

# Pairwise and higher-order epistatic effects among somatic cancer mutations across oncogenesis

Jorge A. Alfaro-Murillo<sup>1</sup> and Jeffrey P. Townsend<sup>1,2,3</sup>

<sup>1</sup>Department of Biostatistics, Yale School of Public Health, New Haven, CT

<sup>2</sup>Department of Ecology and Evolutionary Biology, Yale University, New Haven, CT

<sup>3</sup>Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT

January 20, 2022

## Abstract

Cancer occurs as a consequence of multiple somatic mutations that lead to uncontrolled cell growth. Mutual exclusivity and co-occurrence of mutations imply—but do not prove—that they can exert synergistic or antagonistic epistatic effects on oncogenesis. Knowledge of these interactions, and the consequent trajectories of mutation and selection that lead to cancer has been a longstanding goal within the cancer research community. Recent research has revealed mutation rates and scaled selection coefficients for specific recurrent variants across many cancer types. However, estimation of pairwise and higher-order effects—essential to estimation of the trajectory of likely cancer genotypes—has been a challenge. Therefore, we have developed a continuous-time Markov chain model that enables the estimation of mutation origination and fixation (flux), dependent on somatic cancer genotype. Coupling the continuous-time Markov chain model with a deconvolution approach provides estimates of underlying rates of mutation and selection across the trajectory of oncogenesis. We demonstrate computation of fluxes and selection coefficients in a somatic evolutionary model for the four most frequently variant driver genes (TP53, LRP1B, KRAS and STK11) from 585 cases of lung adenocarcinoma. Our analysis reveals multiple antagonistic epistatic effects that reduce the possible routes of oncogenesis, and inform cancer research regarding viable trajectories of somatic evolution whose progression could be forestalled by precision medicine. Synergistic epistatic effects are also identified, most notably in the somatic genotype TP53+LRP1B for mutations in the KRAS gene, and in somatic genotypes containing KRAS or TP53 mutations for mutations in the STK11 gene. Large positive fluxes of KRAS variants were driven by large selection coefficients, whereas the

flux toward LRP1B mutations was substantially aided by a large mutation rate for this gene. The approach enables inference of the most likely routes of site-specific variant evolution and estimation of the strength of selection operating on each step along the route, a key component of what we need to know to develop and implement personalized cancer therapies.

## 1 Introduction

Abnormal cell proliferation and survival can be driven by gene mutations in somatic cells, and can result in cancer. The somatic mutations that lead to cancer can become frequent within tumors either because they are frequently mutated or because they are strongly selected to increase proliferation and survival. However, the selection operating on somatic mutations is complicated by epistatic effects, wherein the effects of one mutation affect the selection operating on other genes [1].

Computational models have been proposed that can indicate pairwise epistatic effects among mutations [2]. Empirically, patterns of mutual exclusivity can be interpreted as a consequence of antagonistic epistasis, and patterns of co-mutation can be interpreted as a consequence of synergistic epistasis. Other approaches have been applied to identify sets of variants that are sufficient to cause cancer [3]. However, to reveal the evolutionary genetic trajectories of cancer that could be intercepted by precision therapeutic approaches to delay or deny cancer morbidity and mortality due to metastasis, quantification of higher order epistatic effects are needed [4]. Models proposed so far do not appropriately condition on underlying mutation rates, nor do they reveal the order of mutations revealed by the epistatic interactions. For example, a mutation in a gene *A* could make it more likely that a mutation in a gene *B* is acquired, whereas when the mutation in the gene *B* occurs first, the mutation in gene *A* could be selected against. These differences in selection can arise from stage-specific physiological divergence or from additional unrecognized epistatic interactions.

The estimation of the order of epistatic effects in cancer is particularly challenging because most large tumor sequence datasets provide only one time point of tumor evolution at the time of tumor biopsy or excision [5]. Most evaluate correlations between the frequencies of mutations [2], yet correlations can arise either because of selective epistasis or because of commonalities of mutation process between selected sites. There are several orders of magnitude of difference in the rates at which somatic mutations occur in different genes and sites within the genome. Therefore, it is vital to distinguish the mutation rate from how much one mutation would confer a selective advantage in the tumor cell population [6]. The scaled selection coefficient or cancer effect size has previously averaged over epistatic effects [6]. Evaluating epistatic effects of selected mutations in cancer evolution is fundamental to illuminating the evolutionary genetic trajectory of tumor evolution, and for the advancement of accurate and personalized predictions of therapeutic responses [7, 8].

Herein, we develop mathematical models enabling the inference of the potential evolutionary trajectories of tumors. Our approaches extend from pairwise epistasis to combinations of three, four, or more cancer drivers. We apply our models to lung cancer, the most frequent cancer in the world [9], and the leading cause of cancer death in the United States [10]. Lung adenocarcinoma represents about 40% of all lung cancers, and has the worst prognosis among all types of lung cancer [11]. It typically has a high tumor mutational burden [12], and thus provides an excellent example in which to test theory on epistatic effects in cancer oncogenesis.

## 2 Theory

### 2.1 Mutational flux for one mutation with no epistasis

For a single mutation,  $\lambda$  can be taken as the exponentially distributed flux from a “normal” tissue state (without a mutation) at time 0 to a tissue state with a mutation fixed throughout a neoplasm, then

$$\text{Prob}\{\text{mutation by time } t\} = \int_0^t \lambda e^{-\lambda u} du = 1 - e^{-\lambda t}.$$

This flux can be decomposed into the mutation rate (the rate at which genetic state is changed in single cells) times the scaled selection coefficient (the consequent increase in survival and proliferation) [6].

### 2.2 Mutational fluxes for two mutations with epistasis

For two mutations labeled  $A$  and  $B$ , the fluxes from a normal state to a state with mutation  $A$  and to a state with mutation  $B$  can be denoted  $\lambda_A$  and  $\lambda_B$ . Assuming a regime of strong selection and weak mutation (SSWM) [13] with no clonal interference on the action of selection, and thereby retaining exponential distributions for the time that it takes for each mutation to appear and be selected to high frequency, the minimum of those distributions is also exponential, with the exponential parameter equal to the sum of the parameters of each distribution. Starting in the normal state at time 0, the probability density function for the time until the first mutation fixes is

$$f(t) = (\lambda_A + \lambda_B)e^{-(\lambda_A + \lambda_B)t}, \quad (1)$$

for  $t \geq 0$  and 0 otherwise. Therefore, the probability of maintaining the normal tissue state through time  $t$  is

$$\begin{aligned} P_0(t) &= \text{Prob}\{X(t) \text{ at normal state} \mid X(0) \text{ at normal state}\} \\ &= e^{-(\lambda_A + \lambda_B)t}, \end{aligned} \quad (2)$$

for  $t \geq 0$ , where  $X(t)$  represents the state at time  $t$ . Additionally, we know that under SSWM, the probability that a specific mutation spreads to fixation

before another mutation is equal to its rate relative to the sum of event rates, thus

$$\text{Prob}\{A \text{ before } B\} = \frac{\lambda_A}{\lambda_A + \lambda_B}. \quad (3)$$

Without epistasis, the fixation probabilities of mutation  $B$  are independent of whether mutation  $A$  has previously risen to fixation in the tumor. Therefore under SSWM the probabilities of fixing either mutation or both can be computed by multiplying the respective probabilities for each event. However, widespread mutual exclusivity among driver mutations indicates that epistatic interactions between them may be commonplace [9, 14, 15]. To model epistatic interactions between two driver mutations, two additional parameters are required: the flux to mutation  $A$  while at a state with mutation  $B$ , denoted  $\lambda_{B \rightarrow AB}$ , and the flux to  $B$  while at  $A$ , denoted  $\lambda_{A \rightarrow AB}$ . Thus, for a tissue fixed for mutation  $A$  at time  $u$  and  $t \geq u$ ,

$$\text{Prob}\{X(t) = A \mid X(u) = A\} = e^{-\lambda_{A \rightarrow AB}(t-u)}. \quad (4)$$

The probability of a tissue fixed with only mutation  $A$  at time  $t$  can be computed by multiplying Equations (1), (3), and (4), and integrating:

$$\begin{aligned} P_A(t) &= \text{Prob}\{X(t) = A \mid X(0) \text{ at normal state}\} \\ &= \int_0^t (\lambda_A + \lambda_B) e^{-(\lambda_A + \lambda_B)u} \frac{\lambda_A}{\lambda_A + \lambda_B} e^{-\lambda_{A \rightarrow AB}(t-u)} du \\ &= \frac{\lambda_A}{\lambda_{A \rightarrow AB} - (\lambda_A + \lambda_B)} \left( e^{-(\lambda_A + \lambda_B)t} - e^{-(\lambda_{A \rightarrow AB})t} \right). \end{aligned} \quad (5)$$

Equations (1), (3), and (4) compose three conditions that together yield the desired probability that a tissue fixed only mutation  $A$  at time  $t$  (Equation(5)). Equation (1) conditions that a mutation fixed at a time  $u$ , Equation (3) conditions for the case that the mutation was  $A$  instead of  $B$ , and Equation (4) conditions for the case that no other mutations were fixed from the time  $u$  that the mutation was fixed to the time  $t$ .

