ of patients. These results are shown in Fig. S7. The memory usage values are each
 414 for a single instance, and the runtimes are averaged over 20 instances. The mean running time ranged from

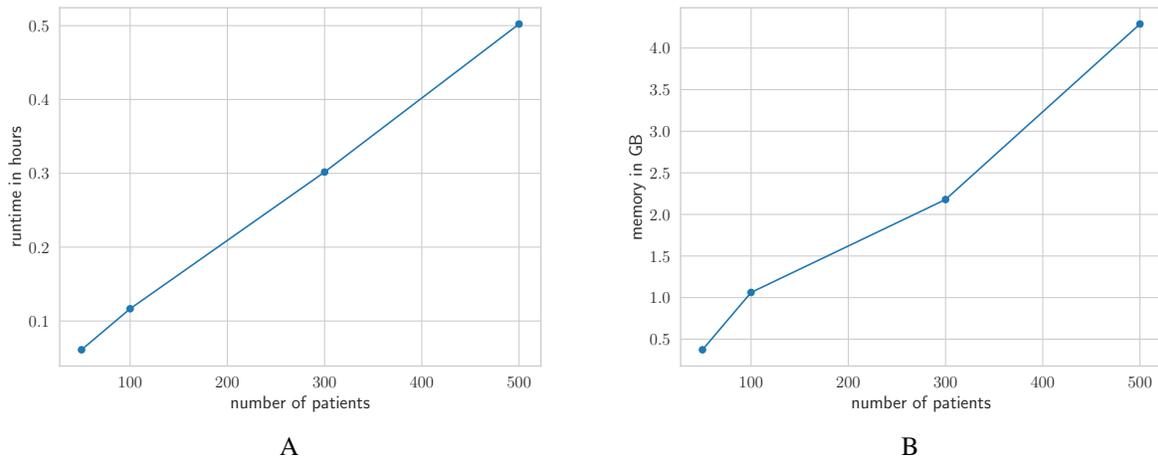


Figure S7: **The runtimes and memory usage of CloMu ran on instances based on Simulations (I-a) with differing numbers of patients.** Simulations were generated following Simulations (I-a) but with $n \in \{50, 100, 300, 500\}$ patients. (A) The mean runtime of CloMu across 20 instances for varying number n of patients. (B) The memory usage of CloMu for a single instance for each number n of patients.

415 3 minutes for $n = 50$ patients to 30 minutes for $n = 500$ patients. The memory usage was 382 MB for
 416 $n = 50$ patients and 4.288 GB for $n = 500$ patients. The runtime and memory usage values seem to be
 417 approximately linear in the number of patients.

418 B.3.3 Decreasing numbers of mutations – adaptation of Simulations (I-a)

419 Simulations (I-a) also contain a randomly selected number of mutations per patient, ranging from 5 to 7.
 420 However, in the AML dataset (Morita et al., 2020), which performed panel-based single-cell DNA sequenc-
 421 ing, the number of mutations per patient can be much smaller, reducing the amount of information provided
 422 by each patient. To assess performance in this regime, we generated new 20 simulations based on Simula-
 423 tions (I-a) with $n = 100$ patients and the number of mutations per patient generated based on the distribution
 424 observed in the AML dataset. Since there are only $m = 10$ total mutations in Simulations (I-a), we did not
 425 allow the number of mutations in any patient to reach too close to 10. Thus, we set the maximum number of
 426 mutations per patient to 7. Fig. S8A shows this distribution of the number of mutations per patient (where
 427 all values above 7 are set to 7). Fig. S8B shows that CloMu performance is similar to the original Simulation
 428 (I-a) instances with $n = 100$ patients.

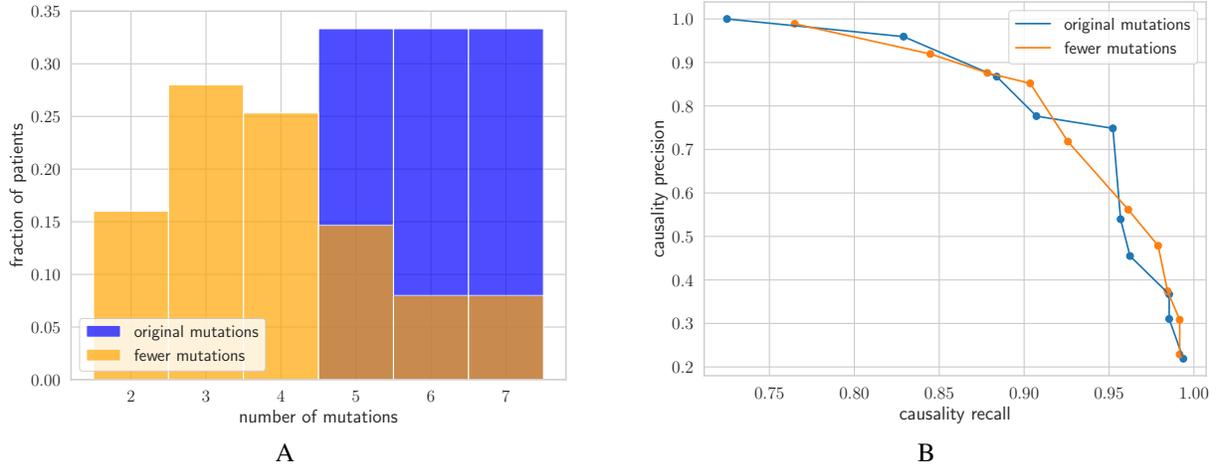


Figure S8: **CloMu’s performance remains high for simulations with a reduced number of mutations per patient.** We assessed the effect of a reduced number of mutations per patient for simulations with $n = 100$ patients. (A) The original (blue) and modified distributions (orange); the latter are based on the AML dataset (Morita et al., 2020). (B) The causality precision and recall remain high.

429 B.3.4 Causation strength based on real data – adaptation of Simulations (I-a)

430 In Simulations (I-a), the strength of the absolute causality was set to $\log(11)$ such that the rate of a mutation
 431 is increased by a factor of 11 if it is caused by another mutation. However, in real cancer evolution, the
 432 strength of all causal relationships are not identical. To improve realism, we determined a distribution of the
 433 strength of causal relationships using the real data. On the AML dataset (Morita et al., 2020), we considered
 434 the strongest causal relationship caused by each mutation as determined by CloMu. Then, we picked the top
 435 5 mutations with the highest value for their strongest (absolute causality) causal relationship. For these 5
 436 mutations, we observed the absolute causality value for their top 5 strongest causal relationships. This gives
 437 a total of 25 causal relationships from the AML dataset. A histogram of these values is shown in Fig. S9A.
 438 We generated 20 new simulation instances based on Simulations (I-a), but with causal relationship strengths
 439 sampled from this distribution and with the number n of patients reduced to 300. As shown in Fig. S9B,
 440 CloMu is still able to give accurate causal relationship predictions.

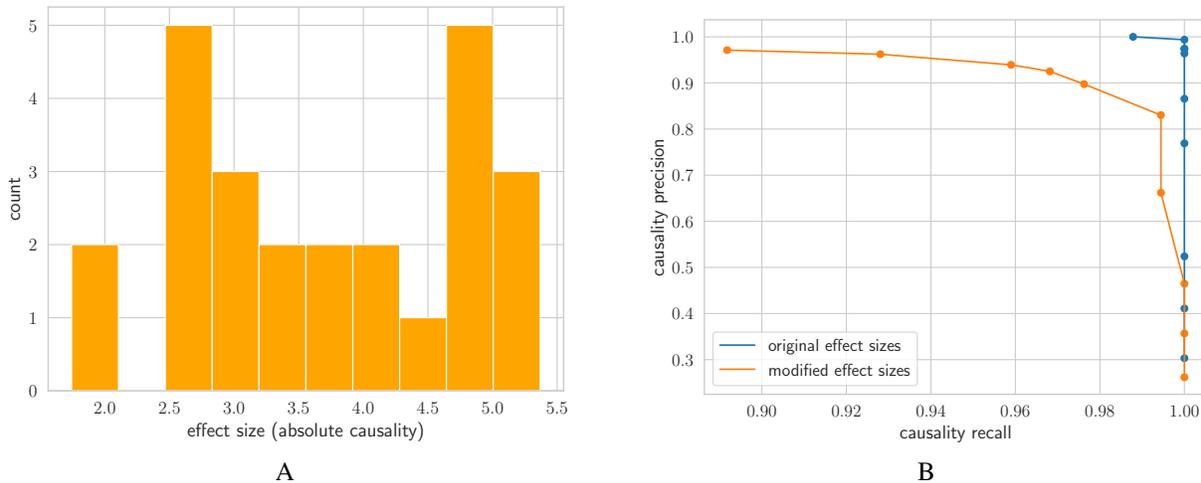


Figure S9: **CloMu’s performance is slightly reduced when using varying effect sizes in causality.** We assessed the impact of varying effect sizes in causality using simulations with $n = 300$ patients. (A) The used distribution of effect sizes, again based on the AML dataset. (B) CloMu has a small reduction in causality precision and recall when the effect size is modified to match real data.

441 B.3.5 Large numbers of passenger mutations – adaptation of Simulations (I-a)

442 In this section, we introduce simulations with a large number of passenger mutations that do not affect the
 443 tumor. In all of the other simulations CloMu was tested on, the total number of mutations is at most 20.
 444 However, in real cancer data, there is often a much larger number of mutations, of which the vast majority are
 445 passengers. Fig. S10A shows a histogram of the fraction of $n = 1250$ breast cancer patients (with at most 10
 446 mutations) for which each of the $m = 410$ mutations occurs (Razavi et al., 2018). We generated simulations
 447 that are identical to Simulation (I-a), but with the same number of mutations as the breast cancer dataset
 448 and with each mutation occurring with a baseline rate proportional to the frequency at which they occur
 449 in the breast cancer set. Additionally, we set the 5 driver mutations to be the 5 highest baseline likelihood
 450 mutations and set all other mutations as passengers. Fig. S10B shows the causality precision and recall
 451 of this simulation. The accuracy was somewhat reduced by the high number of passenger mutations but
 452 remains acceptable. The primary reason for the reduction in accuracy is the difficulty of differentiating the
 453 tiny number of positive relationships from the vast number of pairs of mutations with no causal relationships.
 454 Specifically, there are at most 20 true causal relationships and at least $410^2 - 410 - 20 = 167,670$ non-
 455 causal relationships. If non-causal relationships were falsely identified as causal in even 1/1000 of cases,
 456 there would be over 8 times as many false positives as true positives (resulting in a precision below 0.11).

457 In addition to testing how the number of patients affects the computational requirements as discussed

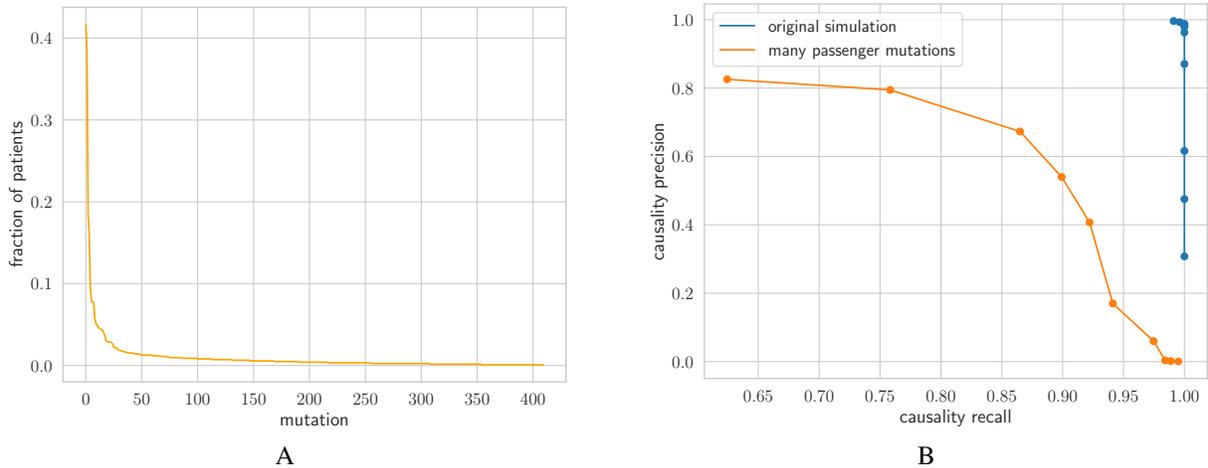


Figure S10: CloMu’s performance is reduced when including a large number of passenger mutations. We assessed the impact of including a large set of passenger mutations among $n = 500$ patients. (A) The frequency of different mutations in the breast cancer dataset (Razavi et al., 2018) that we used to determine the default rate at which mutations occur (noting that the rates can then be modified by the presence of driver mutations). (B) CloMu has somewhat reduced accuracy when 405 passenger mutations are included. However, the performance is still acceptable, given that even a 99.9% accuracy at identifying non-causal relationships would give a precision below 0.11.

458 in Supplementary Material B.3.2, we recorded the effect of the number of mutations. Specifically, for the
 459 above set of simulations discussed in this section, we dramatically increased the number of mutations from
 460 10 to 410. The new runtime for a single instance is 3 hours and 5 minutes, and the new memory usage is
 461 5.378 GB compared to 30 minutes and 4.289 GB for $n = 300$ patients and $m = 10$ mutations.

