

---

# MODELING EARLY-ONSET CANCER KINETICS TO STUDY CHANGES IN UNDERLYING RISK, DETECTION, AND IMPACT OF POPULATION SCREENING

---

Navid Mohammad Mirzaei<sup>1,\*</sup>, Chin Hur<sup>1,2,3</sup>, Mary Beth Terry<sup>1,3,4</sup>, Piero Dalerba<sup>5,6,7</sup>, and Wan Yang<sup>1,3,\*</sup>

<sup>1</sup>Department of Epidemiology, Mailman School of Public Health, Columbia University, New York, New York, USA

<sup>2</sup>Department of Medicine, Columbia University Irving Medical Center, New York, New York, USA

<sup>3</sup>Herbert Irving Comprehensive Cancer Center (HICCC), Columbia University Irving Medical Center, New York, New York, USA

<sup>4</sup>Silent Spring Institute, Newton, Massachusetts, USA

<sup>5</sup>Center for Discovery and Innovation (CDI), Hackensack Meridian Health (HMH), Nutley, New Jersey, USA

<sup>6</sup>Department of Medical Sciences, Hackensack Meridian School of Medicine (HMSOM), Nutley, New Jersey, USA

<sup>7</sup>Lombardi Comprehensive Cancer Center (LCCC), Georgetown University, Washington, DC, USA

\*Corresponding Authors: [nm3519@cumc.columbia.edu](mailto:nm3519@cumc.columbia.edu) (NM); [wy2202@cumc.columbia.edu](mailto:wy2202@cumc.columbia.edu) (WY)

November 21, 2024

## ABSTRACT

Recent studies have reported increases in early-onset cancer cases (diagnosed under age 50) and call into question whether the increase is related to earlier diagnosis from other medical tests and reflected by decreasing tumor-size-at-diagnosis (apparent effects) or actual increases in underlying cancer risk (true effects), or both. The classic Multi-Stage Clonal Expansion (MSCE) model assumes cancer detection at the emergence of the first malignant cell, although later modifications have included lag-times or stochasticity in detection to more realistically represent tumor detection requiring a certain size threshold. Here, we introduce an approach to explicitly incorporate tumor-size-at-diagnosis in the MSCE framework and account for improvements in cancer detection over time to distinguish between apparent and true increases in early-onset cancer incidence. We demonstrate that our model is structurally identifiable and provides better parameter estimation than the classic model. Applying this model to colorectal, female breast, and thyroid cancers, we examine changes in cancer risk while accounting for detection improvements over time in three representative birth cohorts (1950-1954, 1965-1969, and 1980-1984). Our analyses suggest accelerated carcinogenic events and shorter mean sojourn times in more recent cohorts. We further use this model to examine the screening impact on the incidence of breast and colorectal cancers, both having established screening protocols. Our results align with well-documented differences in screening effects between these two cancers. These findings underscore the importance of accounting for tumor-size-at-diagnosis in cancer modeling and support true increases in early-onset cancer risk in recent years for breast, colorectal, and thyroid cancer.

**Keywords:** Early-onset cancer, Cancer kinetics, Tumor size, Sojourn time, Screening impact

## Significance

This study models recent increases in early-onset cancers, accounting for both true factors contributing to cancer risk and those caused by improved detection. We show that while advancement in detection has led to earlier detection, our model estimates shorter sojourn times and more aggressive carcinogenic events for recent cohorts, suggesting faster tumor progression. Further, a counterfactual analysis using this model reveals the known statistically significant reduction in colorectal cancer incidence (supporting a robust modeling approach), likely due to screening and timely removal of precancerous polyps. Overall, we introduce an enhanced model to detect subtle trends in cancer risk and

demonstrate its ability to provide valuable insights into cancer progression and highlight areas for future refinement and application.

## Introduction

Cancer is the second most common cause of death in the world and is predicted to be responsible for about 2 million new cases and more than 600 thousand deaths in 2024, according to the American Cancer Society Statistics [1]. Moreover, globally, early-onset (less than 50 years of age) cancer incidence has increased by 79.1% with a 27.7% increase in death between 1990 and 2019 [2]. Colorectal, breast, and thyroid cancers are among those with the most dramatic surges in early-onset incidence [3].

Many studies have aimed to understand the reason for such recent increases in early-onset incidence, with dietary factors, tobacco use, and alcohol consumption recognized as the most common modifiable risk factors [3]. Moreover, screening policies have been modified to respond to this issue. For example, the United States Preventive Task Force (USPSTF) updated their guidelines in 2024, recommending that all women begin screening for breast cancer (BrC) at age 40, as opposed to ages 45 and 50, as recommended previously [4]. The USPSTF also recommended the age of colorectal cancer (CRC) screening at 45, lowered from 50, in 2021 [5]. Given that screening facilitates earlier diagnosis, such population-level screening could result in increases in incident cases. In addition, advancements in and increased use of medical imaging technologies including the increased use of CT scan [6], enable the detection of smaller tumors, even for cancer types with no population screening, such as thyroid cancer (ThC) [7, 8], which could, in part, contribute to the apparent increases in incidence among young adults. Indeed, the tumor sizes recorded at diagnosis in general decreased over time since 1988 when the Surveillance, Epidemiology, and End Results (SEER) program started to report these data [9] (Figure 1(A)). Thus, the apparent increases need to be accounted for to gauge the true risk increases more accurately and to examine the underlying etiology of early-onset cancers.

Mathematical models have been applied to explore the underlying factors contributing to cancer incidence. A prime example of such application is Armitage and Doll's pioneering work in 1954, which introduced the theory of multistage carcinogenesis and set the groundwork for similar cancer models [10]. This theory posits that cancer is the product of the accumulation of genetic mutations in normal cells. Building on this, in a series of papers, Moolgavkar et al. (1979, 1981) showed that a two-stage mutation model is not enough to capture the intricacies of cancer incidence unless stage-wise clonal expansion of mutated cells is also considered [11, 12]. This led to the introduction of the Multistage Clonal Expansion (MSCE) model, which incorporates the proliferation of mutated cells at each stage in addition to the accumulation of genetic mutations. In recent years, various studies have modified and adopted the MSCE model to capture different details that contribute to cancer incidence. These include multiple studies examining the number of mutations (i.e., the number of stages in the MSCE model) needed to capture the age-specific incidence patterns for different cancer types [13, 14, 15, 16]; the development of different hazard functions to better capture the phases and transitions of tumor growth kinetics over the life course [17, 18]; additional model components to account for cancer detection [19, 20]; and the incorporation of population- or individual- level risk factors in the mutation rates [21, 22]. Another line of cancer kinetics modeling examines the evolutionary dynamics of cancer initiation and progression; for instance, Paterson et al. (2020) [23] and Li et al. (2023) [24] considered the heterogeneity of mutation orders to explore the genetic pathways that drive these dynamics.

This study focuses on modeling early-onset cancers, particularly three cancer types with reported incidence increases in the United States (US): BrC, CRC, and ThC. Although BrC and CRC diagnostics have improved over time, case detection prior to the recommended screening age for these cancer types is typically due to symptoms. On the other hand, ThC is not recommended for population screening in the US. However, as noted above, the recent increased use of medical imaging and subsequent incidental detection of ThC has been speculated as a factor for its rising incidence [7]. These differences in detection improvements and screening practices allow model testing and comparison of the estimated apparent effects due to increased detection (more prominent for ThC compared to the other two cancer types) and true cancer risks, as well as examining the impact of population screening. To account for such changing detection, we propose a model that captures the clonal expansion of malignant cells via a simple birth process in the MSCE framework while accounting for tumor-size-at-diagnosis. By adding this factor, our model encompasses both the apparent and true effects in the observed early-onset cancer incidence. Unlike previous works that numerically test the fitness of MSCE models with different numbers of stages [13, 14, 15, 16], we adopt a general modeling framework with three stages of mutation and clonal expansion based on findings from recent genetic analyses [25]. To capture the heterogeneity of cancers, we use incidence data to estimate the rate of progression during each stage and the relative contribution of each stage, such that the model can be applied to different cancer types. We apply this model to the three key early-onset cancers (BrC, CRC, and ThC) over multiple birth cohorts. Our results show that our model is fully structurally identifiable, has a better parameter estimation capability, and explains the increase in early-onset cancer better than the classic model. By accounting for changing detection, we are able to estimate the increases in tumor

growth rates and the key stages of such increases in different birth cohorts for the three cancers. In addition, using the model estimates, we calculate the mean sojourn times (the time from the occurrence of the first malignant cell to cancer detection) and conduct a counterfactual analysis to investigate the impact of population screening on BrC and CRC.