By symmetry,

$$\begin{aligned} P_B(t) &= \text{Prob}\{X(t) = B \mid X(0) \text{ at normal state}\} \\ &= \frac{\lambda_B}{\lambda_{B \rightarrow AB} - (\lambda_A + \lambda_B)} \left( e^{-(\lambda_A + \lambda_B)t} - e^{-(\lambda_{B \rightarrow AB})t} \right), \end{aligned} \quad (6)$$

and a formula for the probability that the neoplasm will be in a state fixed for both mutations follows from Equations (2), (5), and (6):

$$\begin{aligned} P_{AB}(t) &= \text{Prob}\{X(t) = AB \mid X(0) \text{ is normal}\} \\ &= 1 - P_0(t) - P_A(t) - P_B(t). \end{aligned} \quad (7)$$

For two mutations, Equations (2), (5), (6), and (7) provide probabilities for all possible genotypic states of the evolving neoplasm.

If three or more mutations are considered, an equation for the probability of having two mutations fixed cannot be obtained by subtracting the other probabilities as in Equation (7). Alternatively,  $P_{AB}(t)$  can be directly computed with the formula

$$P_{AB}(t) = \int_0^t (\lambda_A + \lambda_B) e^{-(\lambda_A + \lambda_B)u} \times \left( \frac{\lambda_A}{\lambda_A + \lambda_B} P_{A \rightarrow AB}(t-u) + \frac{\lambda_B}{\lambda_A + \lambda_B} P_{B \rightarrow AB}(t-u) \right) du,$$

where

$$\begin{aligned} P_{A \rightarrow AB}(t) &= \text{Prob}\{X(t) = AB \mid X(0) = A\}, \\ P_{B \rightarrow AB}(t) &= \text{Prob}\{X(t) = AB \mid X(0) = B\}, \end{aligned}$$

are probabilities that can be obtained by a similar argument as the one used to obtain Equation (5). However, as we only require  $P_0(t), P_A(t), P_B(t)$  and  $P_{AB}(t)$  for the likelihood (Section 2.4), it would be better to obtain a formula where  $P_{AB}(t)$  is in terms of  $P_A(t)$  and  $P_B(t)$ . We will obtain such an equation for the general case with  $M$  mutations.

### 2.3 Mutational fluxes for $M$ mutations with epistasis

To solve the general case of  $M$  possible somatic mutations, we can model the somatic genetic state of a neoplasm with respect to time  $t$  as a continuous-time Markov chain  $X(t)$ ,  $t \geq 0$ . We define the set of all possible states of the system as the  $M$ -ary Cartesian product  $\mathcal{S} = \{0, 1\}^M$ . Any state in the system is represented by a vector in  $\mathcal{S}$ :

$$\mathbf{x} = (x_1, \dots, x_M),$$

where  $x_i$  is 1 if the state has fixed the  $i$ -th mutation and 0 otherwise. Modeling two mutations  $A$  and  $B$ ,  $M = 2$ ; the normal state would be represented by  $(0, 0)$ , the state with only mutation  $A$  fixed by  $(1, 0)$ , the state with only  $B$  fixed by  $(0, 1)$ , and the state with both  $A$  and  $B$  fixed ( $AB$ ) with  $(1, 1)$ .

Under an SSWM regime, mutations occur and spread one at a time. Consequently, the flux from  $\mathbf{x}$  to  $\mathbf{y}$  is 0 unless  $\mathbf{y}$  has exactly one more mutation than  $\mathbf{x}$ , or  $\mathbf{y}$  is  $\mathbf{x}$ . The infinitesimal parameters  $\lambda_{\mathbf{x} \rightarrow \mathbf{y}}$  that determine the flux from state  $\mathbf{x}$  to  $\mathbf{y}$  are such that  $\lambda_{\mathbf{x} \rightarrow \mathbf{y}} = 0$  unless there is an  $i \in \{1, \dots, M\}$  with  $x_i = 0$ ,  $y_i = 1$  and  $x_j = y_j$  for all  $j \neq i$ , or  $\mathbf{y} = \mathbf{x}$ . The  $\mathbf{y} = \mathbf{x}$  case is relevant because—as is customary for continuous-time Markov chains—we define

$$\lambda_{\mathbf{x} \rightarrow \mathbf{x}} = - \sum_{\substack{\mathbf{y} \in \mathcal{S} \\ \mathbf{y} \neq \mathbf{x}}} \lambda_{\mathbf{x} \rightarrow \mathbf{y}}. \quad (8)$$

If  $\mathbf{x} \neq (1, \dots, 1)$ , Equation (8) can be simplified to

$$\lambda_{\mathbf{x} \rightarrow \mathbf{x}} = - \sum_{\substack{1 \leq i \leq M \\ \mathbf{e}_i \cdot (\mathbf{x} + \mathbf{e}_i) = 1}} \lambda_{\mathbf{x} \rightarrow \mathbf{x} + \mathbf{e}_i},$$

where  $\mathbf{e}_i$  represents the  $i$ -th vector of the standard basis of  $\mathbb{R}^M$ , that is, the vector with zeros everywhere except for a 1 in its  $i$ -th entry. The state with all mutations fixed  $\mathbf{x} = (1, \dots, 1)$  is an absorbing state, therefore the infinitesimal departure rate from the state is zero, i.e.  $\lambda_{(1, \dots, 1) \rightarrow (1, \dots, 1)} = 0$ .

Modeling  $M = 2$  mutations, all infinitesimal parameters can be represented by the continuous Markov matrix:

$$Q = \begin{pmatrix} -\lambda_{(0,0) \rightarrow (1,0)} - \lambda_{(0,0) \rightarrow (0,1)} & 0 & 0 & 0 \\ \lambda_{(0,0) \rightarrow (1,0)} & -\lambda_{(1,0) \rightarrow (1,1)} & 0 & 0 \\ \lambda_{(0,0) \rightarrow (0,1)} & 0 & -\lambda_{(0,1) \rightarrow (1,1)} & 0 \\ 0 & \lambda_{(1,0) \rightarrow (1,1)} & \lambda_{(0,1) \rightarrow (1,1)} & 0 \end{pmatrix},$$

with the ordered rows and columns representing the states  $(0, 0)$ ,  $(1, 0)$ ,  $(0, 1)$  and  $(1, 1)$ . Equivalently, this matrix can be written with the notation of Section 2.2 as:

$$Q = \begin{pmatrix} -\lambda_A - \lambda_B & 0 & 0 & 0 \\ \lambda_A & -\lambda_{A \rightarrow AB} & 0 & 0 \\ \lambda_B & 0 & -\lambda_{B \rightarrow AB} & 0 \\ 0 & \lambda_{A \rightarrow AB} & \lambda_{B \rightarrow AB} & 0 \end{pmatrix},$$

with ordered rows and columns representing tissue in a normal state, a state with only the first mutation  $A$ , a state with only the second mutation  $B$ , and a state with both mutations  $AB$ .