462 B.3.6 Violations of independent clonal evolution – adaptation of Simulations (I-a)

463 CloMu is built off of the independent clonal evolution assumption that only the mutations on a clone will
 464 affect the probability of new mutations occurring on that clone. This assumption may be generally realistic
 465 but could have exceptions in real tumor evolution. Therefore, ideally, we would want CloMu to be somewhat
 466 robust to this possibility. We tested this by using a simulation where the causal relationships between
 467 mutations also occur to a smaller extent between mutations in different clones. In other words, if mutation
 468 s being in a clone causes mutation t to be much more likely to occur in that clone, then mutation s being in
 469 any clone will cause mutation t to be somewhat more likely to occur in any, not necessarily the same, clone.

470 We found that CloMu still performs well on a causal relationship simulation when clonal interactions
 471 are included in the simulation. Specifically, our 20 simulation instances result from a modified version
 472 of Simulations (I-a) with no interchangeable mutations, 5 driver mutations, and 5 passenger mutations.

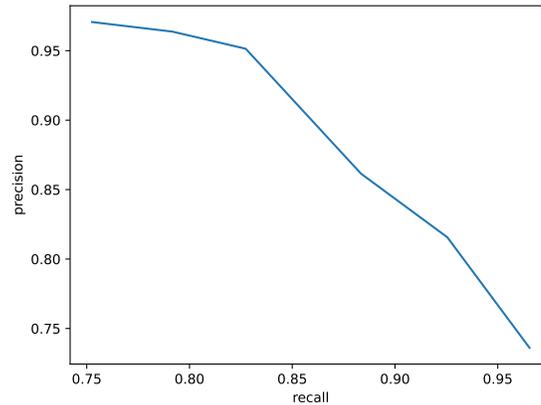


Figure S11: **CloMu is robust to violations of the independent clonal evolution assumption.** We consider a modification of Simulations (I-a) that results in a violation of the assumption. The resulting causality precision and recall values are slightly reduced but remain very high, as for the original Simulations (I-a) results shown in Main Text Fig. 2B).

473 If mutation s causes mutation t , then as a higher proportion of clones have mutation s , the probability of
 474 mutation t on all clones will increase. Specifically, if a proportion p of clones have some mutation s , then the
 475 log rate of mutation t will be increased by $p \log(2)$. Applying CloMu on this simulation gives a causality
 476 precision-recall curve shown in Fig. S11, with values close to 1, similar to the original Simulations (I-
 477 a) results shown in Main Text Fig. 2B). Hence, for this scenario, CloMu was robust to violations of the
 478 independent clonal evolution.

479 **B.3.7 Including highly incorrect possible trees – adaptation of Simulations (I-a)**

480 One important ability of CloMu is accounting for the fact that the set of possible trees for each patient
 481 can contain completely incorrect trees. In all of our simulations, we simulate bulk sequencing to add tree
 482 uncertainty. However, here we explicitly test CloMu’s ability to deal with entirely random incorrect trees
 483 being included in the set of possible trees. We generate a set of simulations with the same setup as Simu-
 484 lations (I-a), with the exact same procedure for generating the true tree for each patient. However, now the
 485 set of possible trees for each patient contains both the true tree and a randomly generated tree with the same
 486 set of mutations as the true tree. The random tree is generated by starting with a normal clone, and then
 487 iteratively adding mutations (from the true tree) to clones in order to form a tree. Mutations are added to
 488 clones uniformly at random, with no causal relationships. As a more extreme case, we also generated 20
 489 simulation instances in which there exists 2 random trees per patient in addition to the one true tree. Fig. S12

490 shows the precision recall curve on these new sets of simulations. In the case of one random tree per patient,
491 the accuracy is only slightly reduced by half of the trees being random. In the case of two random trees per
492 patient, the accuracy is noticeably reduced by the majority of trees being random. We also tested CloMu's
493 tree selection capability. In the case with one random tree per patient, CloMu achieved an average tree
494 selection accuracy of 81.82% across 20 simulation instances. In the case with two random trees per patient,
495 CloMu's average tree selection accuracy was 57.50 % (in comparison to 33.33% that would be achieved
496 by random selection). To investigate tree selection further, we define the *relative probability of a candidate*
497 *tree* for a patient as the probability of the candidate tree divided by the sum of probabilities of all candidate
498 trees for the patient. Fig. S13A shows a histogram of the relative probability of the correct tree and incorrect
499 tree on all 20 simulation instances with one random tree per patient. Note that a relative probability above
500 0.5 indicates the correct tree would be selected. Fig. S13B shows a histogram of the relative probability of
501 the correct tree and one incorrect tree on all 20 simulation instances with two random trees per patient. Note
502 that, the average relative probability of all possible trees is 1/3 in this case, since there are now 3 possible
503 trees per patient.

504 As a control case, we also generated simulations in which both trees are randomly generated. In this
505 case, no causal relationships are predicted for the cutoff value $\tau = 0.8$, which allows for perfect 100% recall
506 on the previously described simulations with one random tree per patient (and one true tree generated as in
507 Simulations (I-a)). Additionally, out of 1800 possible causal relationships on 20 simulation instances, only
508 5 causal relationships exceed an absolute causality of 0.4. This cutoff value of 0.4 also allows for precision
509 and recall values of 0.9515 and 0.9121, respectively, on the simulations where each patient has two random
510 trees and one correct tree generated as in Simulations (I-a). The few false positive weak causal relationships
511 that occur may be caused by weak patterns that occur in the trees by random chance. A histogram of the
512 relative probability of the (random) correct tree is shown in Fig. S13C. This is much more uniform since
513 the correct tree and non-correct tree are both drawn from the same random distribution. However, note that
514 the process of randomly adding mutations to clones does not generate all trees with equal probability. In
515 particular, branching topologies are generated by this random process with a higher likelihood than chain-
516 like trees (since there are more unique ways of generating branching trees). Consequently, our model, trained
517 on this random process, does not assign all trees exactly equal probability. Despite this, the probabilities
518 of different trees are much more similar to each other than under a model of evolution in which there exist
519 causal relationships.

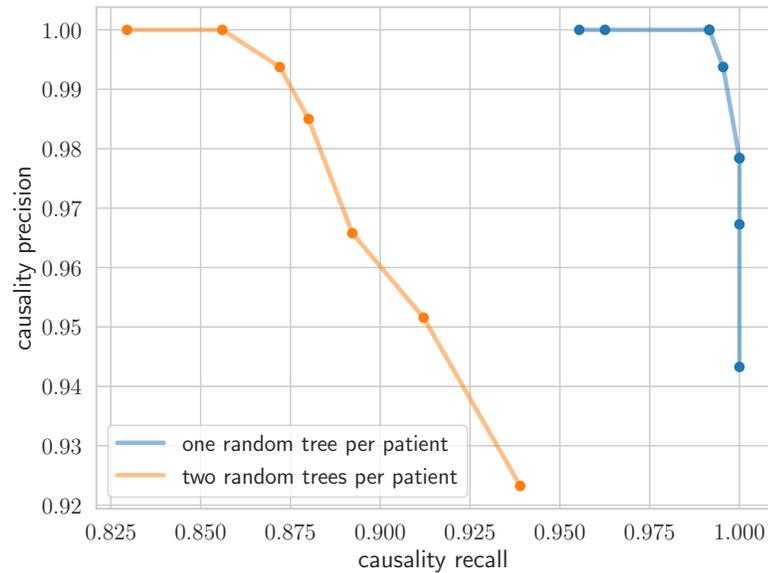


Figure S12: **CloMu is robust to random trees in the set of possible trees.** In these simulations based on Simulations (I-a), each patient contains one true tree and either one or two trees with a completely random topology on the same set of mutations as the true tree. When half of the trees are random (one random tree per patient, blue), the precision and recall remain very high. However, when the majority of the trees are random (two random trees per patient, orange), CloMu’s accuracy is noticeably diminished.

520 These simulations demonstrate the robustness of CloMu on datasets where the set of possible trees is
 521 derived in a manner that could commonly include highly incorrect trees.

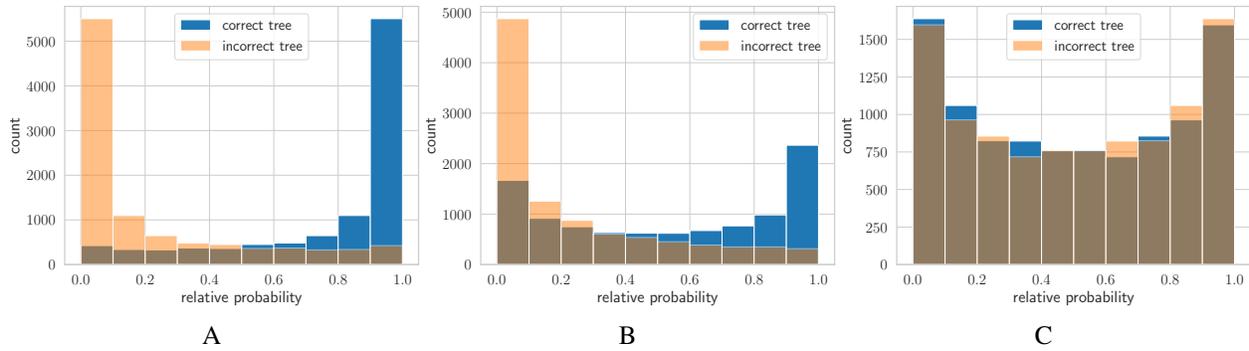


Figure S13: **CloMu assigns reasonable relative tree probabilities on simulations with random trees included.** The relative probability of correct and incorrect trees on all 20 simulation instances is shown in histograms for several simulation setups. (A) In this set of simulations, the set of possible trees for each patient contains one true tree generated as in Simulations (I-a) and one random tree. There is frequently a much higher probability assigned to the true tree than to the random tree. (B) In these simulations, the set of possible trees for each patient contains one true tree and two random trees. The probability of the true tree and one random tree is shown (since the two random trees are generated identically). The true tree typically has a higher probability than the random tree. (C) As a control case, simulations are created in which both the correct tree and incorrect tree are random. In this case, both random trees typically have a fairly similar probability.

522 B.3.8 Inhibition between interchangeable mutations – adaptation of Simulations (I-b,c)

523 Interchangeable mutations that belong to the same functional pathway have been observed to exhibit mutual
 524 exclusivity and thus inhibit each other's occurrence (Yeang et al., 2008). This effect was not assumed to
 525 be the case in Simulations (I-b) and (I-c). We generated new simulations exactly mimicking Simulations
 526 (I-b) and (I-c), but with all pairs of interchangeable mutations having negative causal relationships with
 527 each other. The performance of TreeMHN slightly improved, but the general pattern in the performance of
 528 CloMu and TreeMHN remains unchanged as shown in Main Text Fig. 2C and Fig. S14. Specifically, CloMu
 529 outperforms TreeMHN in the presence of two interchangeable mutations in each set of interchangeable
 530 mutations. Additionally, CloMu vastly outperforms TreeMHN in the presence of three interchangeable
 531 mutations in each set of interchangeable mutations.

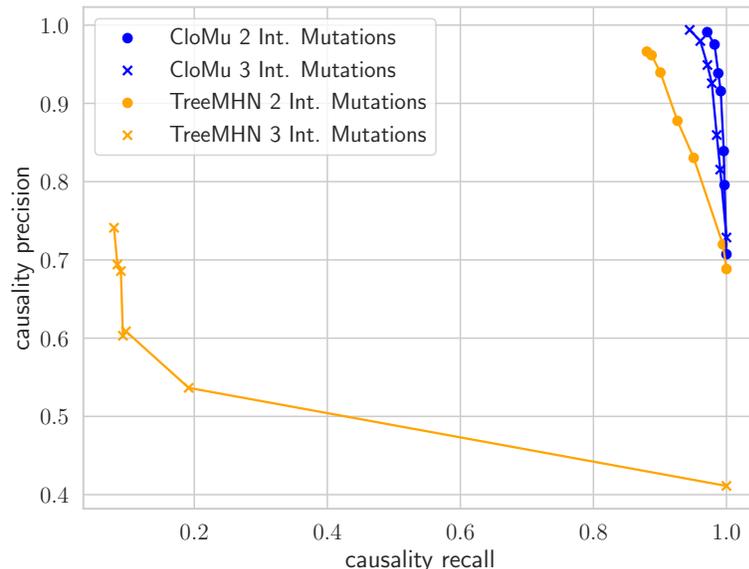


Figure S14: **Precision recall curves for simulations based on Simulations (I-b) and (I-c) still demonstrate the advantage of CloMu when inhibition between interchangeable mutations is added.** Simulations with 2 interchangeable mutations in each set of interchangeable mutations and 3 interchangeable mutations in each set of interchangeable mutations are shown. The addition of inhibition somewhat improved the performance of TreeMHN. However, CloMu still outperforms this baseline method in the presence of two interchangeable mutations in each set and greatly outperforms it in the presence of three interchangeable mutations in each set.