## The model: framework, identifiability and validation

**Model framework.** It is a widely accepted theory that cancer arises from the accumulation of genetic mutations that turn normal stem cells into malignant ones [26]. Key mutations occur in genes that regulate cell proliferation and death [27], such as APC, TP53, and KRAS in colorectal cancer; TP53, PIK3CA, and GATA3 in breast cancer; and BRAF, CHEK2, and RET in thyroid cancer [28, 29, 30]. Research by Tomasetti et al. suggests that as few as three driver mutations may be sufficient for cancer to develop, based on studies of colorectal and lung cancers, which have the highest number of somatic mutations (i.e., can serve as an upper bound for other cancer types) [25]. Therefore, in this study, we consider three mutations before malignancy, but note the model provided here can be easily expanded to account for more mutations. Figure 1 (B) shows a schematic of the biological assumption considered for the model. An Ordinary Differential Equation (ODE) system can be derived to describe the changes in survival (i.e., the probability that an individual's tumor size is less than the reported size-at-diagnosis at age  $t$ ) and the associated hazard rate given the three transitions in this MSCE process.

$$\text{Survival probability 1: } \frac{dx_1}{dt} = \mu_0 N_0 x_1 (x_3 - 1), \quad (1)$$

$$\text{Hazard rate of cancer incidence: } \frac{dx_2}{dt} = -\mu_0 N_0 x_4 \quad (2)$$

$$\text{Survival probability 2: } \frac{dx_3}{dt} = \beta_1 - (\alpha_1 + \beta_1 + \mu_1)x_3 + \mu_1 x_3 x_5 + \alpha_1 x_3^2, \quad (3)$$

$$\text{Derivative of survival probability 2: } \frac{dx_4}{dt} = -(\alpha_1 + \beta_1 + \mu_1)x_4 + \mu_1 x_4 x_5 + \mu_1 x_3 x_6 + 2\alpha_1 x_3 x_4, \quad (4)$$

$$\text{Survival probability 3: } \frac{dx_5}{dt} = \beta_2 - (\alpha_2 + \beta_2 + \mu_2)x_5 + \mu_2 f(t)x_5 + \alpha_2 x_5^2, \quad (5)$$

$$\text{Derivative of survival probability 3: } \frac{dx_6}{dt} = -(\alpha_2 + \beta_2 + \mu_2)x_6 + \mu_2 f'(t)x_5 + \mu_2 f(t)x_6 + 2\alpha_2 x_6 x_5. \quad (6)$$

The value  $N_0$  denotes the number of stem cells. For breast we set  $N_0 = 1.74 \times 10^{10}$  [31], for colon and rectum  $N_0 = 2 \times 10^8$  [32], and for thyroid  $N_0 = 6.5 \times 10^7$  [32]. Table 1 gives all the variable and parameter descriptions. The derivation process entails tracking how stem cells transition from normal to malignant as they accumulate genetic mutations. The model considers a birth-death process for each stage, where cells can divide, die, or mutate at specific rates. By incorporating the probabilities of these transitional events, we estimate the likelihood that a certain number of malignant cells will develop and form a detectable tumor by a given age. Refer to the supplementary materials for more details on the model and derivation. When  $f(t) = 0$ , equations (1)-(6) give the classic MSCE model as first introduced in [12] (referred to as the MSCE model hereafter), which counts cancer incidence at the occurrence of the first malignant cell. Here, we introduce a model that counts a tumor as a case only if it surpasses the reported size-at-diagnosis in the data. This is an attempt to segregate the apparent effects (caused by improved diagnosis over time) and true effects (caused by mutagenic factors) contributing to the recent increases in early-onset cancer. To include the tumor size in the model, we derive the probability of having less than  $N_t$  malignant cells given one malignant cell at  $t = 0$  through a linear birth-death process (See Bailey chapter 8 [33]):

$$f(t) = 1 - \frac{(\alpha_3 - \beta_3)\alpha_3^{N_t-1} \{1 - e^{(\beta_3 - \alpha_3)t}\}^{N_t-1}}{\{\alpha_3 - \beta_3 e^{(\beta_3 - \alpha_3)t}\}^{N_t}}. \quad (7)$$

Briefly, to derive equation (7), we first solve the moment-generating partial differential equation of the linear homogeneous birth-death process. For the solution, we then obtain the probability generating function and use that function to generate a general formula for the probability of having exactly  $n$  malignant cells. Using an infinite sum over  $n$  and the geometric series formula, we arrive at equation (7); see the details in the supplementary material. Combining this function with (1)-(6) gives a general model accounting for tumor-size-at-diagnosis (referred to as the General MSCE-T model in this paper).

Table 1: Models, model variables and parameters definition.

Variable/Parameter/Model	Description
$x_1$	Survival* given $N_0$ normal cells and no other cell types at time zero.
$x_2$	Hazard rate** of cancer incidence, i.e., $-\frac{d}{dt}[\ln x_1]$
$x_3$	Survival given $N_0$ normal cells, one FSMC <sup>†</sup> , and no other cell types at time zero.
$x_4$	Auxiliary variable defined as $\frac{d}{dt}[x_3]$
$x_5$	Survival given $N_0$ normal cells, one SSMC <sup>‡</sup> , and no other cell types at time zero.
$x_6$	Auxiliary variable defined as $\frac{d}{dt}[x_5]$
$\mu_0$	Mutation rate of normal/stem cells
$\mu_1$	Mutation rate of FSMCs
$\mu_2$	Mutation rate of SSMCs
$\alpha_1$	Birth rate of FSMCs
$\alpha_2$	Birth rate of SSMCs
$\alpha_3$	Birth rate of malignant cells used in Models 2 and 3
$\beta_1$	Death rate of FSMCs
$\beta_2$	Death rate of SSMCs
$\beta_3$	Death rate of malignant cells used only in Models 2
$N_t$	Number of malignant cells at diagnosis time $t$
$\mu_1 \times (\alpha_1 - \beta_1)$	Aggressiveness <sup>◊</sup> of FSMCs
$\mu_2 \times (\alpha_2 - \beta_2)$	Aggressiveness of SSMCs
MSCE model	Equations (1)-(6) with $f(t) = 0$ , $[x_1(0), x_2(0), x_3(0), x_4(0), x_5(0), x_6(0)] = [1, 0, 1, 0, 1, -\mu_2]$ , cancer incidence := the occurrence of the first malignant cell, No $\alpha_3$ and $\beta_3$ .
General MSCE-T model	Equations (1)-(6) with $f(t)$ given by (7), $[x_1(0), x_2(0), x_3(0), x_4(0), x_5(0), x_6(0)] = [1, 0, 1, 0, 1, 0]$ , cancer incidence := having $N_t$ or more malignant cells, Both $\alpha_3$ and $\beta_3$ require estimation.
MSCE-T model	Equations (1)-(6) with $f(t)$ given by (8), $[x_1(0), x_2(0), x_3(0), x_4(0), x_5(0), x_6(0)] = [1, 0, 1, 0, 1, 0]$ , cancer incidence := having $N_t$ or more malignant cells, $\alpha_3$ is known and $\beta_3$ is not required.

\* Survival: The probability of having less than  $N_t$  malignant cells. For the MSCE model,  $N_t = 1$ .

\*\* Hazard rate: the instantaneous rate at which a specific event occurs.

† FSMC: First-Stage Mutated Cell, ‡ SSMC: Second-Stage Mutated Cell,

◊ Aggressiveness: Likelihood of newly proliferated cells to mutate to the next stage.

A well-accepted simple model of cancer growth is the logistic growth with the rate of  $\alpha N(t)(1 - \frac{N(t)}{K})$  [34]. Parameters  $\alpha$  and  $K$  are net proliferation rate and carrying capacity, respectively, and  $N(t)$  is the number of cancer cells at a given time  $t$ . If the carrying capacity is much larger than  $N(t)$ , then it can be approximated by a simple birth process  $\alpha N(t)$ . In this study, as tumor sizes were recorded at the time of diagnosis, the sizes would be much smaller than the carrying capacity. For example, for CRC, the largest size at diagnosis reported in SEER data is around 52 mm, while tumors as big as 17 cm are possible [35]. Hence, by setting  $\beta_3 = 0$  in equation (7), we obtain the probability of having less than  $N_t$  malignant cells assuming a simple birth process:

$$f(t) = 1 - \{1 - e^{-\alpha_3 t}\}^{N_t - 1}. \quad (8)$$

Combining function (8) with equations (1)-(6) leads to a simplified model to account for tumor-size-at-diagnosis and the main model in this study, which we refer to as the MSCE-T model. The MSCE-T model can be further simplified using the reported doubling times for different cancer types in the literature. The average doubling time is 193 days for BrC [36], 211 days for CRC [37, 38], and 967 days for ThC [39]. This allows us to set  $\alpha_3 = \frac{365 \times \ln 2}{\text{doubling time}}$  in (8). This simplification is not feasible for the General MSCE-T model, as  $\alpha_3 - \beta_3$  (the net proliferation rate) is not uniquely determined, and the term  $\beta_3$  appears with an exponential factor in the denominator. Thus, we can instead use the reported doubling times to establish bounds for estimating  $\alpha_3$  and  $\beta_3$  in the General MSCE-T model.

Finally, to solve (1)- (6) along with (7) or (8) we consider the initial conditions  $[x_1(0), x_2(0), x_3(0), x_4(0), x_5(0), x_6(0)] = [1, 0, 1, 0, 1, 0]$ . However, solving the classic MSCE model (i.e.,  $f(t) = 0$ ) will entail a different set of initial conditions namely  $[x_1(0), x_2(0), x_3(0), x_4(0), x_5(0), x_6(0)] = [1, 0, 1, 0, 1, -\mu_2]$ . The initial values for  $x_1, x_3$ , and  $x_5$  are one because they represent the survival (i.e., the probability of having less than  $N_t$  malignant cells) at time zero given different initial scenarios for mutated cells. The initial values for  $x_2, x_4$  and  $x_6$  are a consequence of their definition (see Table 1) which can be directly determined from equations (1), (3) and (5). See Table 1 for a summary of each model.