Because  $X(t)$  is a continuous-time Markov chain, the probabilities that a neoplasm starts in a state  $\mathbf{y}$  at time  $u$  and is in a state  $\mathbf{x}$  at time  $t + u$  are independent of  $u$ , so we can denote them as:

$$P_{\mathbf{y} \rightarrow \mathbf{x}}(t) = \text{Prob} \{X(t + u) = \mathbf{x} \mid X(u) = \mathbf{y}\},$$

for any  $t, u \geq 0$ . Applying Kolmogorov's backward equation [16] to the matrix  $P(t)$  with entries  $p_{ij} = P_{\mathbf{x}_i \rightarrow \mathbf{x}_j}(t)$ , for any ordering  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{2^M}$  of all possible states, we obtain the differential equation

$$P'(t) = P(t)Q, \tag{9}$$

where  $Q$  is the continuous Markov matrix for the same ordering of all possible states. Thus we could find the solution for  $P(t)$  by computing the matrix exponential for  $Q$  and applying the Fundamental Theorem for Linear Systems [17, Section 1.4], then:

$$P(t) = e^{Qt}X(0).$$

If all eigenvalues of the matrix  $Q$  are real and distinct then the exponential matrix can be found by finding its diagonalization [17, Section 1.2]. Performing this operation and solving for the probabilities of being in the mutation states associated with  $M = 2$  ( $0$ ,  $A$ ,  $B$ , and  $AB$ ) yields Equations (2), (5), (6), and (7). However, a direct solution for the exponential matrix is exhaustive to compute for large matrices [17, Section 1.8]. The size of the matrix  $Q$  grows exponentially with  $M$ . Even for the case of  $M = 3$  an  $8 \times 8$  matrix is required, and each

entry of  $e^Q$  is prohibitively complicated. To develop an alternate approach toward quantification of the probabilities of state for  $M$  mutations, a main property of our estimation problem may be capitalized upon: it can be assumed that all individuals start at the normal state, that is, without relevant somatic mutations. Therefore, only one column in the matrix  $P$  is of interest, the one that includes the entries  $P_{\mathbf{0} \rightarrow \mathbf{x}}$ . To simplify notation, we will write

$$P_{\mathbf{x}}(t) = P_{\mathbf{0} \rightarrow \mathbf{x}}(t),$$

for any  $t \geq 0$ . The relevant column in Equation (9) reduces to

$$P'_{\mathbf{x}}(t) = \sum_{\mathbf{y} \in \mathcal{S}} P_{\mathbf{y}}(t) \lambda_{\mathbf{y} \rightarrow \mathbf{x}}, \quad (10)$$

where the specific ordering of state transitions is no longer important as it was in Equation (9). By the properties of the continuous-time Markov chain, we know the initial condition for this differential equation:

$$P_{\mathbf{x}}(0) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{0}, \\ 0 & \text{if } \mathbf{x} \neq \mathbf{0}. \end{cases} \quad (11)$$

For the case of  $\mathbf{x} = \mathbf{0} = (0, \dots, 0)$ , that is, the normal tissue state, Equation (10) becomes

$$P'_{\mathbf{0}}(t) = \lambda_{\mathbf{0} \rightarrow \mathbf{0}} P_{\mathbf{0}}(t).$$

Solving this differential equation with the initial condition specified by Equation (11), we have

$$P_{\mathbf{0}}(t) = e^{\lambda_{\mathbf{0} \rightarrow \mathbf{0}} t} = e^{-(\sum_j \lambda_{\mathbf{0} \rightarrow \mathbf{e}_j}) t}, \quad (12)$$

where the sum goes for all  $j = 1, \dots, M$ .

For the case that  $\mathbf{x} \neq \mathbf{0}$ , using Equation (10) and the definition of the infinitesimal parameters,

$$P'_{\mathbf{x}}(t) = \lambda_{\mathbf{x} \rightarrow \mathbf{x}} P_{\mathbf{x}}(t) + \sum_{\substack{1 \leq i \leq M \\ \mathbf{e}_i \cdot \mathbf{x} = 1}} \lambda_{\mathbf{x} - \mathbf{e}_i \rightarrow \mathbf{x}} P_{\mathbf{x} - \mathbf{e}_i}(t). \quad (13)$$

Solving Equation (13) with the integral factor  $e^{-(\lambda_{\mathbf{x} \rightarrow \mathbf{x}}) t}$  and the initial condition in Equation (11) provides a recursive formula

$$P_{\mathbf{x}}(t) = \sum_{\substack{1 \leq i \leq M \\ \mathbf{e}_i \cdot \mathbf{x} = 1}} \lambda_{\mathbf{x} - \mathbf{e}_i \rightarrow \mathbf{x}} \int_0^t e^{\lambda_{\mathbf{x} \rightarrow \mathbf{x}}(t-u)} P_{\mathbf{x} - \mathbf{e}_i}(u) du. \quad (14)$$

This formula enables computation of  $P_{\mathbf{x}}(t)$  for all states  $\mathbf{x}$  by starting with the normal state using Equation (12), then proceeding to compute all states with one mutation (where Equation (14) depends on  $P_{\mathbf{0}}(t)$ ), then all states with two mutations (where Equation (14) depends on  $P_{\mathbf{x}}(t)$  for states  $\mathbf{x}$  with only one mutation), and so on.

A straightforward validation case of the Equation (14) comes from its application when  $\mathbf{x}$  already features exactly one mutation. Applied to that case,

$$\begin{aligned} P_{\mathbf{e}_i}(t) &= \frac{\lambda_{\mathbf{0} \rightarrow \mathbf{e}_i}}{\lambda_{\mathbf{0} \rightarrow \mathbf{0}} - \lambda_{\mathbf{e}_i \rightarrow \mathbf{e}_i}} (e^{\lambda_{\mathbf{0} \rightarrow \mathbf{0}} t} - e^{\lambda_{\mathbf{e}_i \rightarrow \mathbf{e}_i} t}) \\ &= \frac{\lambda_{\mathbf{0} \rightarrow \mathbf{e}_i}}{\sum_{j \neq i} \lambda_{\mathbf{e}_i \rightarrow \mathbf{e}_i + \mathbf{e}_j} - \sum_j \lambda_{\mathbf{0} \rightarrow \mathbf{e}_j}} \left( e^{-(\sum_j \lambda_{\mathbf{0} \rightarrow \mathbf{e}_j}) t} - e^{-(\sum_{j \neq i} \lambda_{\mathbf{e}_i \rightarrow \mathbf{e}_i + \mathbf{e}_j}) t} \right), \end{aligned}$$

which agrees with Equation (5) for the case of  $M = 2$ .

## 2.4 Likelihood of observed frequencies of tumors

To quantify the flux associated with one mutation with no epistasis (Section 2.1), if we have a sample of  $N$  tumors within which  $n$  tumors exhibited the variant site, we can assume that the samples are taken at a similar time  $T$  from an initial state without the mutation assessed. The likelihood is binomial:

$$\mathcal{L}(n \text{ tumors with mutations} \mid \lambda) \propto (1 - e^{-\lambda T})^n (e^{-\lambda T})^{N-n},$$

and can be maximized to obtain an estimate of  $\lambda$ .

With two mutations and accounting for epistatic effects (Section 2.2), the genotypes of  $N$  tumors can be subdivided into  $n_{\mathbf{0}}$ ,  $n_A$ ,  $n_B$ , and  $n_{AB}$ , representing those tumors without any mutations, with only  $A$ , with only  $B$ , and with both  $A$  and  $B$ . To estimate all the fluxes, we set the time again at  $t = T$  and the likelihood is multinomial:

$$\mathcal{L}(n_{\mathbf{0}}, n_A, n_B, n_{AB} \mid \lambda_A, \lambda_B, \lambda_{B \rightarrow AB}, \lambda_{A \rightarrow AB}) \propto \prod_{x \in \{\mathbf{0}, A, B, AB\}} (P_x(T))^{n_x}.$$

The maximization of this likelihood will provide estimators of the fluxes.