532 B.3.9 CloMu with increased number L of hidden neurons – Simulations (I-c)

533 CloMu’s usage of a neural network with $L = 5$ hidden neurons allows it to have less parameters than
 534 existing methods and thus resistance to overfitting. The exact sensitivity of CloMu to the value of L is
 535 therefore worth investigating. In theory, increasing L creates the potential for overfitting. However, our use
 536 of L1 regularization may result in only a subset of hidden neurons being used, thus avoiding overfitting.
 537 For instance, on our breast cancer cohort and AML cohort, only a small number of hidden neurons were
 538 used as shown in our latent representation plots (Fig. S15). To test this, we applied CloMu with $L = 20$
 539 hidden neurons to our Simulations (I-c) in which there are 3 interchangeable mutations per type and 5 types
 540 of mutations ($m = 15$). As $L = 20 > m = 15$, CloMu technically has more parameters than TreeMHN
 541 in this case – i.e., $O(Lm)$ for CloMu vs. $O(m^2)$ for TreeMHN. However, due to CloMu’s neural network
 542 structure, only a small number of CloMu’s parameters need to be used to fit the data. Fig. S15 demonstrates
 543 that CloMu with $L = 20$ performs vastly outperforms TreeMHN.

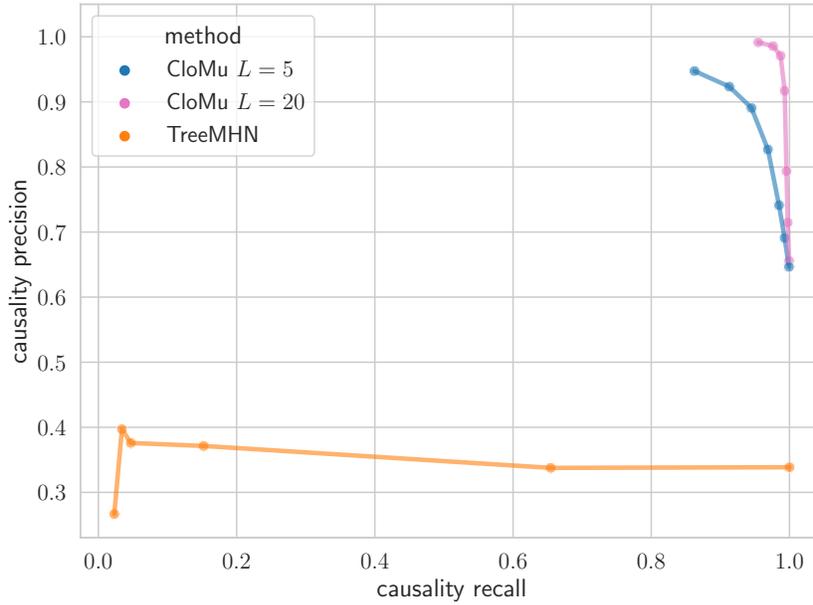


Figure S15: **CloMu is robust to increasing the number L of hidden neurons.** CloMu is ran on Simulations (I-c) with $L = 20$ as well as the original $L = 5$ hidden neurons. CloMu greatly outperforms TreeMHN for both values of L , demonstrating CloMu’s robustness to differing numbers of hidden neurons.

544 B.3.10 Mutations with shared properties

545 On Simulations (II) it was demonstrated that interchangeable mutations can be detected by determining
 546 mutations with similar latent representations. However, in real cancer evolution, two mutations may have
 547 similar or related effects on tumors without being exactly interchangeable. We generated an additional set
 548 of simulations to demonstrate the ability of our latent representations to pick up on these shared properties
 549 of mutations. Additionally, on Simulations (I-b) and (I-c) it was demonstrated that CloMu predicts causal
 550 relationships more accurately than TreeMHN in the presence of interchangeable mutations. The simulations
 551 generated in this section also demonstrate CloMu’s ability to more accurately predict causal relationships in
 552 the presence of mutations with shared properties.

553 In these simulations, there exist $m = 20$ mutations. Each mutation has two values associated with it
 554 called the value of property 1 and the value of property 2. These values range from -1 to 1 . Generally
 555 speaking, mutations with similar values of the two properties tend to cause each other. One can imagine
 556 that there exists some capability the tumor must develop in order to grow effectively (such as avoiding being
 557 stopped by the immune system) and that there are multiple independent and mutually exclusive methods of
 558 accomplishing this strategy. If property 1 represents some required feature of the tumor, then high values

559 and low values of property 1 indicate mutations that develop the tumor towards independent and mutually
560 exclusive strategies for accomplishing this feature. This results in inhibition between mutations with oppo-
561 site values of property 1 and causation between mutations with similar values of property 1. The same is
562 true for property 2. Define $P_L(s)$ to be the 2-dimensional vector containing the values of property 1 and 2
563 for some mutation s . If mutation s exists in a clone, then the rate of mutation t is multiplied by a factor of 10
564 to the power of $P_L(s) \cdot P_L(t) / \|P_L(t)\|_2$. Equivalently, the value $\log(10) P_L(s) \cdot P_L(t) / \|P_L(t)\|_2$ is added to
565 the log rate. As in previous simulations, the rate of mutation t occurring on a clone c is simply the product
566 of these rate multipliers and 0 if c already contains t .

567 For each mutation, we generate $P_L(s)$ by first generating a random ϕ between 0 and 2π , then generating
568 a random R between $1/2$ and 1, and then setting $P_L(s) = [R \cos(\phi), R \sin(\phi)]^\top$. This simulation setup
569 is intentionally somewhat more complex than our other simulations for the purpose of creating a complex
570 landscape of related mutations with partially shared yet partially differing effects. We refer to Fig. S16 for a
571 worked example.

572 After training the model, one goal is for similar mutations (with similar values of $P_L(s)$) to have similar
573 latent representations. Fig. S17A shows a scatter plot of the Euclidean distance between latent representa-
574 tions and the distance between mutation property vectors ($P_L(s)$). The Pearson correlation between these
575 two distance measurements is 0.8316, showing high concordance. Additionally, we wish to reconstruct the
576 properties $P_L(s)$ of a mutation s exclusively through observing the latent representation. This more ambi-
577 tious goal requires some post-processing to be possible. Specifically, we normalized $P_L(s)$ to have a mean
578 of zero and a norm of 1 across all mutations s . Then we applied PCA to the latent representations of all
579 mutations, reducing the dimension to 2 and normalized this representation to also have an average norm of 1.
580 This allows $P_L(s)$ and the latent representations to be of the same dimension and on the same scale. Then,
581 these two-dimensional processed latent representations are rotated by some angle and/or flipped in order to
582 reduce the Euclidean distance between these representations and P_L . This only requires two parameters (a
583 rotation angle and a Boolean of possibly flipping the first component) shared across all 20 mutations. This
584 step is necessary since there is no reason to assume by default that the first PCA component of the latent
585 representation corresponds to the first mutation property in $P_L(s)$. We define this processed version of the
586 latent representation to be the predicted mutation property.

587 We compare the predicted mutation property derived from the latent representation vector with the true
588 property vector. Fig. S17B shows CloMu gives highly accurate predictions on an example simulation.

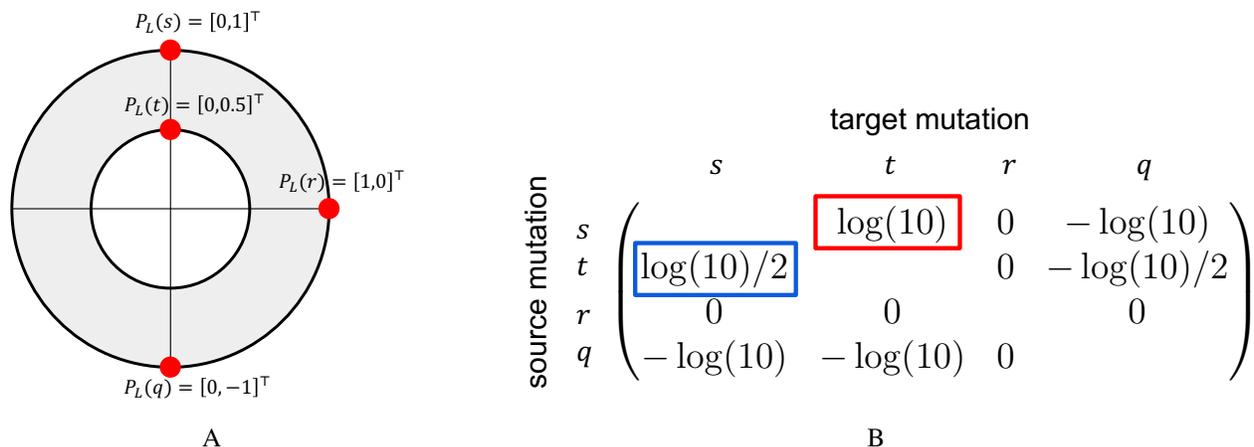


Figure S16: **A demonstration of mutations with shared properties.** (A) Hypothetical mutations with their property value vectors are shown. The area in grey shows possible mutation property vectors. (B) The log rates of each new mutation on each clone with a single other mutation are shown. The log rate $f(\mathbf{c}_s, t) = \log(10)$ is circled in red and demonstrates the strongest possible positive causal relationship. This log rate is calculated as $f(\mathbf{c}_s, t) = \log(10)[0, 1]^\top \cdot [0, 0.5]^\top / \|[0, 0.5]\|_2 = \log(10)0.5/0.5 = \log(10)$. The log rate $f(\mathbf{c}_t, s) = \log(10)/2$ is circled in blue and demonstrates how our causal relationships are not symmetric. This log rate is calculated as $f(\mathbf{c}_t, s) = \log(10)[0, 0.5]^\top \cdot [0, 1]^\top / \|[0, 1]\|_2 = \log(10)0.5/1 = \log(10)/2$.

589 Fig. S17C shows CloMu generally reconstructs the mutation properties well, evaluated on all simulation
 590 instances. We also applied CloMu to predict causal relationships. In these simulation instances, causal rela-
 591 tionships can be arbitrarily weak, which increases the difficulty of predicting causal relationships. Fig. S18A
 592 shows an example of true causal relationships for one simulation instance. Fig. S18B shows a precision
 593 recall curve for CloMu and TreeMHN. This demonstrates that CloMu's advantage in simulations with in-
 594 terchangeable mutations also applies to simulations with shared mutation properties in general. In contrast,
 595 Simulations (IV) show that CloMu does not have an advantage over TreeMHN on simulations with 20
 596 mutations but random independent linear causal relationships.

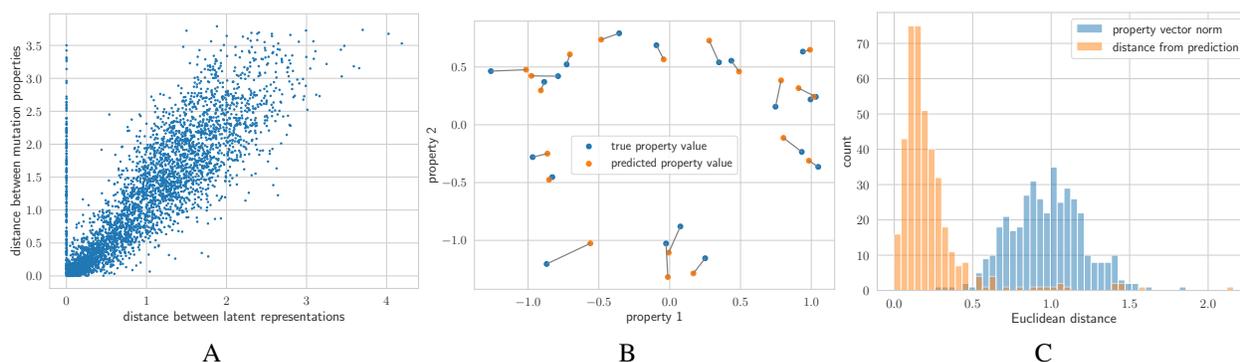


Figure S17: CloMu accurately determines similar mutations and reconstructs mutation properties. (A) Similar mutations (with close property vectors) tend to have close latent representations. A scatter plot of Euclidean distances between all pairs of mutations in all simulations is shown, where the x -axis shows distances between latent representations and the y -axis shows distances between mutation property vectors. The Pearson correlation is 0.8316. (B) Predicted property vectors are reconstructed from the latent representations. A plot of these true property vectors and predicted property vectors is shown for one simulation instance. Grey lines connect the predicted and true property values for each mutation. (C) The Euclidean distance between predicted property vectors and true property vectors is plotted in orange. This is compared against a baseline of predicting the mean normalized property vector of zero for all mutations. The error of this baseline is exactly equal to the norm of the property vector by definition.