**Model identifiability.** A model is structurally identifiable if model parameters can be determined uniquely given perfect (noise-free and error-free) data, indicating that there is a unique mechanism as represented by the estimated parameter set that can explain the observation. Conversely, a model is termed structurally unidentifiable if several parameter sets yield the same data, indicating unreliable parameter estimation. As such, structural identifiability is important for parameter estimation.

The MSCE model, despite its simplicity, is not structurally identifiable. Using a differential algebraic approach, Brouwer et al. prove that for the MSCE model, the system is unidentifiable and the parameters need to be combined into groups to retrieve full structural identifiability [40].

The General MSCE-T model does not resolve the identifiability issue. Moreover, due to highly nonlinear terms, including two extra unknown parameters  $\alpha_3$  and  $\beta_3$ , recovering identifiable groups like that of the classic MSCE model is difficult. Please refer to the supplementary material for a justification.

For the MSCE-T model, when  $\alpha_3$  in equation (8) is known (based on values reported in the literature), and  $N_t$  is obtained from the data, equation (8) can be treated as a known input to the model. Introducing additional inputs is a recognized approach to mitigate structural unidentifiability, and incorporating (8) instead of (7) makes the model fully structurally identifiable. A detailed proof is provided in the supplementary material.

**Model validation using synthetic data.** We validate the model and parameter estimation approach using synthetic data. This synthetic dataset is generated through Poisson sampling, using the model-derived incidence as the mean, to simulate the imperfect observations in real-world settings (see the "Materials and methods" section for more details). As the parameters are prescribed, we can compare the model estimates with the true parameter values (typically unknown for real-world data) in addition to model fit. Here, we assess i) the model fit to the data, ii) the accuracy of parameter estimation from the best-fit model run, and iii) the consistency based on the ensemble from 100 runs. As shown in Figure 2, the MSCE-T model performs best in all three aspects – It generates the narrowest incidence output range from the parameter ensemble, its estimated parameters have the lowest error compared to the true values, and it has high consistency as shown by tighter parameter distributions. These results are consistent with the structural identifiability of the MSCE-T model.

## Data analysis results

**Parameter estimates accounting for tumor size.** We next apply the MSCE-T model to examine the tumor growth kinetics of three early-onset cancers – BrC (under age 40), CRC (under age 50), and ThC (all ages included), based on cancer incidence across three representative cohorts (1950-1954, 1965-1969, and 1980-1984) in the US. All three cancer types saw increases in incidence among young adults in more recent cohorts (Figure 3; darker colors for more recent cohorts). The MSCE-T model is able to capture the observed incidence trends for all three cohorts (Figure 3). For ease of interpretation, we group the parameters into three sets representing three key stages of carcinogenesis considered in the model: the mutation rate of normal stem cells (i.e., initial mutation), the first-stage mutated cell (FSMC) aggressiveness (combining the stage-specific mutation-, birth-, and death rate; Table 1), and the second-stage mutated cell (SSMC) aggressiveness. However, we note these parameter sets are computed directly using individual parameter estimates from the MSCE-T model given it is fully identifiable, rather than through reparametrization using these groupings. Across the three cancer types, the estimated mutation rate of normal stem cells is higher in CRC than BrC and ThC, while estimated aggressiveness for both the first-stage and second-stage mutated cells is the highest in ThC.

Importantly, for all three cancer types, sensitivity analysis indicates the parameter combination  $\mu_1 \times (\alpha_1 - \beta_1)$ , which measures the aggressiveness of FSMCs, incurs the highest sensitivity such that per-unit change in this parameter set leads to the largest change in cancer incidence (Figure 4), sensitivity plots, red curves). The MSCE-T model estimates an increase in this parameter set in more recent birth cohorts (Figure 4, second column of box plots, higher red bars in later cohorts) for all three early-onset cancers. For female BrC, compared to women born in 1950-1954, the estimated FSMC aggressiveness increased by 10.4% (95% CI: 1-21; 1965-1969 cohort) and 23.3% (95% CI: 11-35; 1980-1984 cohort) in a span of 30 years. In addition, the model estimates higher values for this parameter in females than males for

both CRC and ThC (Figure 4, see the second column of box plots, higher darker red bars for females). Particularly, for CRC, compared to the 1950-1954 birth cohort, estimated FSMC aggressiveness increased by 3.9% (95% CI: 0.1-9; female) and 7.8% (95% CI: 5-10; male) in the 1965-1969 birth cohort and by 22.2% (95% CI: 16-29; female) and 10.4% (95% CI: 7-14; male) in the 1980-1984 birth cohort. For ThC, the estimated increases were even higher, by 80.7% (95% CI: 20-147; female) and 22.2% (95% CI: 12-32; male) in the 1965-1969 birth cohort and by 111.9% (95% CI: 32-191; female) and 180.9% (95% CI: 158-206; male) in the 1980-1984 birth cohort. Estimates for the other two parameter sets are less consistent across cohorts and sexes, and as noted, the model is less sensitive to these parameter sets.

**Sensitivity and supplemental analyses.** To assess the robustness of our model estimates, we conducted a sensitivity analysis by varying the values of  $\alpha_3$  for all three cancer types, given that this parameter is derived from the literature. Figure S2 compares parameter estimation results for both increased and decreased values of  $\alpha_3$ . Despite the variation in parameter values, the qualitative trends remain consistent with the results reported above (and shown in Figure 4) across all parameter sets and cancer types. Particularly, the most sensitive parameter (i.e., the aggressiveness of FSMC) consistently increases by birth cohort for different values of  $\alpha_3$ .

As noted in the Introduction, improvement in cancer detection and incidental detection due to increased use of medical imaging could, in part, contribute to the apparent increases in cancer incidence. Such biases in the data are more profound for ThC (see the large increases in Figure 3)[7]. To examine the ability of the MSCE-T model to account for such detection-related data biases, we conducted two additional analyses. First, we compare parameter estimates using the MSCE model, which does not account for tumor size, with estimates from the MSCE-T model. As shown in Figure S3, estimates by the MSCE model do not show clear changes by birth cohort for the three cancer types examined. In the second analysis, we performed parameter estimation for ThC by subtype. About 90% of ThC cases are papillary carcinomas, 1% are anaplastic carcinomas, and the remaining cases are follicular, Hürthle, and medullary carcinomas [41]. Despite its scarcity, anaplastic ThC accounts for over 30% of all ThC-related deaths [41]. Thus, we conducted the analysis for papillary (most prevalent) and anaplastic (most lethal) ThC separately. Given the low incidence of anaplastic ThC, for this analysis, we aggregate the data for 15-year cohorts (i.e., 1940-1954, 1955-1969, and 1970-1984) to reduce observational noises. However, we note the mean incidence of anaplastic ThC is still low (Figure S4) and caution the greater uncertainty in model estimates for this subtype. As shown in Figure S4 (A), the MSCE-T model estimates increases in the FSMC aggressiveness for papillary ThC over the three birth cohorts, similar to the main analysis combining all ThC subtypes. In contrast, for anaplastic ThC, there are no clear changes in incidence over the three cohorts for both sexes, and the model does not estimate an increase in FSMC aggressiveness (Figure S4 (B)). Together, these results demonstrate the ability of the MSCE-T model to account for changing detection and identify changes in tumor growth kinetics and suggest there are increases in the risk of the three early-onset cancers (for ThC, such increases are mostly related to papillary ThC) in more recent cohorts, independent of changing detection.

**Sojourn time.** We use the MSCE-T model to calculate the mean sojourn time (the time from the emergence of first malignant cell to cancer detection) of BrC, CRC, and ThC for the three cohorts (1950-1954, 1965-1969, and 1980-1984). An illustration of the sojourn time based on the MSCE-T model incidence curves is provided in Figure S5. The mean sojourn times (Table 2) estimated here are comparable to results from the Cancer Intervention and Surveillance Modeling Network (CISNET), for example, 2-4 years for breast cancer [42] and 10.6 years (with an interquartile range of 5-14 years) from adenoma incidence to cancer diagnosis for CRC [43]. Here, we provide more detailed estimates by sex and birth cohort. The estimates are similar for females and males of the same cohort. However, we notice a decreasing trend for more recent cohorts for all cancer types. To assess the influence of the likely improved detection over time on the estimates, we fixed  $N_i$  as the average across all cohorts to compute the sojourn times controlling for changes in tumor-size-at-diagnosis. The declining trend in sojourn time persisted but with a slight reduction (though more noticeable for ThC) in each value. Together, these results suggest increases in tumor aggressiveness (hence the shorter sojourn times), consistent with the model parameter estimates reported above.

**Impact of population screening on cancer incidence.** Given the mechanistic design, the MSCE-T model also affords prediction to help assess the impact of population screening on cancer incidence. Particularly, for BrC and CRC, the parameters are estimated using incidence data before the recommended screening ages. These estimates thus represent tumor growth kinetics without screening, and the model projections would represent incidence under a counterfactual scenario where there is no population screening. We thus use these parameter estimates to project the incidence of BrC for ages above 40 (Figures 5 (A)) and CRC for ages above 50 (Figures 5 (B) & (C)) to assess the impact of screening for these two cancer types. For this analysis, we examine two cohorts with high rates of screening (>~60%) for each cancer type (i.e., the 1950-1954 and 1955-1959 cohorts for BrC and the 1955-1959 and 1960-1964 cohorts for CRC), and an earlier cohort (1945-1949) with limited screening for comparison (for more details, refer to the "Materials and methods" section).