For the general case with an arbitrary number of mutations (Section 2.3), a sample of  $N$  tumors that have been sequenced can be divided according to the somatic genotype in  $\mathcal{S}$  attributed to each tumor. Let  $n_{\mathbf{x}}$  be the total number of samples that have the somatic genotype  $\mathbf{x}$  for each  $\mathbf{x} \in \mathcal{S}$ . We assume that the samples are taken at a similar time  $T$  (of potentially arbitrary unit) from an initial state without any of the  $M$  mutations assessed (i.e. the “normal” state  $\mathbf{0}$ ). Consequently, the probability  $P_{\mathbf{x}}(T)$  reflects the fraction of the cases observed with the somatic genotype  $\mathbf{x}$ . The likelihood  $\mathcal{L}$  is then multinomial:

$$\mathcal{L}(\{n_{\mathbf{x}}\}_{\mathbf{x} \in \mathcal{S}} \mid \{\lambda_{\mathbf{x} \rightarrow \mathbf{y}}\}_{\mathbf{x}, \mathbf{y} \in \mathcal{S}}) \propto \prod_{\mathbf{x} \in \mathcal{S}} (P_{\mathbf{x}}(T))^{n_{\mathbf{x}}}, \quad (15)$$

where  $\sum_{\mathbf{x} \in \mathcal{S}} n_{\mathbf{x}} = N$ , and  $P_{\mathbf{x}}(T)$  is dependent on the fluxes  $\lambda_{\mathbf{x} \rightarrow \mathbf{y}}$  according to Equations (12) and (14).

## 3 Methods

We obtained single-nucleotide variants (SNV) from 585 cases of lung adenocarcinoma from The Cancer Genome Atlas, and classified them into genes



following human genome coordinates from the Ensembl Project [18]. All mutations were used to determine the baseline mutation rate. However, synonymous mutations were removed for the purpose of tallying prevalence of selected mutations because they are typically not selected. To estimate epistatic effects among four genes using the quadruplewise case  $M = 4$ , we restricted our analysis to the four genes most commonly fixed for mutations among the 585 samples.

We set  $T = 1$  in Equation (15) to evaluate the likelihood, equating one unit to the average duration between the clonal origin of the tumor and tumor sampling; all rates and fluxes derived with  $T = 1$  are thus in units of 1 over this unit. To compute the integral on the right-hand side of the recursive formula in Equation (14), we used a trapezoidal rule with a resolution of 1,000 points between 0 and  $T$ . We tested higher resolutions, and our results were unchanged.

We estimated the fluxes  $\lambda_{\mathbf{x} \rightarrow \mathbf{y}}$  by maximizing the likelihood in Equation (15). We computed asymptotic confidence intervals for each of the flux estimates by computing the log-likelihood ratio and using Wilk's theorem [19]. To validate model fit, we compared for each state  $\mathbf{x}$  the probability  $P_{\mathbf{x}}(1)$  evaluated with the flux estimates Equation (14) to the observed fraction of the samples in each category. The values were equal for each somatic genotype.

To factor each flux into a mutation rate and a scaled selection coefficient, we assumed that the mutation rate per gene did not change with the acquisition of the somatic mutations of interest. We used `cancereffectsizeR` to obtain gene-specific mutation rates [6].

## 4 Results

The genes that were most frequently mutated were tumor protein p53 (TP53,  $n = 278$ ), low-density lipoprotein receptor-related protein 1B (LRP1B,  $n = 183$ ), Kirsten rat sarcoma (KRAS,  $n = 150$ ) and the tumor suppressor serine/threonine kinase 11 (STK11,  $n = 82$ ). The numbers of patients with each somatic genotype informed the likelihood of our model, which provided estimates and confidence intervals for the flux and scaled selection coefficient of each mutation in the context of each somatic genotype for these four genes (Table 1).

Fixed-value estimates for the mutation rates for each gene over oncogenesis led to a gene-specific linear relationship between the flux and the scaled selection coefficient (Figure 1). A consequence of the multiplicative relationship between mutation rate and scaled selection coefficient is that mutation rate determined the slope of response of flux to scaled selection coefficient. Consequently, underlying mutation rate had a substantial effect on flux and prevalence of somatic genotypes, especially in the context of less-strongly selected mutations such as those in LRP1B. Mutations to KRAS were the least frequent of any gene (Table 1), and exhibited the lowest slope of flux in response to scaled selection coefficient.

Despite the lower slope of response to scaled selection coefficient evident for somatic mutations of TP53 and KRAS, the scaled selection coefficients for

Table 1: Fluxes, mutation rates, and scaled selection coefficients for a four-gene model of lung adenocarcinoma oncogenesis.

| Genotype         | Mutation | Flux (CI) <sup>a</sup> | Mutation rate         | SSC <sup>b</sup> (CI) <sup>a</sup> |
|------------------|----------|------------------------|-----------------------|------------------------------------|
| normal           | TP53     | 0.84 (0.73, 0.95)      | $1.98 \times 10^{-6}$ | 422 (368, 482)                     |
| normal           | KRAS     | 0.31 (0.25, 0.38)      | $1.18 \times 10^{-6}$ | 267 (216, 325)                     |
| normal           | STK11    | 0.08 (0.05, 0.11)      | $2.09 \times 10^{-6}$ | 37 ( 24, 54)                       |
| normal           | LRP1B    | 0.19 (0.13, 0.26)      | $9.86 \times 10^{-6}$ | 19 ( 14, 26)                       |
| TP53             | KRAS     | 0.29 (0.19, 0.41)      | $1.18 \times 10^{-6}$ | 243 (160, 351)                     |
| TP53             | STK11    | 0.13 (0.07, 0.23)      | $2.09 \times 10^{-6}$ | 64 ( 33, 108)                      |
| TP53             | LRP1B    | 0.84 (0.64, 1.07)      | $9.86 \times 10^{-6}$ | 85 ( 65, 109)                      |
| KRAS             | TP53     | 0 ( 0, 0.3)            | $1.98 \times 10^{-6}$ | 0 ( 0, 150)                        |
| KRAS             | STK11    | 0.82 (0.54, 1.19)      | $2.09 \times 10^{-6}$ | 392 (261, 569)                     |
| KRAS             | LRP1B    | 0.56 (0.34, 0.86)      | $9.86 \times 10^{-6}$ | 56 ( 35, 87)                       |
| STK11            | TP53     | 0 ( 0, 0.6)            | $1.98 \times 10^{-6}$ | 0 ( 0, 305)                        |
| STK11            | KRAS     | 0 ( 0, 0.69)           | $1.18 \times 10^{-6}$ | 0 ( 0, 585)                        |
| STK11            | LRP1B    | 0.41 (0.14, 0.91)      | $9.86 \times 10^{-6}$ | 42 ( 15, 93)                       |
| LRP1B            | TP53     | 1.02 (0.51, 1.68)      | $1.98 \times 10^{-6}$ | 516 (258, 849)                     |
| LRP1B            | KRAS     | 0 ( 0, 0.38)           | $1.18 \times 10^{-6}$ | 0 ( 0, 319)                        |
| LRP1B            | STK11    | 0 ( 0, 0.23)           | $2.09 \times 10^{-6}$ | 0 ( 0, 109)                        |
| TP53+KRAS        | STK11    | 0.52 (0.16, 1.21)      | $2.09 \times 10^{-6}$ | 249 ( 79, 581)                     |
| TP53+KRAS        | LRP1B    | 0.21 ( 0, 1.17)        | $9.86 \times 10^{-6}$ | 21 ( 0, 119)                       |
| TP53+STK11       | KRAS     | 0 ( 0, 1.21)           | $1.18 \times 10^{-6}$ | 0 ( 0, 1022)                       |
| TP53+STK11       | LRP1B    | 0.52 ( 0, 2.01)        | $9.86 \times 10^{-6}$ | 53 ( 0, 204)                       |
| KRAS+STK11       | TP53     | 0 ( 0, 0.59)           | $1.98 \times 10^{-6}$ | 0 ( 0, 298)                        |
| KRAS+STK11       | LRP1B    | 0.82 (0.35, 1.63)      | $9.86 \times 10^{-6}$ | 84 ( 36, 166)                      |
| TP53+LRP1B       | KRAS     | 0.72 (0.45, 1.08)      | $1.18 \times 10^{-6}$ | 608 (384, 911)                     |
| TP53+LRP1B       | STK11    | 0.21 (0.07, 0.42)      | $2.09 \times 10^{-6}$ | 100 ( 34, 203)                     |
| KRAS+LRP1B       | TP53     | 0 ( 0, 1.0)            | $1.98 \times 10^{-6}$ | 0 ( 0, 504)                        |
| KRAS+LRP1B       | STK11    | 0.08 ( 0, 0.85)        | $2.09 \times 10^{-6}$ | 36 ( 0, 408)                       |
| STK11+LRP1B      | TP53     | 0 ( 0, 2.99)           | $1.98 \times 10^{-6}$ | 0 ( 0, 1513)                       |
| STK11+LRP1B      | KRAS     | 0 ( 0, 2.95)           | $1.18 \times 10^{-6}$ | 0 ( 0, 2500)                       |
| TP53+KRAS+STK11  | LRP1B    | 0.92 ( 0, 5.92)        | $9.86 \times 10^{-6}$ | 93 ( 0, 601)                       |
| TP53+KRAS+LRP1B  | STK11    | 0.23 ( 0, 0.86)        | $2.09 \times 10^{-6}$ | 110 ( 0, 415)                      |
| TP53+STK11+LRP1B | KRAS     | 0 ( 0, 1.59)           | $1.18 \times 10^{-6}$ | 0 ( 0, 1352)                       |
| KRAS+STK11+LRP1B | TP53     | 0 ( 0, 1.77)           | $1.98 \times 10^{-6}$ | 0 ( 0, 892)                        |