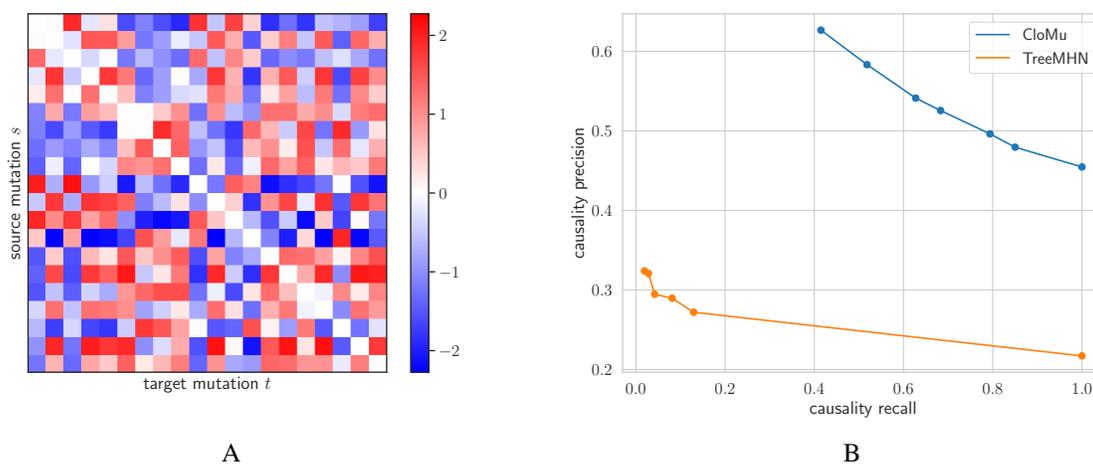


Figure S18: CloMu accurately predicts causal relationships on simulations with shared mutation properties. (A) A plot of causal relationships for one simulation instance. (B) CloMu's advantage over existing methods on simulations with interchangeable mutations also applies to simulations where mutations share properties.

597 **B.3.11 Isolating multi-mutation effects**

598 One advantage of CloMu is that it allows for complex interactions between mutations while utilizing a small
599 number of parameters. The ability of CloMu to detect the effects of combinations of mutations was demon-
600 strated in Simulations (II). However, Simulations (II) also introduced other concepts such as interchange-
601 ability and pathways. We generated an additional set of simulation instances to specifically demonstrate the
602 combination of two mutations affecting the likelihood of other mutations. These simulations contain two
603 driver mutations and eight passenger mutations, totaling $m = 10$ mutations. The combination of the two
604 driver mutations together in a clone affects the rate of all eight passenger mutations. For each of the eight
605 passenger mutations, there is a 0.5 probability of being inhibited and a 0.5 probability of being caused by the
606 combination of the two driver mutations. Inhibited mutations have their rate reduced by a factor of 10, and
607 positively caused mutations have their rate increased by a factor of 10. There is no effect of any individual
608 mutation and the driver mutations only impact the clone when they occur together. The baseline rate of the
609 driver mutations is 5 times the rate of the passenger mutations. Each simulation instance contains $n = 500$
610 patients, and 20 simulation instances were generated.

611 Given the definitions of multi-mutation causality outlined in Supplemental Material [A.2.2](#), we predict
612 causal relationships from pairs of mutations to third mutations. The precision-recall curve is shown in
613 Fig. [S19A](#), which uses the definition of precision and recall for signed causal relationships stated in Sup-
614 plemental Material [B.3](#). The average F1 score across simulation instances reaches 0.9286 for the optimal
615 threshold value of τ .

616 Without additional analysis, it is theoretically possible that CloMu is only picking up on causal relation-
617 ships from pairs in a linear manner by attributing half of the effect to each individual mutation. For instance,
618 if s and t cause r it is possible CloMu is falsely concluding that s and t both individually cause r each with
619 half the total effect size. To disprove this, we compare the predicted causal effect of a pair of mutations with
620 the sum of the individual effects of the two mutations in the pair. Specifically, we define the *linear predicted*
621 *causality* from s and t to r as $f_{\theta}(\mathbf{c}_s, r) + f_{\theta}(\mathbf{c}_t, r) - f_{\theta}(\mathbf{c}_0, r)$. In Fig. [S19](#) we compare this with the predicted
622 multi-mutation causality $f_{\theta}(\mathbf{c}, r)$ where \mathbf{c} contains exactly mutations s and t . Specifically, Fig. [S19B](#) shows
623 this for positive causal relationships and Fig. [S19C](#) shows this for negative inhibitory relationships. This
624 demonstrates that CloMu is successfully determining non-linear effects.

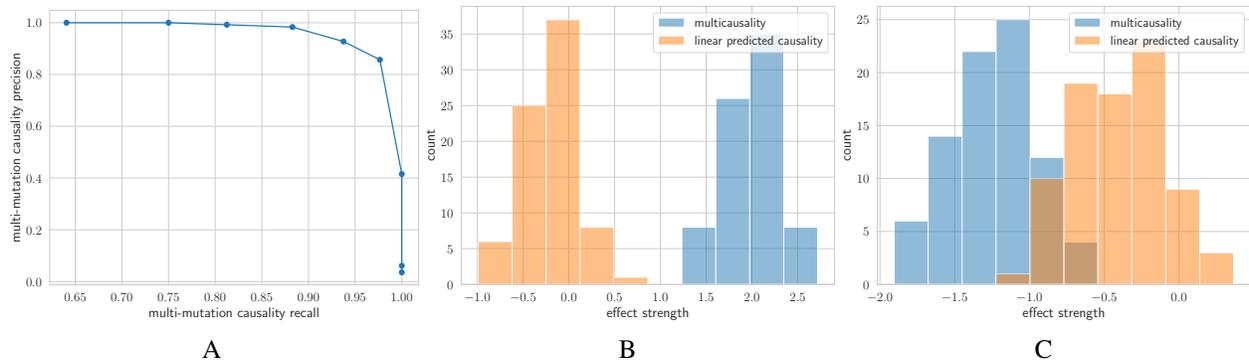


Figure S19: **CloMu detects non-linear multi-mutation causality caused by pairs of mutations** (A) A precision recall curve of CloMu detecting multicausality. The average F1 score across simulation instances reaches 0.9286 for the optimal threshold value of τ . (B) For positive relationships, CloMu predicts multi-mutation causality from a pair of mutations beyond the causality from each individual mutation in the pair. (C) For negative relationships, CloMu predicts negative multi-mutation causality from a pair of mutations beyond the causality from each individual mutation in the pair.

625 **B.3.12 Results on RECAP simulations – Simulations (III)**

626 We compared CloMu to RECAP as well as REVOLVER. As in the Main Text, we consider the tree selection
 627 accuracy, which is the fraction of patients for which the correct true tree is predicted. Fig. S20 shows that
 628 CloMu performed almost identically to RECAP and far outperforms REVOLVER. In fact, our method and
 629 RECAP both achieved a 98.6% accuracy on average for simulations with 5 mutations and an average 100%
 630 accuracy for 12 mutations (Fig. S20A,C). We additionally investigated the predicted number of patient
 631 clusters for our method and RECAP, finding perfect performance for CloMu and near perfect performance
 632 for RECAP across varying numbers of mutations (Fig. S20B,D).

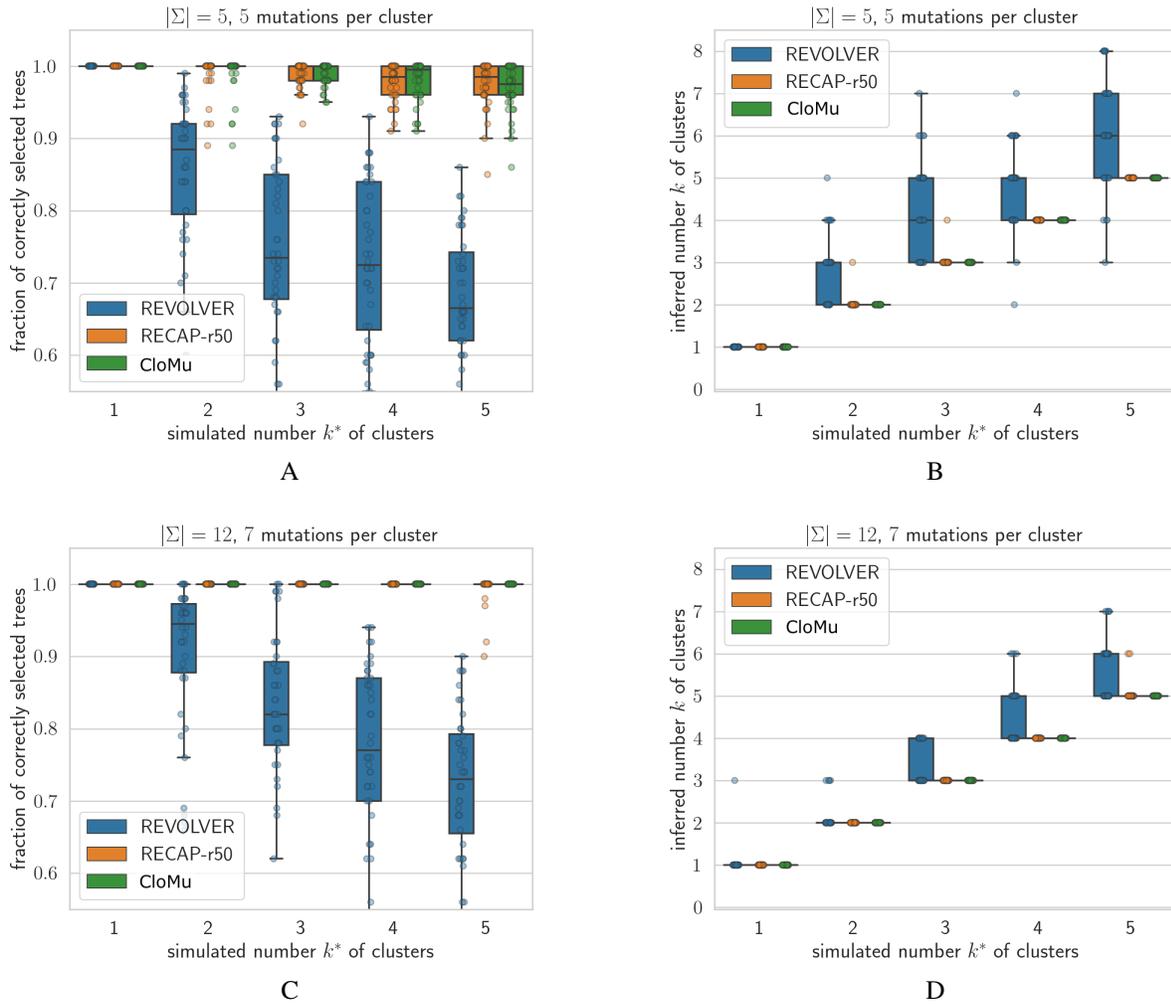


Figure S20: **CloMu’s matches RECAP’s performance on RECAP’s simulated data (Simulations (III)).** These plots show RECAP results with 50 restarts and additionally include REVOLVER results reported in [Christensen et al. \(2020\)](#). (A,C) Tree selection accuracy. (B,D) The number of ground-truth patient clusters vs. the predicted number of patient clusters. Panels (A-B) show results for $|\Sigma| = 5$ total mutations and $m = 5$ mutations per patient whereas panels (C-D) show results for $|\Sigma| = 12$ total mutations and $m = 5$ mutations per patient.

633 **B.3.13 Results on TreeMHN simulations – Simulations (IV)**

634 CloMu is generally competitive with but slightly less effective than TreeMHN on datasets from the
635 TreeMHN paper. Datasets in the TreeMHN paper use completely independent causal relationships be-
636 tween all mutations, with no related or interchangeable mutations. In this case, using CloMu with a neural
637 network function creates a slight disadvantage on datasets with a very poor signal-to-noise ratio. Therefore,
638 we included a version of CloMu where the neural network function is replaced with a simple linear model,
639 denoted as “Linear CloMu”. One virtue of the training setup of CloMu is that it is completely independent
640 of the function chosen for f_θ . Therefore, this only requires a minor tweak to the code. We compared CloMu
641 to TreeMHN with stability selection in the cases of $n = 300$ patients, and $m \in \{10, 15, 20\}$ mutations.
642 The TreeMHN setup and hyperparameters were set exactly as in the TreeMHN paper. Like in the TreeMHN
643 paper, we plot the precision and recall curves. Unlike the TreeMHN paper, we report the precision and recall
644 on all mutations, not only the top half most frequent mutations.

645 Fig. S21 shows the precision and recall curves averaged over the 20 simulations per simulation setup.
646 One can see that TreeMHN outperformed CloMu on these data, but CloMu remained competitive especially
647 when a linear function is used for f_θ . For instance, it seems CloMu with a linear model may extend to higher
648 recall values more effectively than TreeMHN, but this is unclear since the range of hyperparameters used
649 in the TreeMHN paper does not extend into high recall. From these results, we conclude the following. If
650 one knows for certain that there are no shared properties between mutations and that all causal relationships
651 are independent, then TreeMHN is a slightly more effective method than CloMu. However, if there exists
652 shared properties between mutations, CloMu ranges from slightly more effective to vastly more effective
653 than TreeMHN.