For BrC, the model projected incidence rates better match with the observed values for the comparison cohort born in 1945-1949 (with less than 30% above the age of 40 screened [44]; Figure 5 (A) left). For the two main screening-

Table 2: Mean sojourn time for different cancer types and cohorts

Cancer Type	Sex	Cohort (Years)	Mean sojourn time with the actual $N_t$ (Years)	Mean sojourn time with fixed $N_t^*$ (Years)
BrC	Female	1950-1954	3.9	3.8
		1965-1969	3.8	3.8
		1980-1984	3.6	3.5
CRC	Female	1950-1954	10.7	10.5
		1965-1969	10.1	10
		1980-1984	8.8	8.7
CRC	Male	1950-1954	10.7	10.5
		1965-1969	10.1	10
		1980-1984	8.8	8.8
ThC	Female	1950-1954	17.1	16.3
		1965-1969	14.7	14.4
		1980-1984	12.3	12.2
ThC	Male	1950-1954	17	16.4
		1965-1969	14.7	14.5
		1980-1984	12.2	12

\*  $N_t$  is fixed as the average tumor cell counts at diagnosis across all cohorts.

affected cohorts (1950-1954 and 1955-1959), the model projections align closely with the observed incidence until approximately ages 50-55 (Figures 5 (A) middle and right) – around this age range, there was a brief transition phase (i.e., slow-down) in incidence rates. After age 50-55, the observed incidence rates start to more substantially surpass the model projections (see, e.g., Figure 5 (A) middle vs left; dots above the lines). The larger discrepancy between the model projection and the data for cohorts with higher rates of screening likely reflects the increased earlier detection of BrC through screening. However, as similar discrepancy is also seen for those older than 50-55 years of age in the comparison cohort with limited screening, we suspect the higher-than-projected incidence is in part due to a true risk increase in those ages.

For CRC, the model projected incidence rates in general align with the observed values for the comparison cohort of 1945-1949 (with less than 34% above the age of 50 screened [45]; Figures 5 (B) & (C) left). For the cohorts of 1955-1959 and 1960-1964, among those aged 50 and above, the projected incidence (without screening) is much higher than the observed (with screening) for both females and males (Figures 5 (B) & (C) middle and right). This discrepancy likely reflects the preventive impact of CRC screening using colonoscopy, which not only detects but also facilitates the removal of precancerous polyps, thereby preventing malignancies from progressing to clinically observable cases [46].

Overall, the observed differences in the contribution of screening for BrC and CRC might be attributed to the nature of the procedure. The removal of breast precancerous lesions is more complex compared to CRC, often requiring lumpectomy followed by radiation therapy as a preventive measure [38]. This may explain the absence of a sudden decline in incidence following the screening age, as observed for CRC. Additionally, BrC screening is recommended biennially [41], potentially contributing to the delayed increase in incidence due to detection (after the transition phase). However, the discrepancy between the model projection and the less screened population (cohort 1945-1949) might indicate additional underlying true risks.

## Discussion

This study introduces an extension to the classic multistage carcinogenesis model by including tumor-size-at-diagnosis data. Our findings suggest that adding this input significantly improves the model's sensitivity and ability to capture important trends in cancer progression. The main MSCE-T model outperformed both the classic model and a more general model (incorporating the birth and death of malignant cells) in terms of structural identifiability and parameter estimation.

To gauge the accuracy of our model estimates, we compare the mean sojourn times estimated using our model and those reported in the literature. Luebeck et al. (2013) used a similar routine and definition for sojourn time as in this study (i.e., from the occurrence of the first malignant cell to cancer diagnosis). However, their estimates for CRC – 5.2 years (CI: 3.6–6.2) for male CRC and 6.5 years (CI: 5.2–7.6) for female CRC [20] – are much shorter than our estimates (8-10 years in Table 2). Two main methodological differences may have contributed to this discrepancy. First, Luebeck et al. (2013) treated detection as a stochastic event, whereas our model explicitly tracks tumor growth from

the first malignant cell to diagnosis based on reported tumor size data. Second, Luebeck et al. (2013) adjusted their incidence data for cohort and period effects, while we used the cohort-specific data directly. The Cancer Intervention and Surveillance Modeling Network (CISNET) MISCAN-Colon project defines sojourn time as the time from preclinical cancer to diagnosis, reporting 3 years (CI: 1-4) for this transition and 10.6 years (CI: 5-14) for the transition from adenoma to diagnosis [43]. The latter transition time is more in line with our estimate for CRC. For BrC, our estimates match the estimates of the CISNET-DFCI model (using data from various sources in the pre-screening era), which suggests a transition time of 2-4 years [42]. However, when comparing to these studies, extra caution is required since CISNET models either calculate the mean sojourn time using a simulated population (with assumptions such as transition probabilities, tumor growth patterns, and detection sensitivity) [43] or consider it as an input estimated through transition probabilities with the possibility of regression [42]. Nonetheless, all these sojourn time estimates are in the order of years (e.g., 5-10 years for CRC), which is substantial compared to the human lifespan and not as negligible as assumed in many previous models.

For ThC, we estimated much longer sojourn times (~15 years; Table 2), despite the earlier lifetime occurrence of ThC than the other two cancer types (i.e., BrC and CRC; see age-specific incidence in Figure 3). Given the elevated incidence rates starting around ages 15-20 (Figure 3), these estimates suggest that the initiation of ThC might have started at a very early stage of life. We are unable to locate a similar model estimation of sojourn time for ThC for comparison (it is less modeled). However, extensive clinical studies have also pointed to the likely initiation of ThC during the infantile period [47]. Particularly, studies of children and adolescents exposed to the Chernobyl nuclear accident reported the highest incidence rate of ThC among those who were under 1 year of age at the time of Chernobyl, and the rates decreased progressively through age 12 [48]. Large population-scale health surveys of Fukushima residents, conducted to monitor the impacts of the TEPCO-Fukushima Daiichi Nuclear Power Plant accident, also revealed high baseline (i.e., not associated with the accident) incidence rates of ThC among children, particularly those aged 15 and older [49, 50]. Further, detailed autopsy data showed the prevalence of ThC (i.e., identified through autopsies of people who died from mostly non-ThC-related causes) increases steeply from age 15 to 34 and then stays roughly constant for the remaining lifetime (see Fig 2 of Takano 2017 [47]). This is possible as ThC is self-limiting (i.e., malignant yet causing very low mortality; [47]). The consistency of our model estimates with these independent clinical observations indicates our model is able to accurately identify the underlying cancer kinetics.

The changes in the kinetics of early-onset cancers are particularly of interest, given the dramatic increases in incidence during recent years in the US and globally [2]. After controlling for potential data biases due to changing detection, estimated carcinogenic aggressiveness still increased substantially for the three studied cancers, i.e., BrC, CRC, and ThC, which indicates a genuine rise in the underlying cancer risk. Further, we identify the increases in the aggressiveness of the first-stage mutated cells (i.e., FSMC aggressiveness) to be the most impactful intermediate step affecting cancer risk. Comparing three cohorts born over a span of 30 years (1950-1954 to 1980-1984), estimated FSMC aggressiveness shows higher values for female cancer than male cancer and has increased significantly over the 30-year study period for both sexes (e.g., by ~23% for female BrC or by ~10% for male CRC; see Figure 4). Given that our model considers only three rate-limiting mutation stages before malignancy, we speculate that the estimated increased carcinogenic effects (i.e., FSMC aggressiveness) correspond to early and mid-adulthood events. However, a separate study is required regarding the timing of driver gene mutations, potentially leveraging the current model and incidence data in conjunction with gene expression data to more precisely estimate the timing of intermediate mutations.

Early detection and treatment through screening have been a key cancer intervention strategy. In the US, population screening is recommended for both BrC (age >40) and CRC (age >50, or >45 from 2021 onward [5]). As clinical trials are difficult to conduct in younger ages when cancer is rare, modeling often serves as a means to evaluate the effectiveness of screening as well as etiology [51]. In particular, several large-scale CISNET models have been developed and applied to test different screening policies and their effectiveness for both cancer types [52, 53, 54, 55]. Here, our model provides an alternative approach to assess the effectiveness. Using counterfactual modeling, we show striking reductions in CRC incidence among those aged >50 (Figure 5 (B) & (C) middle and right), likely thanks to the implemented screening programs and the feasibility of removal of precancerous polyps detected during screening [46]. For BrC, our modeling analysis indicates the absence of an immediate decline in incidence compared to CRC (likely related to the complexity of breast lesion removal [38]) and a delayed increase in the incident cases (possibly due to the biennial screening recommendation [41]); see Figures 5 (A) middle and right. Importantly, unlike the CISNET models assessing screening impact primarily at the population level without explicitly accounting for cohort-specific trends [56, 57, 43, 58], our model accounts for the apparent incidence changes due to changing detection and constructs cohort-specific counterfactuals, based on cohort- and age-specific incidence and tumor-size-at-diagnosis data that are readily available from the SEER program. This simplicity and specificity of our model thus affords a powerful alternative to more explicitly assess the screening impacts, particularly in the context of evolving early-onset cancer risk.