<sup>a</sup> 95% confidence interval.

<sup>b</sup> Scaled selection coefficient (in thousands).

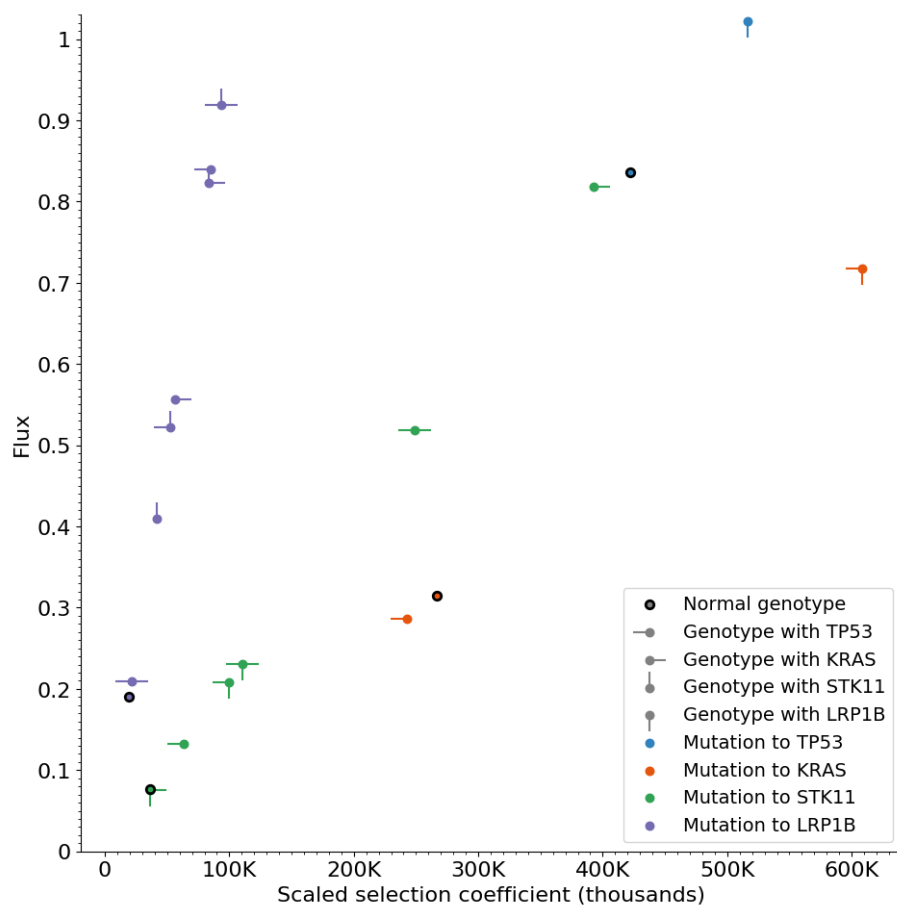


Figure 1: Estimates of positive scaled selection coefficients and positive genotypic fluxes for a four-gene model of lung adenocarcinoma oncogenesis. Each genotype (indicated by the directions of ticks superimposed on each point) has distinct effects on the flux that are mediated by epistasis affecting the scaled selection coefficients of new mutations (TP53, blue; KRAS, red; STK11, green; LRP1B, purple). No points are shown from one genotype to another where the estimated selection coefficient—and consequently the flux—was zero.

the positively selected mutations to these two genes were typically high in lung adenocarcinoma (Figure 1). Especially notably, scaled selection coefficients for TP53 and KRAS mutations were higher than for mutations in other genes at the initiation of oncogenesis. Consequently, from a normal somatic genotype the flux to TP53 was vastly higher than for any other mutation, with KRAS a distant second, and LRP1B and STK11 following. Even in later stages of oncogenesis when the flux to LRP1B was quite large, the magnitude of that larger flux could be attributed to the high mutation rate of LRP1B (Table 1, Figure 2), rather than to its relatively small scaled selection coefficient.

Comparison of the scaled selection coefficient for fixation of the first mutations to the scaled selection coefficient for that mutation in non-normal genotypes frequently quantified an antagonistic epistatic effect of somatic mutations to oncogenic drivers (Table 1, Figure 2C). Antagonistic epistatic effects, in turn, likely explain the low number of patients that had mutations in three or more of the four studied genes (Figure 2A). Antagonistic epistatic effects also partitioned the order of mutation fixation into three classes of routes: (i) STK11 then LRP1B; (ii) KRAS, then STK11 or LRP1B, and then the remainder of LRP1B or (less often) of STK11; and (iii) more complex routes with a first fixation of either LRP1B or TP53 (Figure 2A).

Some synergistic epistatic effects were also evident. For example, any genotype with a KRAS or TP53 mutation substantially increased the scaled selection coefficient on LRP1B mutation compared to when neither gene was mutated (Table 1, Figure 2C). The largest significant synergistic epistatic effects were the presence of TP53 and LRP1B when acquiring the KRAS mutation, the presence of KRAS when acquiring STK11, and the presence of TP53 and KRAS when acquiring the STK11 mutation (Table 1, Figure 2C).

When considering each one of the 24 possible paths of mutation acquisition we observed that fitness increases after an initial TP53 mutation and decreases afterwards (Figure 3). A similar situation occurs with an initial KRAS mutation, except when STK11 mutates afterwards (Figure 3). Initial acquisition of STK11 or LRP1B mutations left relatively low selection coefficients for additional mutations of these genes, with one notable exception: when LRP1B was followed by TP53 (Figure 3).