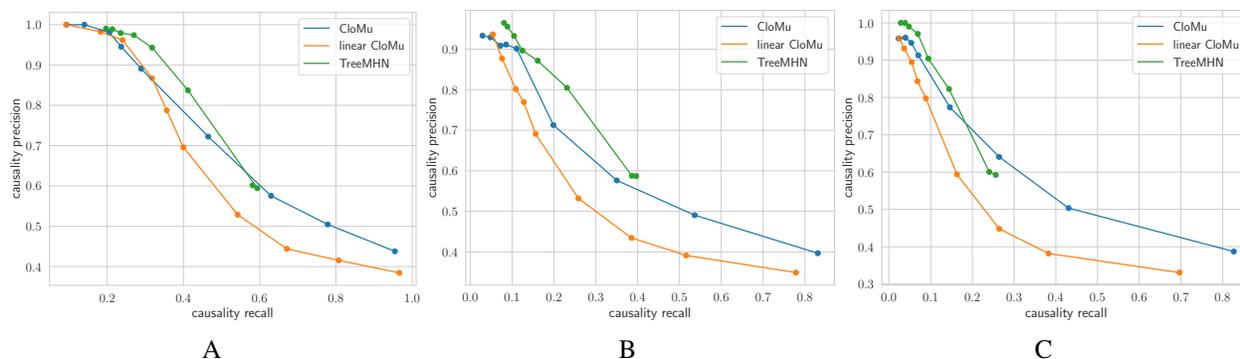


Figure S21: **CloMu performs slightly worse than TreeMHN on TreeMHN’s simulated data (Simulations (IV)).** (A) Simulations with $m = 10$ total mutations. (B) Simulations with $m = 15$ mutations. (C) Simulations with $m = 20$ total mutations.

654 C Supplemental results for real data

655 In Supplemental Material C.1 we describe additional analyses on the breast cancer dataset (Razavi et al.,
 656 2018). In Supplemental Material C.2 we describe additional analyses on the AML dataset (Morita et al.,
 657 2020).

658 C.1 Breast cancer

659 CloMu was ran on the breast cancer dataset (Razavi et al., 2018) using a laptop with a 2.4 GHz CPU and
 660 64 GB of RAM without using a GPU. The runtime was 4 hours and 23 minutes, and the memory usage was
 661 10.777 GB. The method of bootstrapping to demonstrate statistical bounds on predictions is demonstrated
 662 in this dataset in Supplemental Material C.1.1. The connection between latent representations and breast
 663 cancer hormone receptor status is shown in Supplemental Material C.1.2. Relative causality predictions
 664 are systematically compared to TreeMHN’s predictions in Supplemental Material C.1.3. CloMu’s ability to
 665 predict the next mutation to occur on subtrees is demonstrated in Supplemental Material C.1.4.

666 C.1.1 Bootstrapping

667 Supplemental Material A.3 introduced bootstrapping for CloMu. We applied this procedure to the breast
 668 cancer dataset (Razavi et al., 2018), using 100 bootstrapping replicates. Fig. S22 shows box plots of relative
 669 causality predictions on the breast cancer data, showing that our predictions fall within the statistical bounds
 670 established by bootstrapping. Fig S23 shows box plots of the rankings of these relative causality predictions,
 671 again demonstrating the robustness of our predictions. Additionally, a 90% confidence interval of predicted

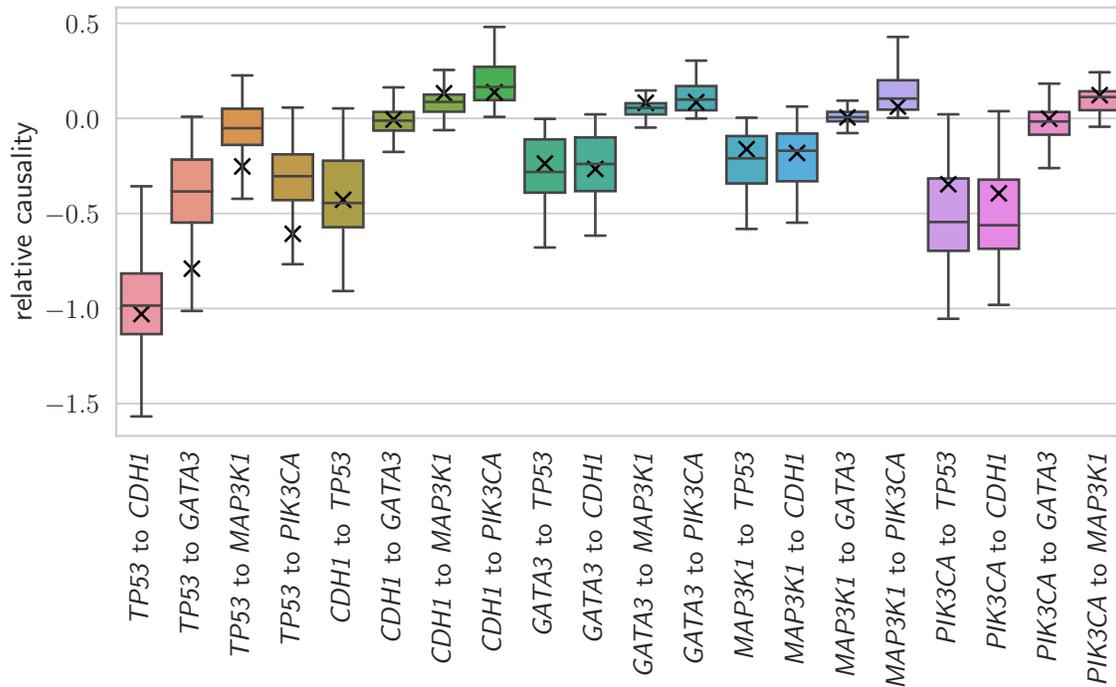


Figure S22: **Bootstrapping provides statistical bounds for relative causality predictions on the breast cancer dataset (Razavi et al., 2018)** . A box plot of bootstrapping predicted relative causality is shown. The value predicted on the original dataset is shown with a black 'X'. This demonstrates statistical bounds on our causal relationship predictions.

672 fitness values on bootstraps is shown in Fig. S24. The set of high-fitness mutations is very stable across
 673 bootstrapping instances.

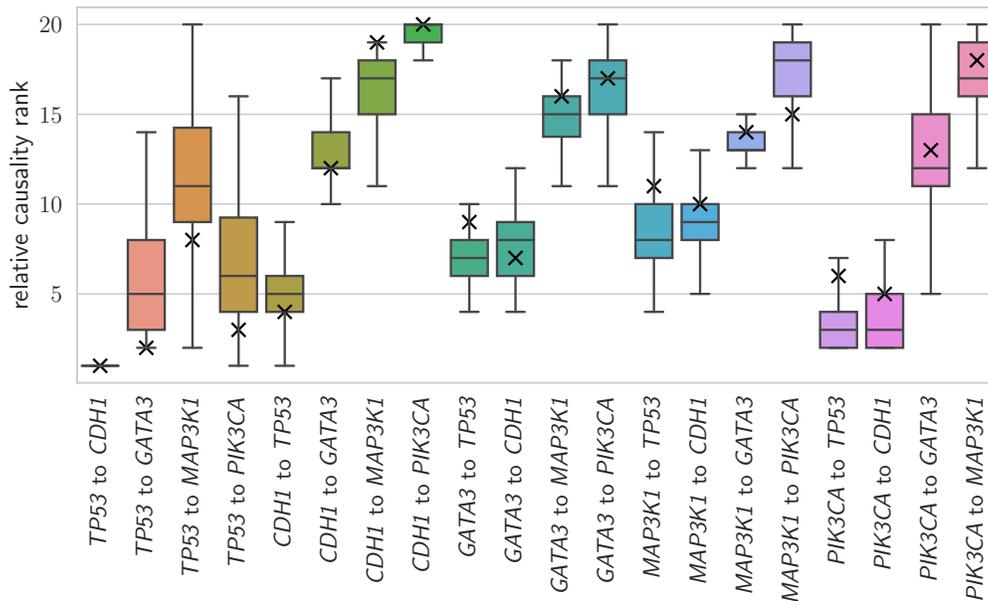


Figure S23: **Bootstrapping provides statistical bounds for the rankings of relative causality predictions on the breast cancer dataset (Razavi et al., 2018)**. The 20 causal relationships being analyzed are ranked from 1 to 20, where 20 indicates the most positive relationship and 1 indicates the most negative relationship. A box plot of bootstrapping predicted relative causality rankings is shown. The ranking predicted on the original dataset is shown with a black 'X'. The strongest positive and negative causal relationships are fairly consistent.

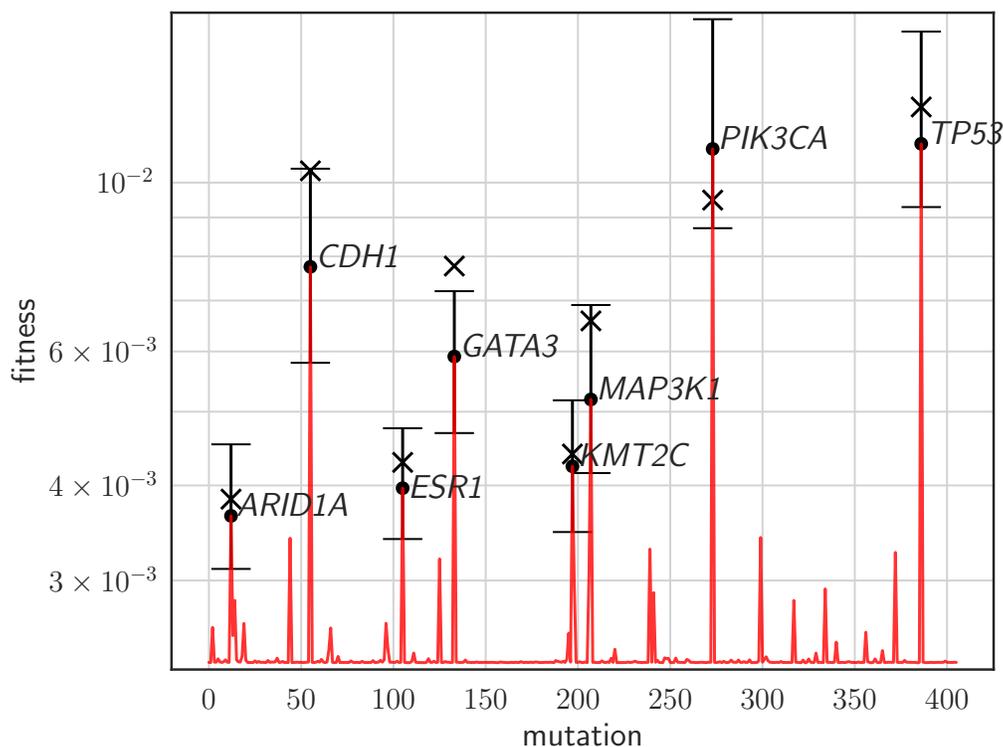


Figure S24: **Bootstrapping gives statistical bounds for fitness predictions on the breast cancer data set (Razavi et al., 2018).** The median fitness values for all mutations are shown in red, with the 90% confidence interval on 100 bootstrapped intervals shown in black. Specifically, this shows the 6th largest and 6th smallest fitness values, therefore excluding the 5 highest and lowest values. The fitness values predicted on the original dataset are shown with a black 'X'.

674 C.1.2 Breast cancer hormone receptor status and latent representations

675 The breast cancer dataset contains patients with different hormone receptor statuses, which include
676 HR+/HER2+, HR+/HER2- and triple negative (Razavi et al., 2018). Specifically, there are 179 patients
677 with HR+ and HER2+, 1501 patients with HR+ and HER2-, 62 patients with HR- and HER2+, and 176
678 patients that are triple negative. Within the subset of patients our model was run on (patients with at most
679 9 mutations), there are 113 patients with HR+ and HER2+, 943 patients with HR+ and HER2-, 39 patients
680 with HR- and HER2+, and 129 patients that are triple negative. We investigated if different latent repre-
681 sentations are associated with different hormone receptor statuses. To do so, we must first define a latent
682 representation of a patient. The latent representation of a clone is directly produced by our neural network
683 by inputting the vector representation of the clone. We define the latent representation of a tree as the mean
684 representation of all clones in the tree. We define the latent representation of a patient as the mean latent
685 representation of all possible trees for that patient (since the true tree is unknown). We found no connection
686 between the latent representation of patients and the HER2+ hormone status, however, we did observe an
687 association with the HR+ hormone status. Fig. S25 shows plots of the first three PCA components of the
688 latent representations of patients with the HR+ and HR- hormone statuses,. We see that there tends to be
689 separation between patients with the HR+ hormone status and those with the HR- hormone status, and that
690 this is mainly driven by the first two components.