We recognize several study limitations. First, based on findings from Tomasetti et al. [25] and for generality, we adopt a three-stage model for all three cancers studied here. While the model estimates are consistent with the literature as noted in the Results, we cannot determine whether three rate-limiting mutations indeed apply to the three cancers (i.e., BrC, CRC, and ThC). Second, while we are able to estimate the changes in cancer kinetics, the current model does not consider the underlying causes for the estimated changes. Future work can incorporate risk factor data (e.g., as done in [21]) to help identify the main underlying causes. Third, for simplicity, we assume constant values for the parameters in this study. Time-dependent parameters (e.g., based on risk factor data) would likely further improve model performance and provide more detailed estimates, particularly the likely causes of the increased cancer risks during different intermediate stages of tumor growth. Fourth, the tumor mass is likely highly complex, encompassing various cell types. Here, for simplicity, we did not consider such heterogeneity when converting the tumor size to the number of malignant cells (i.e.,  $N_t$  in the model). Nonetheless, we note this simplification would not affect the estimated changes by birth cohort, as shown in our analysis of mean sojourn time setting  $N_t$  to a fixed number for all cohorts; Table 2. Finally, the model presented here is deterministic, with inputs and outputs that describe the overall population in each cohort. For simplicity, we did not include stochasticity reflecting the individual-level heterogeneity in data, which may, in turn, underestimate the parameter uncertainties. We intend to improve this model and its applications by addressing these limitations in our future work.

In conclusion, this study introduces an improvement in cancer modeling by integrating tumor-size-at-diagnosis data into the multistage carcinogenesis framework. This approach enhances the model's ability to detect subtle trends in cancer risk, controlling for observational biases due to changing detections. While the study has some limitations, its contributions offer valuable insights into cancer progression and highlight areas for future refinement and application.

## Materials and methods

**Data sources and processing.** We use the SEER cancer incidence registries for patients diagnosed during 1973-2015 with BrC, CRC, and ThC. We use the International Classification of Diseases for Oncology (ICD-O) codes to filter out the three cancer types: 1) C50.0-C50.9 for BrC; 2) C180.0-C180.9, C19.9, and C20.9 for CRC to include the colon, rectosigmoid junction, and rectum; and 3) C73.9 for ThC. We divide the data into 5-year birth cohorts and stratify them by sex.

The SEER dataset reports the tumor-size-at-diagnosis for 1988-2015, while the year of diagnosis goes as far back as 1973. The size data are reported under 10-digit EOD (1988-2003) or CS tumor size (2004-2015) and describe the largest dimension, or the diameter of the primary tumor, at the time of diagnosis. Using linear extrapolation, we estimate tumor size at detection for years without SEER data (1973-1987); see Figure S1 in the supplementary material. Figure 1 (A) shows the average size data extracted from SEER registry. A decreasing trend is evident for all three cancer types. We use a formula based on Kepler's conjecture to compute the number of cells from the tumor size (see details in the supplementary material).

**Model validation using synthetic data.** We compare the parameter estimation for each of the three models. Models MSCE and MSCE-T each have seven parameters, and the General MSCE-T model, with the unknown  $\alpha_3$  and  $\beta_3$ , has nine parameters. To make the comparison fair, we add a constraint for this model's parameter estimation, forcing the algorithm to find  $\alpha_3$  and  $\beta_3$  values such that  $\alpha_3 - \beta_3$  is equal to the pure birth rate of the MSCE-T model. We generate synthetic data for the incidence curves using the models and an arbitrary parameter set - as the true parameter values are known here, we can assess the accuracy of model estimates directly. To mimic noise in observations, we use Poisson random sampling with the mean set to the model simulated incidence. We use Hybrid Genetic Algorithm (HGA) optimization for parameter estimation [59], employing MATLAB's HGA from the Global Optimization Toolbox. We run the algorithm 100 times and record the fittest set of parameters and their fitness value in each iteration. The fitness value is the least square distance of the model output from the data (cost). We consider the fittest (i.e., lowest least square distance) of the 100 parameter sets as the best-fit parameter set. For illustration purposes, we calculate and plot the distribution of the relative difference percentage between the best-fit parameter set and the true value, see Figure 2.

**Parameter estimation for BrC, CRC, and ThC incidence data.** We estimate parameters to identify trends in cancer progression across different cohorts and cancer types. In this study, we estimated parameters for BrC, CRC, and ThC based on the data for cohorts born in 1950-1954, 1965-1969, and 1980-1984. These cohorts were chosen to represent distinct historical periods, allowing us to capture temporal trends in cancer biology and treatment advancements. Additionally, these cohorts provide a significant number of early-onset incident cases. To minimize bias introduced by cancer screening for BrC and CRC, we restrict our analysis to incident cases diagnosed before age 40 and 50, respectively. As in the synthetic testing, parameter estimation for each cohort is carried out 100 times using MATLAB's HGA toolbox. We acquire a distribution of the parameter values from 100 iterations. We summarize the results via three biologically meaningful groups containing all the model-estimated parameters. These groups represent the

three key stages of carcinogenesis considered in the model: the initial mutation rate of normal stem cells  $\mu_0$  and the aggressiveness of the 1st and 2nd stage mutated cells,  $\mu_1 \times (\alpha_1 - \beta_1)$  and  $\mu_2 \times (\alpha_2 - \beta_2)$ , respectively (Table 1).

**Sojourn Time.** To calculate the mean sojourn time, first, we acquire an estimation for parameter values by fitting the MSCE-T model to the incidence of CRC (under age 50), BrC (under age 40), and ThC (all ages) for cohorts 1950-1954, 1965-1969, and 1980-1984, as explained in the previous section. Note that this fitting is done considering the tumor cell count at diagnosis ( $N_t$ ). Using the same parameters in the model but setting  $N_t = 1$  will result in cancer incidence being recorded at first malignancy. Hence, to find the sojourn time, we calculate the difference between the time the MSCE-T model with varying  $N_t > 1$  produces the same incidence as the MSCE-T model with  $N_t = 1$ . We obtain the mean sojourn time by averaging the sojourn times restricted to cases for which relevant clinical data are available. We repeat the procedure by fixing  $N_t$  (average over all cohorts) to explore the effect of detection improvement.

**Impact of population screening on BrC and CRC incidence: a counterfactual analysis.** As noted above, we estimate model parameters for ages under 40 for BrC and 50 for CRC, who were not subject to population screening for these cancers during the study period. To examine the impact of population BrC and CRC screening, we then use these parameter estimates and the MSCE-T model to project cancer incidence for those older than the recommended screening ages. That is, as the parameters do not include the impact of screening, these projections represent cancer incidence under a counterfactual scenario with no screening. For female BrC, according to CDC National Center for Health Statistics, mammogram screening for women over the age of 40 reached 59.7% by 1993 and increased to over 70% by 2000 [44]. Therefore, for analysis of BrC screening, we pick cohorts of 1950-1954 (reaching age 40 in 1990-1994) and 1955-1959 (reaching age 40 in 1995-1999). For CRC, about 59-62% of adults over the age of 50 underwent some colorectal screening procedures between 2005 and 2010 [45]. Therefore, for analysis of CRC screening, we pick cohorts of 1955-1959 (reaching age 50 in 2005-2009) and 1960-1964 (reaching age 50 in 2010-2015). These cohorts contain a significant number of data points before and after the screening age. Moreover, for both cancer types, we include the cohort of 1945-1949 (less than 30% of the population for BrC [44] and 34% for CRC [45] were screened) to compare the data trends when screening was much more scarce.

**Author Contributions.** NM and WY conceived the study, performed the analysis, and wrote the first draft. CH, MBT, and PD contributed to study design and result interpretation. All authors contributed to the final draft.

**Funding information.** This work was supported by the National Institutes of Health (R01CA257971).

**Acknowledgement.** We would like to thank The Mailman School of Public Health and The Center for Computational Biology and Bioinformatics (C2B2) at Columbia University for access to high-performance computing resources.

**Competing interest statement.** Piero Dalerba is listed as a co-inventor on patents owned by the University of Michigan (US-07723112), Stanford University (US-09329170, US-09850483, US-10344094, US-11130813) and Columbia University (US-12115140-B2), and related to: 1) the discovery of surface markers for the differential purification of cancer stem cell populations from human malignancies; 2) the use of single-cell genomics technologies for the identification of pharmacological targets expressed in cancer stem cell populations; 3) the combination of anti-CD47 and anti-EGFR monoclonal antibodies for the treatment of human colon cancer; and 4) the use of inverse agonists of RAR/RXR signaling as anti-tumor agents against adenoid cystic carcinomas. Some of the patents listed above were licensed to pharmaceutical companies, resulting in the award of royalties and/or stock from Oncomed Pharmaceuticals, Quantice Pharmaceuticals and Forty Seven Inc. (Gilead). Piero Dalerba recently owned stock of Eli Lilly and Company. Piero Dalerba's spouse is employed by Regeneron Pharmaceuticals Inc., and owns (or recently owned) stock of the following pharmaceutical companies: AbbVie, Amgen, AstraZeneca, Eli Lilly and Company, Gilead Sciences Inc., GlaxoSmithKline (GSK), Johnson & Johnson, Merck & Co., Novartis, Organon & Co., Pfizer, Teva Pharmaceutical Industries Ltd and Viatrix. Chin Hur has served as a consultant for Guardant Health. Other authors declare no potential conflicts of interest.