## 5 Discussion

Here we have shown how to estimate high-order epistatic effects on the somatic selection of cancer mutations. Our model enables computation of the flux in somatic genetic state of tissue from one genotype to another by single mutations. Application of our model to 585 lung adenocarcinoma samples provided estimates of 32 fluxes to the 4 most commonly mutated genes (TP53, LRP1B, KRAS and STK11), and showed that the flux to each of those four genes depended on the current somatic genotype. Many genotypes exhibited an antagonistic epistatic effect that resulted in zero or near-zero fluxes out of that genotype. Antagonistic epistatic effects likely partly explain the relatively low

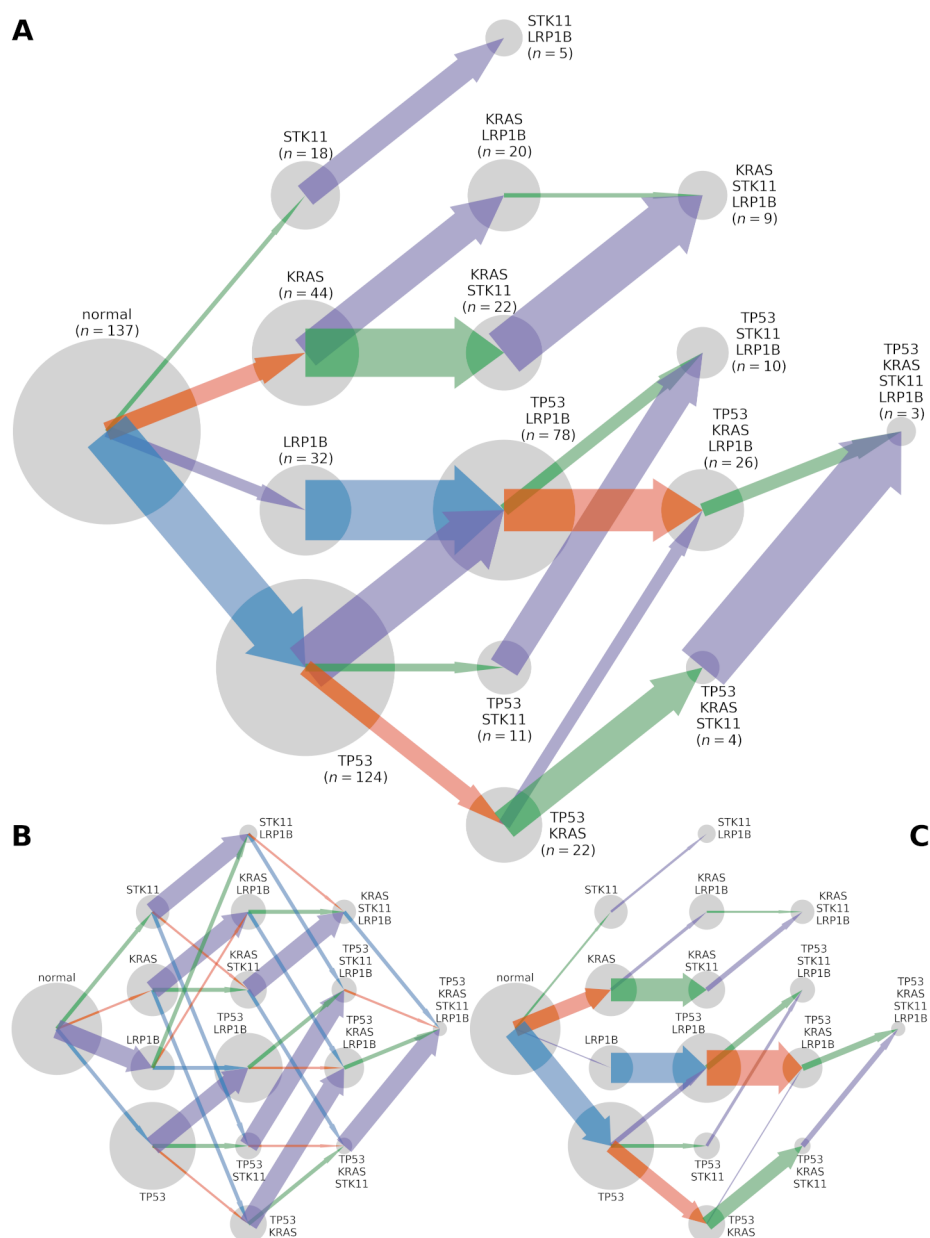


Figure 2: Trajectories of the somatic evolution by mutation of TP53, KRAS, LRP1B, and STK11, inferred from a total of 585 whole-exome sequenced lung adenocarcinoma tumors. Genotypes (grey circles; areas are proportional to observed  $n$  for the genotype) evolve at (A) fluxes, (B) mutation rates, and (C) scaled selection coefficients that are proportional to the width of arrows pointing from one genotype to another, colored by the gene in which the mutation occurs (TP53, blue; KRAS, orange; STK11, green; LRP1B: purple).

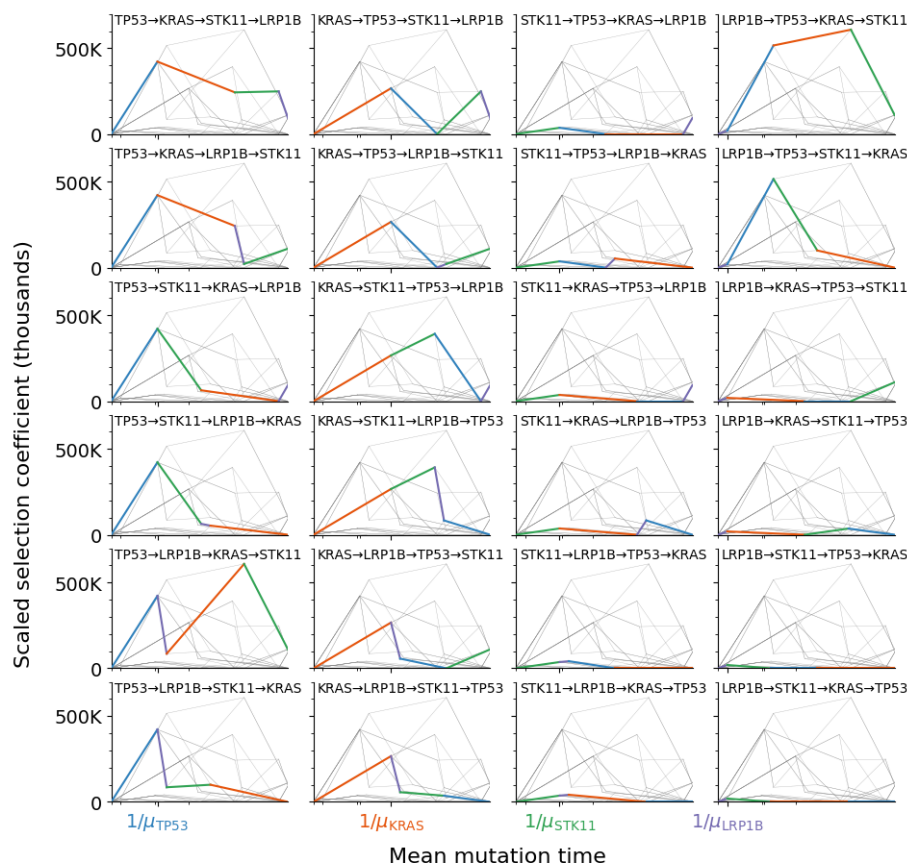


Figure 3: Mutation landscape depicting all paths of mutation acquisition in a four-gene model for lung adenocarcinoma. The four mutations incorporated are TP53 (blue, initial mutation in column one), KRAS (orange, initial mutation in column 2), STK11 (green, initial mutation in column three), and LRP1B (purple, initial mutation in column four). The mean time to acquire a mutation in a cancer-competent cell lineage ( $x$ -axis) quantifies how quickly each mutation will occur on a cellular level in the landscape of mutations (shorter is quicker), whereas the scaled selection coefficient ( $y$ -axis) is a measure of the benefit of the mutation to lineage proliferation and survival (the higher the selection coefficient, the more likely a mutation, once it occurs, will spread to high frequency in tumor tissue). In every subplot, one path is highlighted (colored curve, which corresponds to the order indicated by the gene names above the curve) and contrasted with all other possible paths (gray curves).

number of tumor samples that contain mutations in more than three of these commonly-mutated genes.

By estimating the neutral mutation rates of each gene, we also computed the corresponding scaled selection coefficients quantifying the degree to which the mutations increased survival and proliferation. We found that KRAS mutations exhibited scaled selection coefficients that were especially large and were a major reason for their high fluxes, especially from the TP53+LRP1B genotype. In contrast, most of the flux to LRP1B mutations—and thus the large number of samples with LRP1B mutations—can be explained by LRP1B’s large mutation rate, not so much by a large cancer effect. Despite their lower overall effect size, LRP1B mutations in lung adenocarcinoma have been associated with chronic obstruction pulmonary disease [20], and suggested as predictors of response to immune checkpoint inhibitors [21]. Additionally, LRP1B appears to cooperate with TP53 to induce a large selection for mutant KRAS as a driver of lung adenocarcinoma [3].