691 The main text contains plots showing the full latent representations of all mutations in the breast can-
692 cer dataset (Razavi et al., 2018). These plots have the advantage of showing all relevant information but
693 the disadvantage of potentially being less visually intuitive. PCA allows for the dimension of the latent
694 representation to be reduced, allowing for a scatter plot visualization. Additionally, we compare the latent
695 representation with mutation fitness values. Fig. S26 shows these PCA values. The correlation between
696 PCA component 1 and fitness is 0.8947 with a p-value less than 10^{-142} . The correlation between PCA com-
697 ponent 2 and fitness is 0.4323 with a p-value less than 10^{-19} . There is no significant correlation between
698 PCA component 3 and fitness (p-value of 0.7073). We find that, especially in the first two components, low
699 fitness mutations cluster together, whereas high fitness mutations are at more distinct positions. Another
700 interesting analysis of latent representations is the Euclidean distance between the latent representations of
701 different mutations. In Fig. S27, we see there are large numbers of pairs with a near zero Euclidean distance
702 apart, indicating interchangeable mutations. Many such mutations are low fitness mutations with a latent

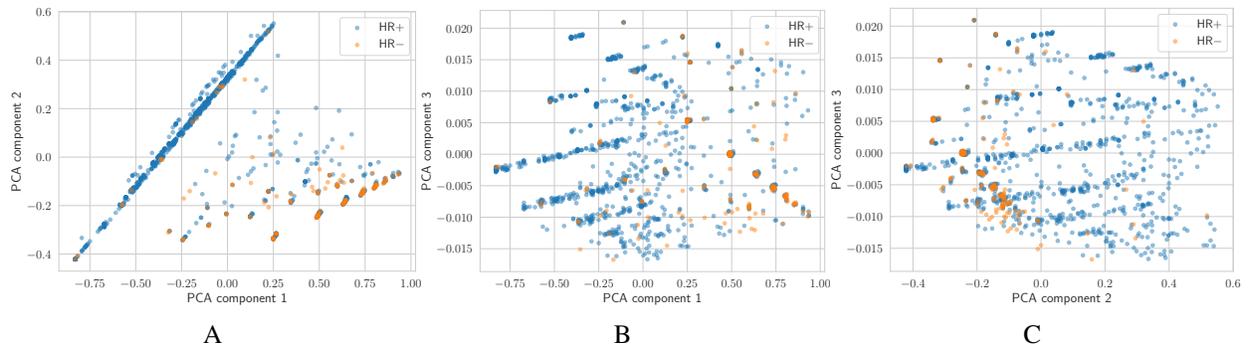


Figure S25: CloMu produces latent representations somewhat associated with hormone receptor status. Scatter plots are shown of latent representations of breast cancer patients with hormone receptor status HR+ an HR-. There appears to be an association between latent representation and hormone receptor status. (A) A plot of PCA components 1 and 2. (B) A plot of PCA components 1 and 3. (C) A plot of PCA components 2 and 3.

703 representation near zero.

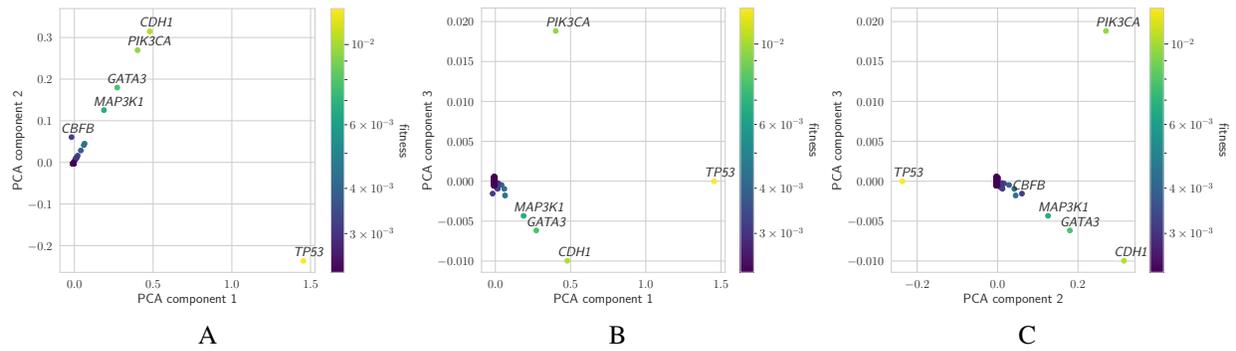


Figure S26: **PCA of latent representations on the breast cancer dataset (Razavi et al., 2018)**. Mutations are colored based on their fitness. All mutations are shown in the scatter plot, however, low-fitness mutations tend to cluster near $(0, 0)$. High-fitness mutations are indicated with a distinct label. The correlation between PCA component 1 and fitness is 0.8947 with a p-value less than 10^{-142} . The correlation between PCA component 2 and fitness is 0.4323 with a p-value less than 10^{-19} . There is no significant correlation between PCA component 3 and fitness. (A) PCA components 1 and 2 are shown. (B) PCA components 1 and 3 are shown. (C) PCA components 2 and 3 are shown.

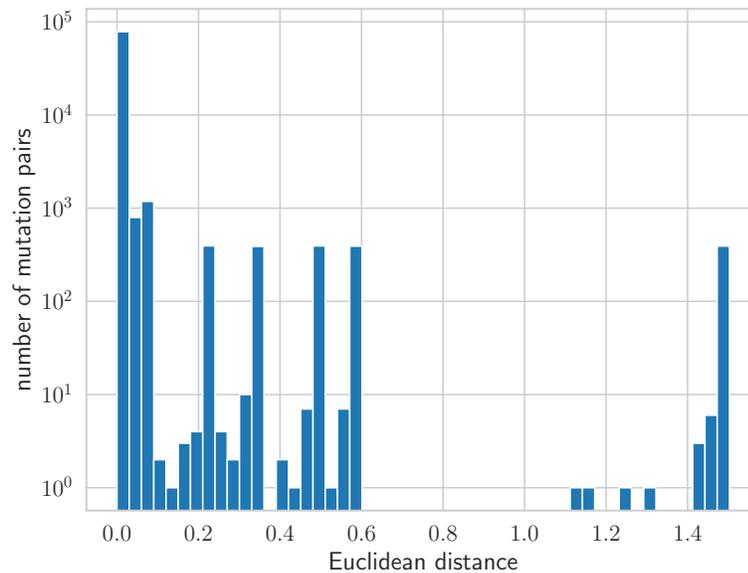


Figure S27: **Many mutations have highly similar latent representations in the breast cancer data.** A histogram of Euclidean distances between pairs of mutation latent representations on the breast cancer data set (Razavi et al., 2018) is shown. There is a clear spike around zero for interchangeable mutations (including low-fitness mutations).

704 C.1.3 Systematic comparison with TreeMHN

705 Unlike TreeMHN’s causality measurement, our definition of relative causality explicitly removes increased
706 mutation likelihoods due to high fitness. Consequently, CloMu should predict lower relative causality values
707 from fit mutations when compared to TreeMHN’s causality predictions. If TreeMHN predicts a negative
708 causal relationship then CloMu should also theoretically predict a negative causal relationship but not vice
709 versa. Similarly, if CloMu predicts a positive causal relationship then TreeMHN should also theoretically
710 predict a positive causal relationship but not vice versa. To systemically compare CloMu and TreeMHN’s
711 causality predictions, we restricted our analysis to causal relationships between pairs of mutations with
712 meaningfully nonzero latent representations when predicting causal relationships as previously described.
713 On this restricted set of mutation pairs, we analyzed all ordered pairs of mutations in which either TreeMHN
714 predicted a negative causal relationship or CloMu predicted a positive causal relationship. We found strong
715 agreement between the causal relationship predictions of CloMu and TreeMHN on these mutations pairs.

716 All of TreeMHN’s negative predictions on the breast cancer cohort were also predicted negative by
717 CloMu (Table S3, four cases). Additionally, all of CloMu’s positive predictions on the breast cancer co-
718 hort were also predicted positive by TreeMHN. However, as expected, there exists disagreements in which
719 CloMu predicts a negative relationship and TreeMHN predicts a positive relationship. Specifically there are
720 four such cases, all of which are from the extremely fit mutations *TP53* and *PIK3CA*.

Source mutation	Target mutation	TreeMHN prediction	CloMu prediction
<i>GATA3</i>	<i>TP53</i>	Negative	Negative
<i>GATA3</i>	<i>CDH1</i>	Negative	Negative
<i>CDH1</i>	<i>PIK3CA</i>	Positive	Positive
<i>GATA3</i>	<i>PIK3CA</i>	Positive	Positive
<i>MAP3K1</i>	<i>PIK3CA</i>	Positive	Positive
<i>CDH1</i>	<i>MAP3K1</i>	Positive	Positive
<i>GATA3</i>	<i>MAP3K1</i>	Positive	Positive
<i>PIK3CA</i>	<i>MAP3K1</i>	Positive	Positive

Table S3: On the breast cancer cohort, all negative causal relationships predicted by TreeMHN are also predicted by CloMu. Additionally, all positive causal relationships predicted by CloMu are also predicted by TreeMHN.

721 C.1.4 Predicting subsequent mutations on subtrees

722 Another way of evaluating performance on real data is observing subtrees of patients' trees and predicting
723 which mutations will subsequently occur. Specifically, if there exists some tree T and some rooted subtree
724 T' of T , methods could predict a new mutation s to occur on some clone in T' . To that end, we define
725 $\text{Extend}(T', T)$ such that $(\mathbf{c}, s) \in \text{Extend}(T', T)$ if and only if \mathbf{c} is a clone in T' , and T includes mutation
726 s being added to \mathbf{c} to produce a new clone in T . More informally, $\text{Extend}(T', T)$ is the set of clones in T'
727 in which additional mutations occur as part of T (together with those mutations).

728 Using CloMu's tree generative process, the probability that s is the next mutation to occur on T' is
729 simply $\Pr((\mathbf{c}_i, s) \mid V(T'), f_\theta) = \Pr((\mathbf{c}_i, s) \mid \mathbf{c}_1, \dots, \mathbf{c}_k, f_\theta)$ where $V(T') = \{\mathbf{c}_1, \dots, \mathbf{c}_k\}$ is the set of clones
730 in T' . An example of this calculation is shown in Fig. S28. As a minor caveat, on the breast cancer dataset
731 (Razavi et al., 2018), the infinite sites assumption is made and so the probability of mutations that violate
732 this assumption is set to zero (and the probability is normalized to sum to 1).

733 We define $\text{Next}_C(T', T)$ as the sum over all $(\mathbf{c}, s) \in \text{Extend}(T', T)$ of these predicted probabilities (of
734 s occurring on \mathbf{c} being the next mutation) according to CloMu. We define $\text{Next}_M(T', T)$ as the analogous
735 computation being performed with TreeMHN's probabilistic model. Finally, we define $\text{Next}_B(T', T)$ as
736 the analogous computation according to a baseline model in which mutations occur at a rate proportional
737 to the percentage of patients in which they occur (independent of the clone they occur on). Specifically,
738 $\text{Next}_B(T', T) = 1/|V(T')| \sum_{(\mathbf{c}, s) \in \text{Extend}(T', T)} \text{freq}(s)$, where $\text{freq}(s)$ is the proportion of patients with
739 mutation s and $|V(T')|$ is the number of clones in T' .

740 Our procedure consists of picking random subtrees of random patients and then evaluating these three
741 models in terms $\text{Next}_C(T', T)$, $\text{Next}_M(T', T)$, and $\text{Next}_B(T', T)$. There are many ways of randomly

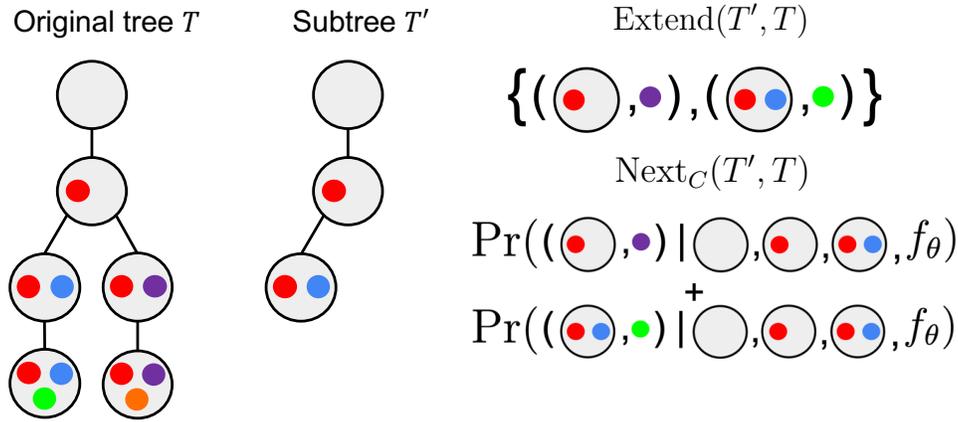
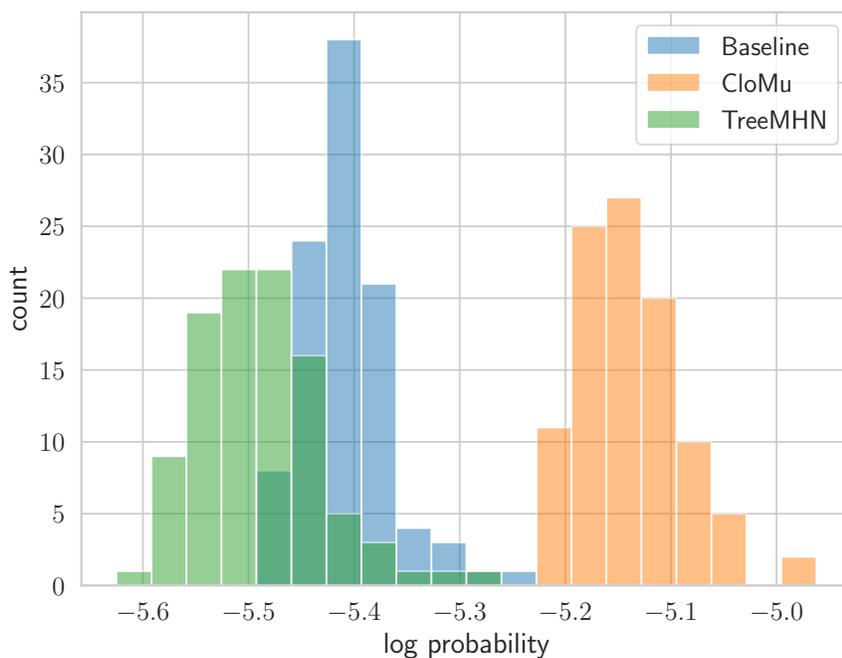


Figure S28: **A diagram of next mutation prediction on subtrees.** An original tree T as well as a subtree T' are shown. The set of extensions $\text{Extend}(T', T)$ is shown. Additionally, the calculation of the probability $\text{Next}_C(T', T)$ is shown as the sum of the probability of extensions in $\text{Extend}(T', T)$ according to CloMu's probability estimates.