## References

- [1] Rebecca L Siegel, Angela N Giaquinto, and Ahmedin Jemal. Cancer statistics, 2024. *CA: a cancer journal for clinicians*, 74(1):12–49, 2024.
- [2] Jianhui Zhao, Liying Xu, Jing Sun, Mingyang Song, Lijuan Wang, Shuai Yuan, Yingshuang Zhu, Zhengwei Wan, Susanna C Larsson, Konstantinos K Tsilidis, et al. The global trends in incidence, death, burden and risk factors of early-onset cancer from 1990 to 2019. *BMJ oncology*, 2:e000049, 2023.
- [3] Tomotaka Ugai, Naoko Sasamoto, Hwa-Young Lee, Mariko Ando, Mingyang Song, Rulla M Tamimi, Ichiro Kawachi, Peter T Campbell, Edward L Giovannucci, Elisabete Weiderpass, et al. Is early-onset cancer an emerging

- global epidemic? current evidence and future implications. *Nature Reviews Clinical Oncology*, 19(10):656–673, 2022.
- [4] Wanda K Nicholson, Michael Silverstein, John B Wong, Michael J Barry, David Chelmow, Tumaini Rucker Coker, Esa M Davis, Carlos Roberto Jaén, Marie Krousel-Wood, Sei Lee, et al. Screening for breast cancer: Us preventive services task force recommendation statement. *JAMA*, 2024.
- [5] Karina W Davidson, Michael J Barry, Carol M Mangione, Michael Cabana, Aaron B Caughey, Esa M Davis, Katrina E Donahue, Chyke A Doubeni, Alex H Krist, Martha Kubik, et al. Screening for colorectal cancer: Us preventive services task force recommendation statement. *Jama*, 325(19):1965–1977, 2021.
- [6] Laura Schöckel, Gregor Jost, Peter Seidensticker, Philipp Lengsfeld, Petra Palkowitsch, and Hubertus Pietsch. Developments in x-ray contrast media and the potential impact on computed tomography. *Investigative radiology*, 55(9):592–597, 2020.
- [7] Riccardo Vigneri, Pasqualino Malandrino, and Paolo Vigneri. The changing epidemiology of thyroid cancer: why is incidence increasing? *Current opinion in oncology*, 27(1):1–7, 2015.
- [8] Kenneth G Nepple, Liu Yang, Robert L Grubb, and Seth A Stroepe. Population based analysis of the increasing incidence of kidney cancer in the united states: evaluation of age specific trends from 1975 to 2006. *The Journal of urology*, 187(1):32–38, 2012.
- [9] National Cancer Institute. The surveillance, epidemiology, and end results (seer) program. <https://seer.cancer.gov>, 2016. Accessed: 2024-06-20.
- [10] Peter Armitage and Richard Doll. The age distribution of cancer and a multi-stage theory of carcinogenesis. *British journal of cancer*, 91(12):1983–1989, 2004.
- [11] Suresh H Moolgavkar and David J Venzon. Two-event models for carcinogenesis: incidence curves for childhood and adult tumors. *Mathematical biosciences*, 47(1-2):55–77, 1979.
- [12] Suresh H Moolgavkar and Alfred G Knudson. Mutation and cancer: a model for human carcinogenesis. *JNCI: Journal of the National Cancer Institute*, 66(6):1037–1052, 1981.
- [13] MP Little. Generalisations of the two-mutation and classical multi-stage models of carcinogenesis fitted to the japanese atomic bomb survivor data. *Journal of Radiological Protection*, 16(1):7, 1996.
- [14] E Georg Luebeck and Suresh H Moolgavkar. Multistage carcinogenesis and the incidence of colorectal cancer. *Proceedings of the National Academy of Sciences*, 99(23):15095–15100, 2002.
- [15] Lingling Li, Tianhai Tian, and Xinan Zhang. Stochastic modelling of multistage carcinogenesis and progression of human lung cancer. *Journal of theoretical biology*, 479:81–89, 2019.
- [16] Jihyou Jeon, E Georg Luebeck, and Suresh H Moolgavkar. Age effects and temporal trends in adenocarcinoma of the esophagus and gastric cardia (united states). *Cancer Causes & Control*, 17:971–981, 2006.
- [17] Rafael Meza, Jihyou Jeon, Suresh H Moolgavkar, and E Georg Luebeck. Age-specific incidence of cancer: Phases, transitions, and biological implications. *Proceedings of the National Academy of Sciences*, 105(42):16284–16289, 2008.
- [18] Rafael Meza, Jihyou Jeon, Andrew G Renehan, and Georg Luebeck. Colorectal cancer incidence trends in the us and uk: evidence of right-to left-sided biological gradients with implications for screening. *Cancer research*, 70(13):5419, 2010.
- [19] Anup Dewanji, Jihyou Jeon, Rafael Meza, and E Georg Luebeck. Number and size distribution of colorectal adenomas under the multistage clonal expansion model of cancer. *PLoS Computational Biology*, 7(10):e1002213, 2011.
- [20] E Georg Luebeck, Kit Curtius, Jihyou Jeon, and William D Hazelton. Impact of tumor progression on cancer incidence curves. *Cancer research*, 73(3):1086–1096, 2013.
- [21] Andrew F Brouwer, Marisa C Eisenberg, and Rafael Meza. Case studies of gastric, lung, and oral cancer connect etiologic agent prevalence to cancer incidence. *Cancer research*, 78(12):3386–3396, 2018.
- [22] Rafael Meza, William D Hazelton, Graham A Colditz, and Suresh H Moolgavkar. Analysis of lung cancer incidence in the nurses’ health and the health professionals’ follow-up studies using a multistage carcinogenesis model. *Cancer causes & control*, 19:317–328, 2008.
- [23] Chay Paterson, Hans Clevers, and Ivana Bozic. Mathematical model of colorectal cancer initiation. *Proceedings of the National Academy of Sciences*, 117(34):20681–20688, 2020.
- [24] Lingling Li, Yulu Hu, Yunshan Xu, and Sanyi Tang. Mathematical modeling the order of driver gene mutations in colorectal cancer. *PLOS Computational Biology*, 19(6):e1011225, 2023.

- [25] Cristian Tomasetti, Luigi Marchionni, Martin A Nowak, Giovanni Parmigiani, and Bert Vogelstein. Only three driver gene mutations are required for the development of lung and colorectal cancers. *Proceedings of the National Academy of Sciences*, 112(1):118–123, 2015.
- [26] LJCRC Foulds. The experimental study of tumor progression: a review. *Cancer research*, 14(5):327–339, 1954.
- [27] Bruce Alberts, Dennis Bray, Julian Lewis, Martin Raff, Keith Roberts, James D Watson, et al. *Molecular biology of the cell*, volume 3. Garland New York, 1994.
- [28] Jiri Jungwirth, Marketa Urbanova, Arnoud Boot, Petr Hosek, Petra Bendova, Anna Siskova, Jiri Svec, Milan Kment, Daniela Tumova, Sandra Summerova, et al. Mutational analysis of driver genes defines the colorectal adenoma: in situ carcinoma transition. *Scientific Reports*, 12(1):2570, 2022.
- [29] Guochun Zhang, Yulei Wang, Bo Chen, Liping Guo, Li Cao, Chongyang Ren, Lingzhu Wen, Kai Li, Minghan Jia, Cheukfai Li, et al. Characterization of frequently mutated cancer genes in chinese breast tumors: a comparison of chinese and tcga cohorts. *Annals of translational medicine*, 7(8), 2019.
- [30] Qiang Wang, Bo Yu, Shuilong Zhang, Dongliang Wang, Zhifu Xiao, Hongjing Meng, Lingxiang Dong, Yuhang Zhang, Jie Wu, Zebin Hou, et al. Papillary thyroid carcinoma: Correlation between molecular and clinical features. *Molecular Diagnosis & Therapy*, pages 1–9, 2024.
- [31] Cristian Tomasetti, Lu Li, and Bert Vogelstein. Stem cell divisions, somatic mutations, cancer etiology, and cancer prevention. *Science*, 355(6331):1330–1334, 2017.
- [32] Cristian Tomasetti and Bert Vogelstein. Variation in cancer risk among tissues can be explained by the number of stem cell divisions. *Science*, 347(6217):78–81, 2015.
- [33] Norman TJ Bailey. *The elements of stochastic processes with applications to the natural sciences*, volume 25. John Wiley & Sons, 1991.
- [34] Vinay G Vaidya and Frank J Alexandro Jr. Evaluation of some mathematical models for tumor growth. *International journal of bio-medical computing*, 13(1):19–35, 1982.
- [35] Weixing Dai, Yaqi Li, Xianke Meng, Sanjun Cai, Qingguo Li, and Guoxiang Cai. Does tumor size have its prognostic role in colorectal cancer? re-evaluating its value in colorectal adenocarcinoma with different macroscopic growth pattern. *International journal of surgery*, 45:105–112, 2017.
- [36] Eun Bi Ryu, Jung Min Chang, Mirinae Seo, Sun Ah Kim, Ji He Lim, and Woo Kyung Moon. Tumour volume doubling time of molecular breast cancer subtypes assessed by serial breast ultrasound. *European radiology*, 24:2227–2235, 2014.
- [37] JR Burke, P Brown, A Quyn, H Lambie, D Tolan, and P Sagar. Tumour growth rate of carcinoma of the colon and rectum: retrospective cohort study. *BJS open*, 4(6):1200–1207, 2020.
- [38] Krastan B Blagoev, Julia Wilkerson, Mauricio Burotto, Chul Kim, Edward Espinal-Dominguez, Pilar García-Alfonso, Meghna Alimchandani, Markku Miettinen, Montserrat Blanco-Codesido, and Tito Fojo. Neutral evolution of drug resistant colorectal cancer cell populations is independent of their kras status. *Plos one*, 12(10):e0175484, 2017.
- [39] R Michael Tuttle, James A Fagin, Gerald Minkowitz, Richard J Wong, Benjamin Roman, Snehal Patel, Brian Untch, Ian Ganly, Ashok R Shaha, Jatin P Shah, et al. Natural history and tumor volume kinetics of papillary thyroid cancers during active surveillance. *JAMA Otolaryngology–Head & Neck Surgery*, 143(10):1015–1020, 2017.
- [40] Andrew F Brouwer, Rafael Meza, and Marisa C Eisenberg. A systematic approach to determining the identifiability of multistage carcinogenesis models. *Risk Analysis*, 37(7):1375–1387, 2017.
- [41] Cari M Kitahara and Arthur B Schneider. Epidemiology of thyroid cancer. *Cancer epidemiology, biomarkers & prevention*, 31(7):1284–1297, 2022.
- [42] Sandra J Lee, Xiaoxue Li, Hui Huang, and Marvin Zelen. The dana-farber cisnet model for breast cancer screening strategies: an update. *Medical Decision Making*, 38(1\_suppl):44S–53S, 2018.
- [43] Karen M Kuntz, Iris Lansdorp-Vogelaar, Carolyn M Rutter, Amy B Knudsen, Marjolein Van Ballegooijen, James E Savarino, Eric J Feuer, and Ann G Zauber. A systematic comparison of microsimulation models of colorectal cancer: the role of assumptions about adenoma progression. *Medical Decision Making*, 31(4):530–539, 2011.
- [44] National Center for Health Statistics. Health, united states, 2020-2021: Table canbrtest, 2021. Available from: <https://www.cdc.gov/nchs/hus/data-finder.htm>.
- [45] National Center for Health Statistics. Health, united states, 2019: Table 035, 2019. Available from: <https://www.cdc.gov/nchs/hus/data-finder.htm>.