KRAS and TP53 have been suggested to play a role in lung adenocarcinoma initiation [22]. KRAS mutations are consistently revealed by sequencing of all tumor grades [23], and because of identification of copy number gains that enable estimation of the relative timing of somatic alterations [24]. Our analysis provided a result consistent with these findings: our results show that KRAS is subject to a positive, synergistic epistatic effect on its selection in the context of LRP1B and STK11 mutations. Interestingly, our quantification of these epistatic effects argues that the order in which KRAS and TP53/STK11 mutations are acquired is relevant to their selective effect. Acquisition of a TP53 mutation before a KRAS mutation does little to change selection on a new KRAS mutation. However, if a KRAS mutation is acquired first, selection on TP53 mutation almost disappears. Conversely, acquisition of a STK11 mutation from a normal state prevents the selection of KRAS—but if KRAS is acquired first, there is strong selection for STK11. TP53 and STK11 have been identified as determinants of distinct subsets of lung adenocarcinomas dominated by KRAS mutants [25]. Importantly, patients whose tumors have KRAS mutations but no TP53 have better overall survival than those with both mutations, especially in the absence of STK11 mutations [26, 27], a situation that arises more often when a LRP1B mutation occurs immediately after KRAS is mutated.

STK11 has previously been identified as a mutation that occurs relatively early in the oncogenesis of lung adenocarcinoma [24]. An early, relatively simple route of mutations starts with a STK11 mutation and continues with an LRP1B mutation, without further mutation of KRAS or TP53 prior to tumor resection. Additionally, we have estimated strong selection for STK11 mutations after KRAS mutation, indicating that if KRAS mutations occur early, then STK11 mutations would be strongly selected and once mutated would fix not long after. However, we have also found a relatively constant selection for STK11 from other genotypes, that—despite being lower than when KRAS mutates first—suggests that STK11 does not occur exclusively early during tumor initiation.

We focused here on single-nucleotide variants that affect the genes studied. However, there are other somatic factors that could have epistatic effects, such

as copy number alterations and structural variants. Our model can easily incorporate those factors and compute the fluxes associated with their fixation. However, the estimation of the underlying mutation rate for copy number alterations and structural variants remains a challenge. Consequently, quantifying the strength of selection operating on them is not feasible. Thus far, research has indicated that copy-number changes and structural variants do not substantially influence selection on simple nucleotide variations, and have instead orthogonal effects in cancers [28].

Confidence intervals for our estimates became wider as the genotypes included mutations in more genes, because the number of tumors with mutations in three or four of these commonly-mutated genes tended to be very low. To reduce the parameter uncertainty, future research should explore larger data sets by aggregating whole-exome sequences and panel data from multiple sources. Alternatively, data that includes tumor samples at multiple time points could better inform the order at which mutations occur, reducing uncertainty. The equations we have derived enable estimation of the time-dependent probability of each somatic genotype given the flux values. Thus, our theory can be extended to applications quantifying cancer effects on tumor samples at metachronous time points.

Using an established approach [6], we were able to estimate the mutation rate in each gene. However, this approach does not evaluate whether somatic genotypes vary in gene mutation rate. Mutation rates that depend on somatic genotype could affect the co-occurrence or mutual exclusivity of future mutations [29–32]. For example, it has long been suggested that mutations in TP53 cause impairment in the response to DNA damage, which might lead to higher mutation rates in other genes [32, 33], and even increases in specific genomic hotspots. Smoking is another major factor that increases the mutation rate of certain oncogenic sites in lung cancer, including KRAS [34]. Incorporation of differences in mutation rates depending on endogenous and exogenous mutational processes associated with somatic genotype and environmental exposures such as smoking, could enable increasingly precise estimation of selection coefficients attributable to each mutation [32].

In summary, we have developed a new approach to estimate fluxes and selection coefficients dependent on genotypes among somatic cancer mutations, and have employed this new approach to obtain epistatic effects among the four most commonly mutated genes in lung adenocarcinoma. We found several antagonistic and synergistic epistatic effects that reduced the space of possible routes of mutation acquisitions substantially and quantified how likely each path is. Determining the most likely trajectories of mutation across the time course of oncogenesis can help to determine optimal personalized therapies that account for the current somatic genotype of tumor tissue from a patient and that are designed to forestall the somatic genotypic trajectories that are most likely to be forthcoming.



## 6 Acknowledgements

Data used for this research were generated by the TCGA Research Network: <https://www.cancer.gov/tcga>. We would like to thank Katherine Brumberg and Vincent Cannataro for early discussions of ideas regarding this research. This research was supported by NSF IOS 1934848, NIH 1R01LM013385, and NIH NIDCR 1P50DE030707.

## References

- [1] Wang, X, Fu, A. Q, Mc Nerney, M. E, & White, K. P. (2014) Widespread genetic epistasis among cancer genes. *Nature communications* **5**, 1–10.
- [2] Manavalan, R & Priya, S. (2021) Genetic interactions effects for cancer disease identification using computational models: a review. *Medical & Biological Engineering & Computing* **59**, 733–758.
- [3] Klein, M. I, Cannataro, V. L, Townsend, J. P, Newman, S, Stern, D. F, & Zhao, H. (2021) Identifying modules of cooperating cancer drivers. *Molecular systems biology* **17**, e9810.
- [4] Baryshnikova, A, Costanzo, M, Myers, C. L, Andrews, B, & Boone, C. (2013) Genetic interaction networks: toward an understanding of heritability. *Annual review of genomics and human genetics* **14**, 111–133.
- [5] Black, J. R & McGranahan, N. (2021) Genetic and non-genetic clonal diversity in cancer evolution. *Nature Reviews Cancer* **21**, 379–392.
- [6] Cannataro, V. L, Gaffney, S. G, & Townsend, J. P. (2018) Effect sizes of somatic mutations in cancer. *JNCI: Journal of the National Cancer Institute* **110**, 1171–1177.
- [7] Wilkins, J. F, Cannataro, V. L, Shuch, B, & Townsend, J. P. (2018) Analysis of mutation, selection, and epistasis: an informed approach to cancer clinical trials. *Oncotarget* **9**, 22243.
- [8] Dasari, K, Somarelli, J. A, Kumar, S, & Townsend, J. P. (2021) The somatic molecular evolution of cancer: Mutation, selection, and epistasis. *Progress in biophysics and molecular biology* **165**, 56–65.
- [9] Bade, B. C & Cruz, C. S. D. (2020) Lung cancer 2020: epidemiology, etiology, and prevention. *Clinics in chest medicine* **41**, 1–24.
- [10] U.S. Cancer Statistics Working Group. (2021) U.S. Cancer Statistics Data Visualizations Tool, based on 2020 submission data (1999–2018) (U.S. Department of Health and Human Services, Centers for Disease Control and Prevention and National Cancer Institute; [www.cdc.gov/cancer/dataviz](http://www.cdc.gov/cancer/dataviz), released in June 2021).