742 generating subtrees from a data set, so we must precisely define our random sampling procedure. First, we
743 uniformly randomly sampled a patient and uniformly randomly sampled one of the possible trees from that
744 patient. Then, we uniformly randomly sampled the size of the subtree in terms of vertices in the range of 1 to
745 the number of vertices in the full tree minus 1. Then, we generated a random subtree by starting with the full
746 tree and uniformly randomly removed leaves until the desired tree size is reached. Theoretically, we wish to
747 maximize the product of probabilities (such as $\text{Next}_C(T', T)$) across samples for a large number of sampled.
748 This is equivalent to maximizing the average value of $\log(\text{Next}_C(T', T))$ across samples. To avoid extreme
749 outliers greatly biasing the data, we instead evaluate an average of $\log(\text{Next}_C(T', T) + 0.001)$ across the
750 data. Specifically, we generate 100 batches of 1000 samples and individual average each of $\text{Next}_C(T', T)$,
751 $\text{Next}_M(T', T)$, and $\text{Next}_B(T', T)$ within each batch. This averaging allows for a plot that is not dominated
752 by noise (due to tumor evolution being a highly random probabilistic process). Histograms of these results
753 on the breast cancer dataset (Razavi et al., 2018) are given in Fig. S29. CloMu greatly exceeds the frequency
754 based baseline and TreeMHN.



A

Figure S29: **CloMu accurately predicts future mutations given a partial tree breast cancer data (Razavi et al., 2018)**. CloMu substantially outperforms TreeMHN and the mutation frequency based baseline on this dataset.

755 **C.2 AML**

756 CloMu was ran on the AML dataset (Morita et al., 2020) using a laptop with a 2.4 GHz CPU and 64 GB
 757 of RAM without using a GPU. The runtime was 41 minutes, and the memory usage was 261 MB. Further
 758 analysis of high fitness mutations and clonal prevalence data is performed in Supplemental Material C.2.1.
 759 An additional analysis of the latent representations on this data set is shown in Supplemental Material C.2.2.
 760 The relative causality predictions of CloMu are compared with TreeMHN’s predictions in Supplemental
 761 Material C.2.3. CloMu’s ability to predict subsequent mutations on subtrees of phylogeny trees is tested in
 762 Supplemental Material C.2.4.

763 **C.2.1 Additional analysis of high fitness mutations and prevalence**

764 We have orthogonally validated our fitness predictions on AML data through the use of clonal prevalence
 765 data. One relevant question is whether higher fitness mutations tend to occur earlier or later in clonal evolu-

766 tion. To do so, we measured the average number of mutations that have already occurred on a clone when
767 some new mutation s is introduced. We averaged this over all patients for which that mutation is intro-
768 duced. A histogram of these data for high-fitness mutations and all other mutations is shown in Fig. S30A.
769 Fig. S30B shows the same for the breast cancer dataset (Razavi et al., 2018). This shows that high-fitness
770 mutations generally but not always occur earlier in cancer evolution. One possible limitation in our fitness
771 predictions is the fact that terminal mutations at the time of sampling do not have observed mutations oc-
772 ccurring after them and thus do not contain any evidence for their fitness. Consequently, if a mutation always
773 occurs terminally our model could not conclude that the mutation is highly fit.

774 Fig. S30C shows a scatter plot of the fitness of a mutation and what proportion of patients have that
775 mutation occur terminally. Generally, higher fitness mutations do not occur terminally the majority of the
776 time. However, we were not able to detect a statistically significant (p-value under 0.05) Pearson correlation
777 between proportion terminal and mutation fitness given the relatively small number of high fitness mutations.
778 For instance, Fig. S30C shows that two mutations (*NPM1* and *DNMT3A*) are vastly more fit than other
779 mutations in the AML data set. Additionally, among mutations that are at most 50% terminal, correlations
780 between proportion terminal and mutation fitness have a p-value above 0.75 on both datasets. It is unknown
781 whether future data sets could uncover a statistically significant connection between fitness and percentage
782 terminal.

783 Another relevant question is whether earlier mutations have higher log prevalence ratios due to the root
784 clone (of normal cells) having a lower prevalence than the initial clones with one additional mutation added
785 to the root. If one assumes a constant percentage rate of clonal growth, the log prevalence ratio of mutations
786 would not be biased by how early or late they occur. Instead, the log prevalence ratio would only be affected
787 by how much earlier the parent clone formed relative to the child clone (in addition to clonal growth rates).
788 However, tumor growth may be more complex than this simple exponential model. Fig. S31A shows a
789 histogram of the prevalence of the root node and the prevalence of clones with only a single mutation.
790 We measured the log prevalence ratio of the healthy root clone and its children clones averaged across all
791 patients. We find this mean log prevalence ratio to be 0.0569 ± 0.0738 . This is sufficiently close to zero that
792 our analysis is not biased by some mutations occurring earlier and thus being more likely to have occurred on
793 the root clone. An additional question is whether terminal mutations tend to have a negative log prevalence
794 ratio in addition to their low fitness. Fig. S31B shows a histogram of the prevalence of a terminal clone
795 and its parent clone. Additionally, Fig. S31C shows a histogram of the log prevalence ratios of terminal

796 mutations. We find this mean log prevalence ratio to be -0.3799 ± 0.0903 . This is meaningfully negative,
797 however, this effect is of a smaller magnitude than the difference in log prevalence ratios between high-
798 fitness mutations (mean: 0.8316) and low-fitness mutations (mean: -0.3703), shown in Main Text Figure
799 4B. To conclude, although our method may have diminished sensitivity in detecting high fitness mutations
800 that typically occur terminally, we believe that specificity is not affected by mutation timing.

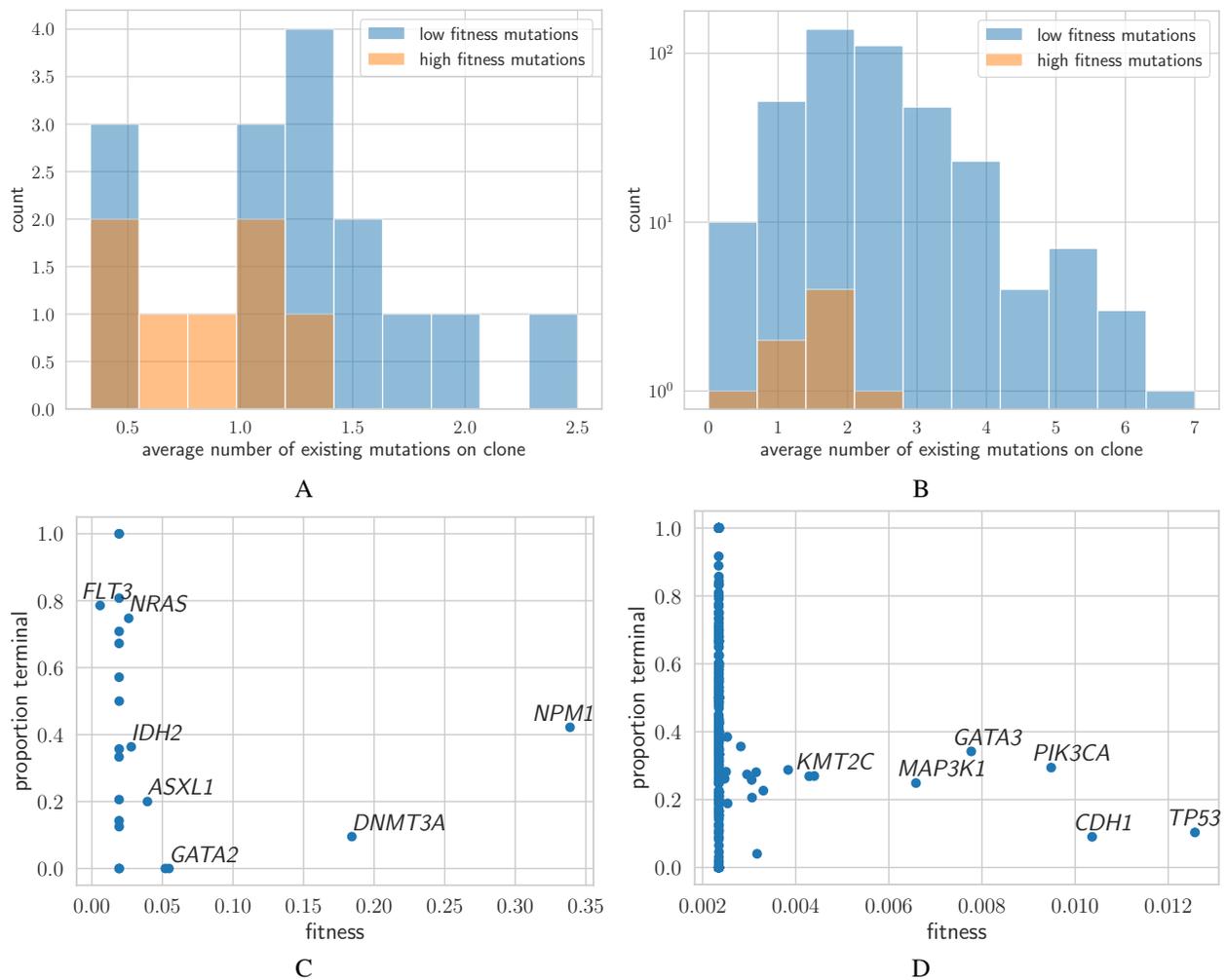


Figure S30: High fitness mutations tend to occur earlier in evolution. (A) On the AML dataset (Morita et al., 2020), higher fitness mutations often but not always occur early in tumor evolution. One exception is the mutation *NRAS* for which the average existing number of mutations on the clone is 1.384 which is above the median value of 1.147 across mutations. (B) On the breast cancer dataset (Razavi et al., 2018), higher fitness mutations often but not always occur early in tumor evolution. One exception is the mutation *ARID1A* with an average existing number of mutations on the clone equal to 2.123, which is above the median value of 2.054 across mutations. (C) On the AML dataset (Morita et al., 2020), mutations that are terminal the majority of the time tend not to be high fitness. However, among the mutations that are frequently non-terminal, there is no clear connection between fitness and the proportion for which it is terminal. (D) On the breast cancer dataset (Razavi et al., 2018), mutations that are terminal the majority of the time are not high fitness. However, among the mutations that are frequently non-terminal, the higher fitness mutations do not have the lowest proportion occurring terminally.

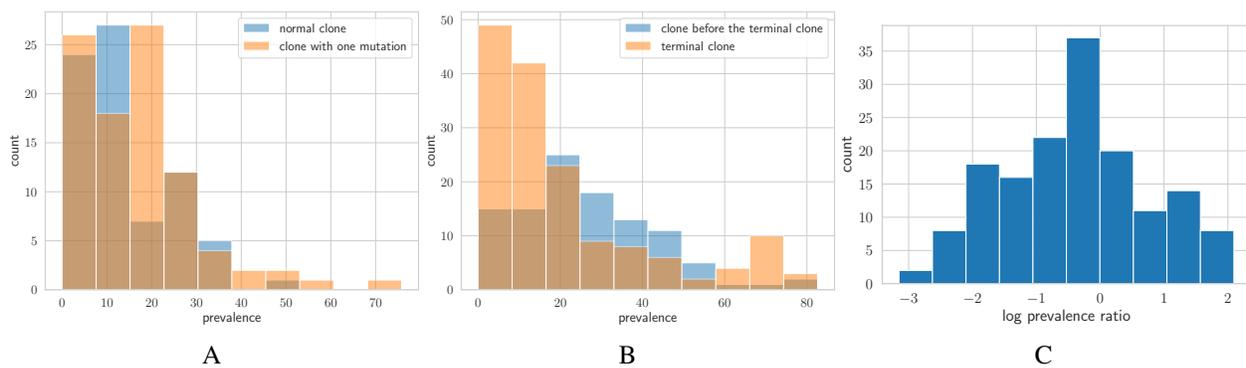


Figure S31: **Analyzing the prevalence of initial clones and terminal clones.** (A) The prevalence of normal clones and clones with one mutation are found to be very similar. (B) The prevalence of terminal clones is found to be smaller on average than clones immediately before the terminal clone. (C) The log prevalence ratios of the terminal mutations are slightly negative on average.