- [46] Karen Simon. Colorectal cancer development and advances in screening. *Clinical interventions in aging*, pages 967–976, 2016.
- [47] Toru Takano. Natural history of thyroid cancer. *Endocrine journal*, 64(3):237–244, 2017.
- [48] YURI Nikiforov, DOUGLAS R Gnepp, and JAMES A Fagin. Thyroid lesions in children and adolescents after the chernobyl disaster: implications for the study of radiation tumorigenesis. *The Journal of Clinical Endocrinology & Metabolism*, 81(1):9–14, 1996.
- [49] Mykola D Tronko, Vladimir A Saenko, Victor M Shpak, Tetiana I Bogdanova, Shinichi Suzuki, and Shunichi Yamashita. Age distribution of childhood thyroid cancer patients in ukraine after chernobyl and in fukushima after the tepco-fukushima daiichi npp accident. *Thyroid*, 24(10):1547, 2014.
- [50] Shinichi Suzuki. Childhood and adolescent thyroid cancer in fukushima after the fukushima daiichi nuclear power plant accident: 5 years on. *Clinical Oncology*, 28(4):263–271, 2016.
- [51] Jianjiu Chen, Mary Beth Terry, Piero Dalerba, Chin Hur, Jianhua Hu, and Wan Yang. Environmental drivers of the rising incidence of early-onset colorectal cancer in the united states. *International Journal of Cancer*, 154(11):1930–1939, 2024.
- [52] Elizabeth Kagan Arleo, R Edward Hendrick, Mark A Helvie, and Edward A Sickles. Comparison of recommendations for screening mammography using cisnet models. *Cancer*, 123(19):3673–3680, 2017.
- [53] Jeanne S Mandelblatt, Natasha K Stout, Clyde B Schechter, Jeroen J Van Den Broek, Diana L Miglioretti, Martin Krapcho, Amy Trentham-Dietz, Diego Munoz, Sandra J Lee, Donald A Berry, et al. Collaborative modeling of the benefits and harms associated with different us breast cancer screening strategies. *Annals of internal medicine*, 164(4):215–225, 2016.
- [54] Ann G Zauber, Iris Lansdorp-Vogelaar, Amy B Knudsen, Janneke Wilschut, Marjolein Van Ballegooijen, and Karen M Kuntz. Evaluating test strategies for colorectal cancer screening: a decision analysis for the us preventive services task force. *Annals of internal medicine*, 149(9):659–669, 2008.
- [55] Amy B Knudsen, Ann G Zauber, Carolyn M Rutter, Steffie K Naber, V Paul Doria-Rose, Chester Pabiniak, Colden Johanson, Sara E Fischer, Iris Lansdorp-Vogelaar, and Karen M Kuntz. Estimation of benefits, burden, and harms of colorectal cancer screening strategies: modeling study for the us preventive services task force. *Jama*, 315(23):2595–2609, 2016.
- [56] Donald A Berry, Kathleen A Cronin, Sylvia K Plevritis, Dennis G Fryback, Lauren Clarke, Marvin Zelen, Jeanne S Mandelblatt, Andrei Y Yakovlev, J Dik F Habbema, and Eric J Feuer. Effect of screening and adjuvant therapy on mortality from breast cancer. *New England Journal of Medicine*, 353(17):1784–1792, 2005.
- [57] Lauren D Clarke, Sylvia K Plevritis, Rob Boer, Kathleen A Cronin, and Eric J Feuer. Chapter 13: A comparative review of cisnet breast models used to analyze us breast cancer incidence and mortality trends. *JNCI Monographs*, 2006(36):96–105, 2006.
- [58] Carolyn M Rutter, Pedro Nascimento de Lima, Jeffrey K Lee, and Jonathan Ozik. Too good to be true? evaluation of colonoscopy sensitivity assumptions used in policy models. *Cancer Epidemiology, Biomarkers & Prevention*, 31(4):775–782, 2022.
- [59] Tarek A El-Mihoub, Adrian A Hopgood, Lars Nolle, and Alan Battersby. Hybrid genetic algorithms: A review. *Eng. Lett.*, 13(2):124–137, 2006.

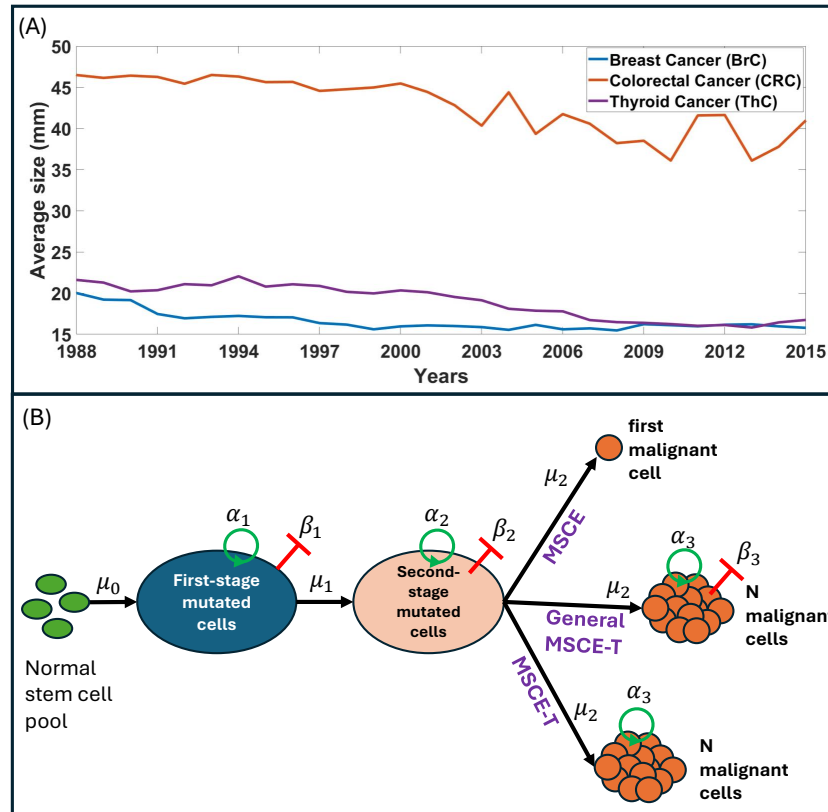


Figure 1: Motivation and schematic of the models. (A) The mean tumor-size-at-diagnosis for breast, colorectal, and thyroid cancer, extracted from the SEER database. (B) The carcinogenesis process from normal stem cells to malignant cells. The classic MSCE model records the cancer incidence at the time of the first malignancy occurrence. The General MSCE-T model considers birth ( $\alpha_3$ ) and death ( $\beta_3$ ) rates for malignant cells and records the cancer incidence when the number of malignant cells reaches  $\geq N$ . The main model (MSCE-T) considers a known proliferation rate ( $\alpha_3$ ) for malignant cells and records the cancer incidence when the number of malignant cells reaches  $\geq N$ .