- [11] Myers, D. J. (2021) Cancer, lung adenocarcinoma (StatPearls Publishing; <https://www.statpearls.com/articlelibrary/viewarticle/24486/>, last updated when retrieved, September 2021).
- [12] Cancer Genome Atlas Research Network. (2014) Comprehensive molecular profiling of lung adenocarcinoma: The Cancer Genome Atlas research network. *Nature* **511**, 543–550.
- [13] Gillespie, J. H. (1983) Some properties of finite populations experiencing strong selection and weak mutation. *The American Naturalist* **121**, 691–708.
- [14] Roy, D. M, Walsh, L. A, & Chan, T. A. (2014) Driver mutations of cancer epigenomes. *Protein & cell* **5**, 265–296.
- [15] van de Haar, J, Canisius, S, Michael, K. Y, Voest, E. E, Wessels, L. F, & Ideker, T. (2019) Identifying epistasis in cancer genomes: a delicate affair. *Cell* **177**, 1375–1383.
- [16] Kolmogoroff, A. (1931) Über die analytischen Methoden in der Wahrscheinlichkeitsrechnung. *Mathematische Annalen* **104**, 415–458.
- [17] Perko, L. (2001) *Differential Equations and Dynamical Systems*. (Springer), Third edition.
- [18] Howe, K. L, Achuthan, P, Allen, J, Allen, J, Alvarez-Jarreta, J, Amode, M. R, Armean, I. M, Azov, A. G, Bennett, R, Bhai, J, Billis, K, Boddu, S, Charkhchi, M, Cummins, C, Da Rin Fioretto, L, Davidson, C, Dodiya, K, El Houdaigui, B, Fatima, R, Gall, A, Garcia Giron, C, Grego, T, Guijarro-Clarke, C, Haggerty, L, Hemrom, A, Hourlier, T, Izuogu, O. G, Juettemann, T, Kaikala, V, Kay, M, Lavidas, I, Le, T, Lemos, D, Gonzalez Martinez, J, Marugán, J. C, Maurel, T, McMahon, A. C, Mohanan, S, Moore, B, Muffato, M, Oheh, D. N, Paraschas, D, Parker, A, Parton, A, Prosovetskaia, I, Sakthivel, M. P, Salam, A. I. A, Schmitt, B. M, Schuilenburg, H, Sheppard, D, Steed, E, Szpak, M, Szuba, M, Taylor, K, Thormann, A, Threadgold, G, Walts, B, Winterbottom, A, Chakiachvili, M, Chaubal, A, De Silva, N, Flint, B, Frankish, A, Hunt, S. E, IIsley, G. R, Langridge, N, Loveland, J. E, Martin, F. J, Mudge, J. M, Morales, J, Perry, E, Ruffier, M, Tate, J, Thybert, D, Trevanion, S. J, Cunningham, F, Yates, A. D, Zerbino, D. R, & Flicek, P. (2020) Ensembl 2021. *Nucleic Acids Research* **49**, D884–D891.
- [19] Wilks, S. S. (1938) The large-sample distribution of the likelihood ratio for testing composite hypotheses. *The annals of mathematical statistics* **9**, 60–62.
- [20] Xiao, D, Li, F, Pan, H, Liang, H, Wu, K, & He, J. (2017) Integrative analysis of genomic sequencing data reveals higher prevalence of LRP1B mutations in lung adenocarcinoma patients with COPD. *Scientific reports* **7**, 1–8.

- [21] Lan, S, Li, H, Liu, Y, Ma, L, Liu, X, Liu, Y, Yan, S, & Cheng, Y. (2019) Somatic mutation of LRP1B is associated with tumor mutational burden in patients with lung cancer. *Lung cancer* **132**, 154–156.
- [22] Grzes, M, Oron, M, Staszczak, Z, Jaiswar, A, Nowak-Niezgoda, M, & Walerych, D. (2020) A driver never works alone—interplay networks of mutant p53, MYC, RAS, and other universal oncogenic drivers in human cancer. *Cancers* **12**, 1532.
- [23] Herbst, R. S, Morgensztern, D, & Boshoff, C. (2018) The biology and management of non-small cell lung cancer. *Nature* **553**, 446–454.
- [24] Gerstung, M, Jolly, C, Leshchiner, I, D'Antonio, S. C, Gonzalez, S, Rosebrock, D, Mitchell, T. J, Rubanova, Y, Anur, P, Yu, K, Tarabichi, M, Deshwar, A, Wintersinger, J, and Ignacio Vázquez-García, K. K, Haase, K, Jerman, L, Sengupta, S, Macintyre, G, Malikić, S, Donmez, N, Livitz, D. G, Cmero, M, Demeulemeester, J, Schumacher, S, Fan, Y, Yao, X, Lee, J, Schlesner, M, Boutros, P. C, Bowtell, D. D, Zhu, H, Getz, G, Imielinski, M, Beroukhi, R, Sahinalp, S. C, Ji, Y, Peifer, M, Markowitz, F, Mustonen, V, Yuan, K, Wang, W, Morris, Q. D, PCAWG Evolution & Heterogeneity Working Group, Spellman, P. T, Wedge, D. C, Looand, P. V, & PCAWG Consortium. (2020) The evolutionary history of 2,658 cancers. *Nature* **578**, 122–128.
- [25] Skoulidis, F, Byers, L. A, Diao, L, Papadimitrakopoulou, V. A, Tong, P, Izzo, J, Behrens, C, Kadara, H, Parra, E. R, Rodriguez Canales, J, Zhang, J, Giri, U, Gudikote, J, Cortez, M. A, Yang, C, Fan, Y, Peyton, M, Girard, L, Coombes, K. R, Toniatti, C, Heffernan, T. P, Choi, M, Frampton, G. M, Miller, V, Weinstein, J. N, Herbst, R. S, Wong, K.-K, Zhang, J, Sharma, P, Mills, G. B, Hong, W. K, Minna, J. D, Allison, J. P, Futreal, A, Wang, J, Wistuba, I. I, & Heymach, J. V. (2015) Co-occurring genomic alterations define major subsets of KRAS-mutant lung adenocarcinoma with distinct biology, immune profiles, and therapeutic vulnerabilities. *Cancer discovery* **5**, 860–877.
- [26] La Fleur, L, Falk-Sörqvist, E, Smeds, P, Berglund, A, Sundström, M, Mattsson, J. S, Brandén, E, Koyi, H, Isaksson, J, Brunnström, H, Nilsson, M, Micke, P, Moens, L, & Botling, J. (2019) Mutation patterns in a population-based non-small cell lung cancer cohort and prognostic impact of concomitant mutations in KRAS and TP53 or STK11. *Lung Cancer* **130**, 50–58.
- [27] Pécuchet, N, Laurent-Puig, P, Mansuet-Lupo, A, Legras, A, Alifano, M, Pallier, K, Didelot, A, Gibault, L, Danel, C, Just, P.-A, Riquet, M, Le Pimpec-Barthes, F, Damotte, D, Fabre, E, & Blons, H. (2017) Different prognostic impact of STK11 mutations in non-squamous non-small-cell lung cancer. *Oncotarget* **8**, 23831.

- [28] Tao, Y, Rajaraman, A, Cui, X, Cui, Z, Chen, H, Zhao, Y, Eaton, J, Kim, H, Ma, J, & Schwartz, R. (2021) Assessing the contribution of tumor mutational phenotypes to cancer progression risk. *PLoS computational biology* **17**, e1008777.
- [29] Youn, A & Simon, R. (2013) Using passenger mutations to estimate the timing of driver mutations and identify mutator alterations. *BMC Bioinformatics* **14**, 1–11.
- [30] Fox, E. J, Prindle, M. J, & Loeb, L. A. (2013) Do mutator mutations fuel tumorigenesis? *Cancer and Metastasis Reviews* **32**, 353–361.
- [31] Hatakeyama, K, Ohshima, K, Nagashima, T, Ohnami, S, Ohnami, S, Serizawa, M, Shimoda, Y, Maruyama, K, Akiyama, Y, Urakami, K, Kusuhashi, M, Mochizuki, T, & Yamaguchi, K. (2018) Molecular profiling and sequential somatic mutation shift in hypermutator tumours harbouring POLE mutations. *Scientific reports* **8**, 1–12.
- [32] Fisk, J. N, Mahal, A. R, Dornburg, A, Gaffney, S. G, Aneja, S, Contessa, J. N, Rimm, D, James, B. Y, & Townsend, J. P. (2022) Premetastatic shifts of endogenous and exogenous mutational processes support consolidative therapy in EGFR-driven lung adenocarcinoma. *Cancer letters* **526**, 346–351.
- [33] Tomlinson, I. P, Novelli, M, & Bodmer, W. (1996) The mutation rate and cancer. *Proceedings of the National Academy of Sciences* **93**, 14800–14803.
- [34] Kadara, H, Choi, M, Zhang, J, Parra, E, Rodriguez-Canales, J, Gaffney, S, Zhao, Z, Behrens, C, Fujimoto, J, Chow, C, Yoo, Y, Kalhor, N, Moran, C, Rimm, D, Swisher, S, Gibbons, D, Heymach, J, Kaftan, E, Townsend, J, Lynch, T, Schlessinger, J, , Lee, J, Lifton, R, Wistuba, I, & Herbst, R. (2017) Whole-exome sequencing and immune profiling of early-stage lung adenocarcinoma with fully annotated clinical follow-up. *Annals of Oncology* **28**, 75–82.