801 C.2.2 Additional analysis of latent representations

802 The main text contains plots showing the full latent representations of all mutations on the AML dataset
803 ([Morita et al., 2020](#)). However, as mentioned in Supplemental Material [C.1.2](#) applying PCA can provide a
804 more intuitive visualization of the results. Fig. [S32](#) shows these PCA values with mutations colored by their
805 fitness value. There is a strong correlation between PCA component 1 and fitness. The correlation between
806 PCA component 1 and fitness is 0.9709 with a p-value under 10^{-13} . There is no statistically significant
807 correlation between either PCA component 2 or component 3 with fitness (p-values of 0.7138 and 0.6175
808 respectively). Another interesting analysis of latent representations is the Euclidean distance between the
809 latent representations of different mutations. Fig. [S33](#) shows a histogram of Euclidean distances between
810 the latent representations of all pairs of mutations on the AML dataset ([Morita et al., 2020](#)). Many pairs of
811 mutations are a very small Euclidean distance apart.

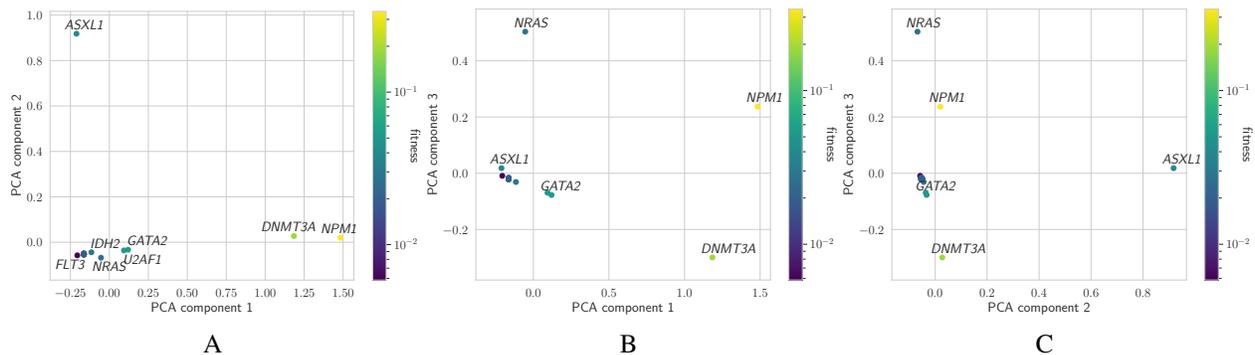


Figure S32: **PCA of latent representations on the AML dataset (Morita et al., 2020)**. All mutations are shown in the scatter plot, however, low-fitness mutations tend to cluster near the bottom left. Mutations are colored based on their fitness value. The correlation between PCA component 1 and fitness is 0.9709 with a p-value under 10^{-13} . There is no statistically significant correlation between either of the other PCA components and fitness. (A) PCA components 1 and 2 are plotted. (B) PCA components 1 and 3 are plotted. (C) PCA components 2 and 3 are plotted.

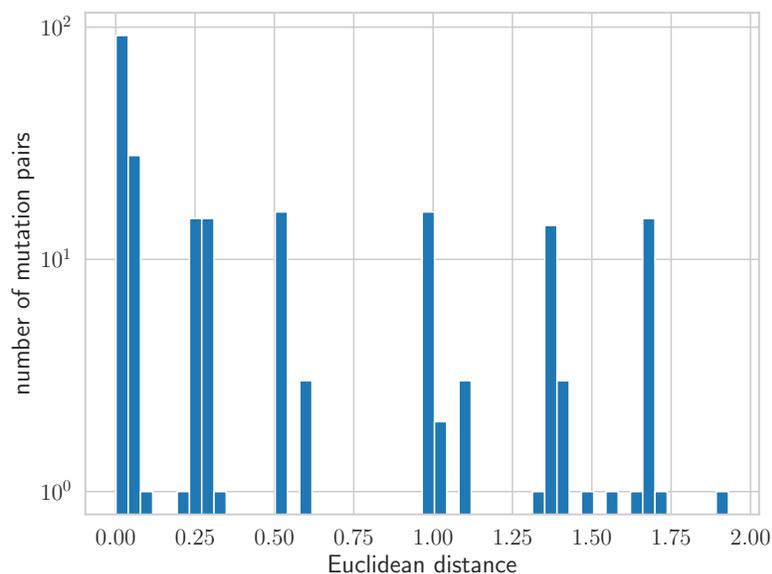


Figure S33: **Many mutations have highly similar latent representations in the AML dataset (Morita et al., 2020)**. A histogram of Euclidean distances between pairs of mutation latent representations is shown. The large spike around 0 demonstrates interchangeable mutations (including low-fitness mutations).

812 C.2.3 Systematic comparison with TreeMHN

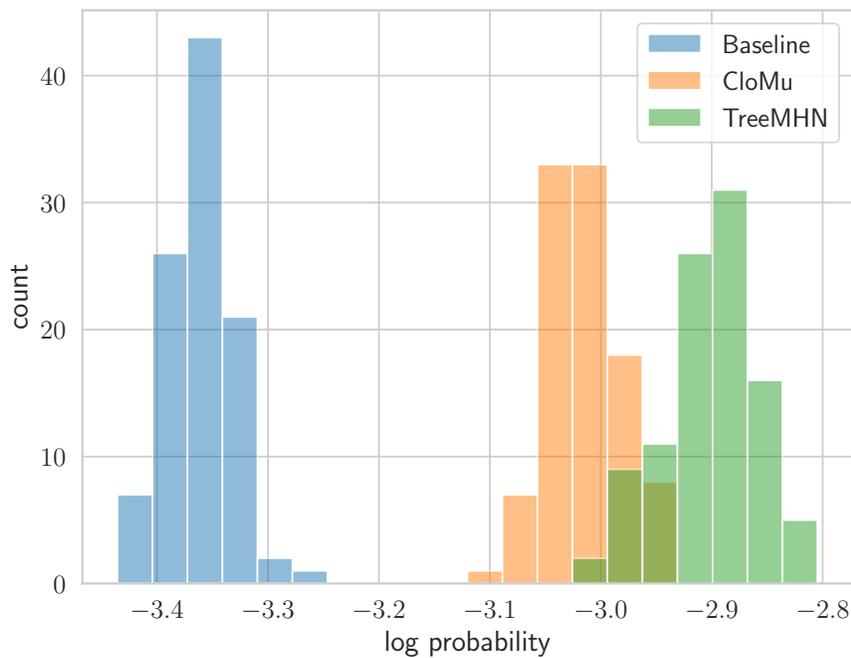
813 As mentioned in Supplemental Material C.1.3, CloMu is expected to predict lower relative causality values
814 from fit mutations when compared to TreeMHN's causality predictions (due to CloMu accounting for fitness
815 differently). Thus, if TreeMHN predicts a negative causal relationship then CloMu should as well and if
816 CloMu predicts a positive causal relationship then TreeMHN should as well. As in Supplemental Material
817 C.1.3, we restricted to pairs of mutations with meaningfully nonzero latent representations when predicting
818 causal relationships. Additionally, we did not consider causal relationships from *FLT3* since *FLT3* is the
819 only low-fitness mutation with a meaningfully nonzero latent representation. However, we still do consider
820 causal relationships where *FLT3* is the target mutation. On this restricted set of mutation pairs, we analyzed
821 all ordered pairs of mutations in which either TreeMHN predicted a negative causal relationship or CloMu
822 predicted a positive causal relationship. We specifically analyzed the 4 strongest positive causal relationships
823 according to CloMu, since all other positive causal relationships were extremely weak. As in Supplemental
824 Material C.1.3, there was a general agreement between the causal relationship predictions of CloMu and
825 TreeMHN. All of TreeMHN's negative predictions were also predicted negative by CloMu. Additionally,
826 all of CloMu's positive predictions were either predicted positive or neutral by TreeMHN (Table S4, 2
827 positive cases, 2 neutral cases). However, there exist 2 cases in which TreeMHN predicts a positive causal
828 relationship and CloMu predicts a negative causal relationship.

Source mutation	Target mutation	TreeMHN prediction	CloMu prediction
<i>ASXL1</i>	<i>NPM1</i>	Negative	Negative
<i>ASXL1</i>	<i>DNMT3A</i>	Negative	Negative
<i>NPM1</i>	<i>ASXL1</i>	Negative	Negative
<i>NPM1</i>	<i>DNMT3A</i>	Negative	Negative
<i>ASXL1</i>	<i>NRAS</i>	Neutral	Positive
<i>NPM1</i>	<i>NRAS</i>	Positive	Positive
<i>NPM1</i>	<i>FLT3</i>	Positive	Positive
<i>NRAS</i>	<i>FLT3</i>	Neutral	Positive

Table S4: On the AML cohort, all negative causal relationships predicted by TreeMHN are also predicted by CloMu. Additionally, all positive causal relationships predicted by CloMu are either positive or neutral according to TreeMHN.

829 C.2.4 Predicting subsequent mutations on subtrees

830 Finally, we evaluate performance by observing subtrees of patients' trees and predicting which mutations
831 will subsequently occur. A description of how we perform this analysis is provided in Supplemental Material
832 C.1.4. A histogram of the results is Fig. S34. Similarly to the breast dataset, CloMu also outperforms the
833 simple mutation frequency based baseline on the AML dataset. On the other hand, TreeMHN performs
834 excellent on the AML dataset (slightly outperforming CloMu), but has worse performance on the breast
835 cancer data set. We believe this is due to advantages of CloMu and the simple frequency based baseline
836 on datasets with vast numbers of passenger mutations. Specifically, the mutation frequency based baseline
837 has the advantage of exactly knowing the frequency of mutations rather than only relying on estimates of
838 mutation frequency stored in the model. Additionally, CloMu has the advantage of modeling how high-
839 fitness mutations increase the rates of all passenger mutations (rather than having to individually model
840 pairwise relationships as in TreeMHN).



A

Figure S34: **CloMu accurately predicts future mutations given a partial tree on the AML dataset (Morita et al., 2020)**. CloMu vastly outperforms the mutation frequency based baseline but slightly underperforms TreeMHN.

841 References

- 842 Caravagna G, Giarratano Y, Ramazzotti D, Tomlinson I, Graham TA, Sanguinetti G, and Sottoriva A. 2018.
 843 Detecting repeated cancer evolution from multi-region tumor sequencing data. *Nature Methods* **15**: 707–
 844 714.
- 845 Christensen S, Kim J, Chia N, Koyejo O, and El-Kebir M. 2020. Detecting evolutionary patterns of cancers
 846 using consensus trees. *Bioinformatics* **36**: i684–i691.
- 847 Jamal-Hanjani M, Wilson GA, McGranahan N, Birkbak NJ, Watkins TB, Veeriah S, Shafi S, Johnson DH,
 848 Mitter R, Rosenthal R, et al.. 2017. Tracking the evolution of non–small-cell lung cancer. *New England*
 849 *Journal of Medicine* **376**: 2109–2121.
- 850 Khakabimamaghani S, Malikic S, Tang J, Ding D, Morin R, Chindelevitch L, and Ester M. 2019. Collaborative
 851 intra-tumor heterogeneity detection. *Bioinformatics* **35**: i379–i388.

852 Kuipers J, Moore AL, Jahn K, Schraml P, Wang F, Morita K, Futreal PA, Takahashi K, Beisel C, Moch H,
853 et al.. 2021. Statistical tests for intra-tumour clonal co-occurrence and exclusivity. *PLOS Computational*
854 *Biology* **17**: e1009036.

855 Luo XG, Kuipers J, and Beerenwinkel N. 2023. Joint inference of exclusivity patterns and recurrent trajec-
856 tories from tumor mutation trees. *Nature Communications* **14**: 3676.

857 Morita K, Wang F, Jahn K, Hu T, Tanaka T, Sasaki Y, Kuipers J, Loghavi S, Wang SA, Yan Y, et al.. 2020.
858 Clonal evolution of acute myeloid leukemia revealed by high-throughput single-cell genomics. *Nature*
859 *communications* **11**: 1–17.

860 Prüfer. 1918. Neuer beweis eines satzes uber permutationen. *Arch Math Phys* **27**: 742–4.

861 Razavi P, Chang MT, Xu G, Bandlamudi C, Ross DS, Vasan N, Cai Y, Bielski CM, Donoghue MT, Jonsson
862 P, et al.. 2018. The genomic landscape of endocrine-resistant advanced breast cancers. *Cancer cell* **34**:
863 427–438.

864 Turajlic S, Xu H, Litchfield K, Rowan A, Horswell S, Chambers T, O’Brien T, Lopez JI, Watkins TB, Nicol
865 D, et al.. 2018. Deterministic evolutionary trajectories influence primary tumor growth: TRACERx renal.
866 *Cell* **173**: 595–610.

867 Williams RJ. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learn-
868 ing. *Machine Learning* **8**: 229–256.

869 Yeang CH, McCormick F, and Levine A. 2008. Combinatorial patterns of somatic gene mutations in cancer.
870 *The FASEB journal* **22**: 2605–2622.