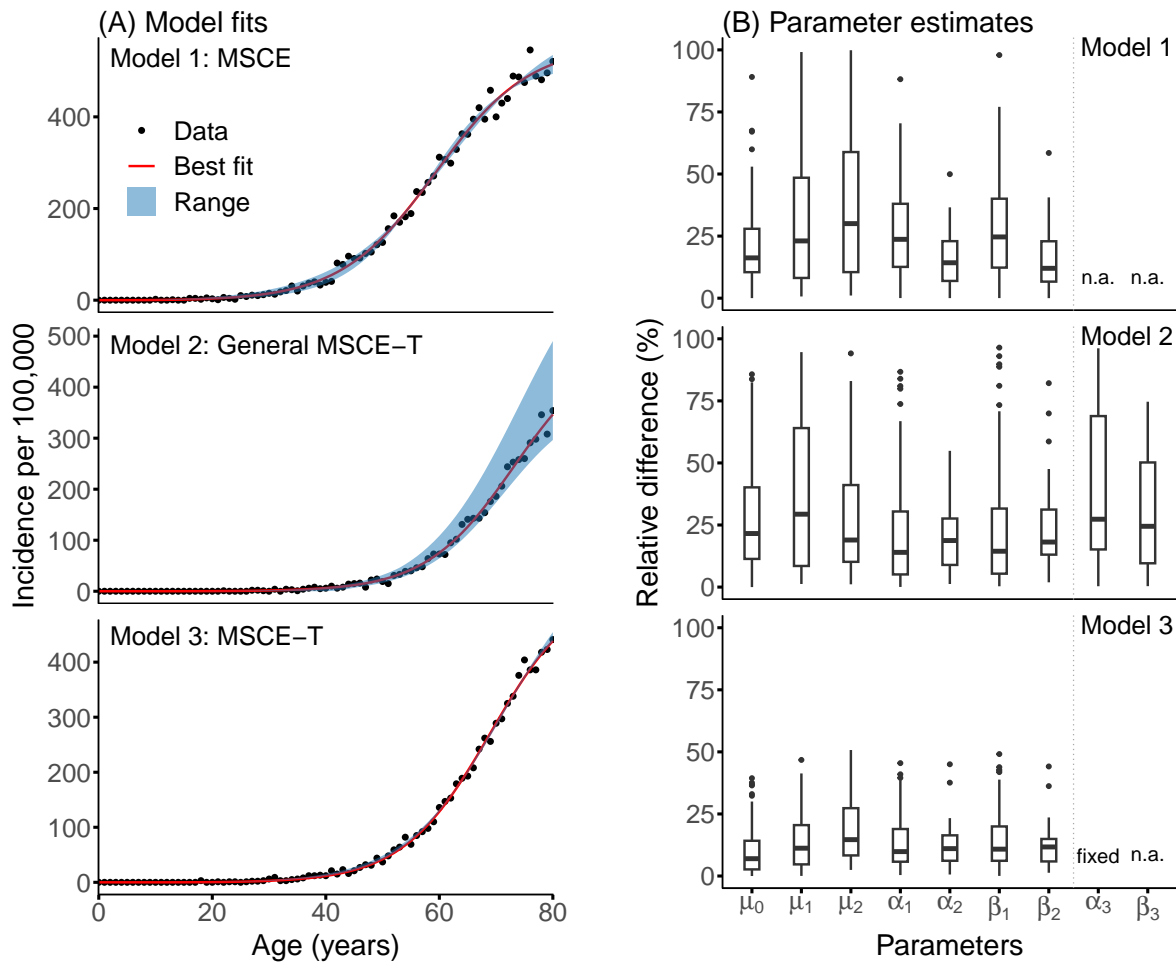


Figure 2: Model validation using synthetic data. (A) Shows the synthetic incidence data (dots) and best model fits (red lines and ranges). (B) Shows the percentage of the relative difference between the estimated parameters and the true values used to generate the synthetic data. The chosen true values for all the models are  $(\mu_0, \mu_1, \mu_2, \alpha_1, \alpha_2, \beta_1, \beta_2) = (5.20 \times 10^{-5}, 3.62 \times 10^{-5}, 1.09 \times 10^{-3}, 2.28, 4.87, 2.16, 4.73)$ . For the General MSCE-T model, we consider  $\alpha_3 = 9$  and  $\beta_3 = 6.43$ . The value of  $\alpha_3$  in the MSCE-T model is fixed and is equal to  $\alpha_3 - \beta_3$  from the General MSCE-T model.

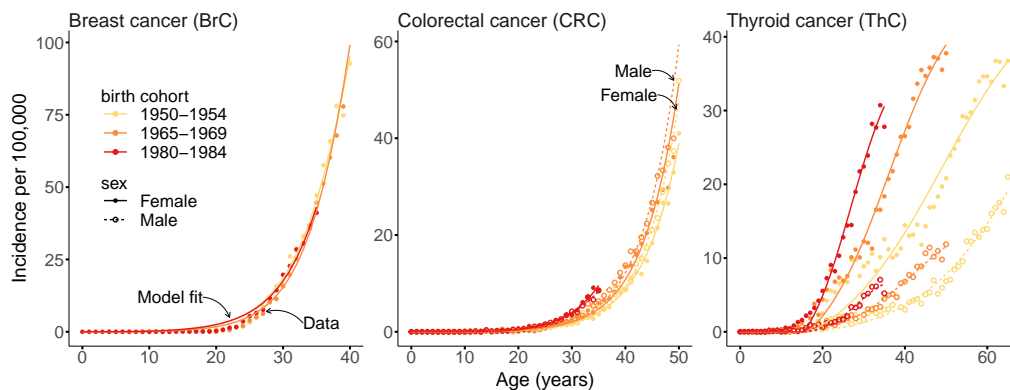


Figure 3: Cancer incidence curves. The curves are produced by fitting the MSCE-T model to incidence data for three different cancer types and three cohorts born in 1950-1954, 1965-1969, and 1980-1984.

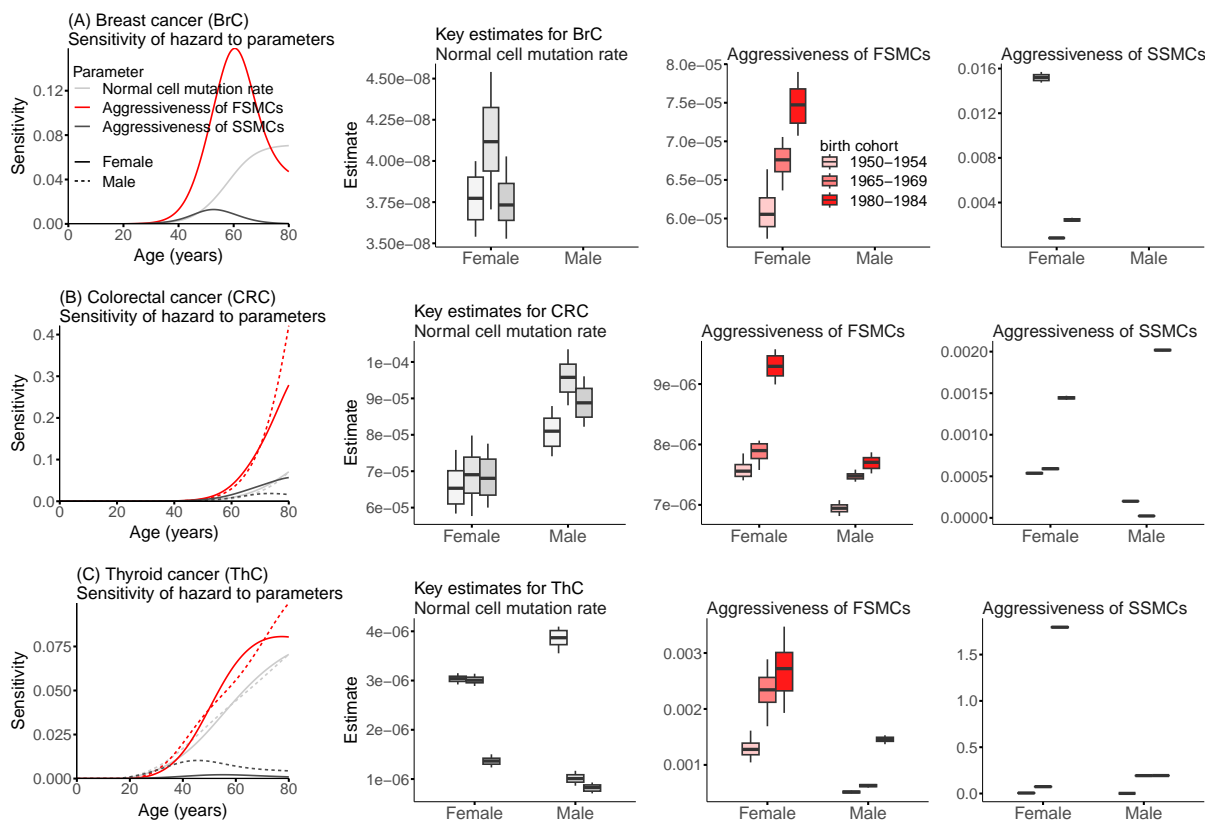


Figure 4: Sensitivity analysis and estimated parameters distribution, for (A) Breast cancer, (B) Colorectal cancer, and (C) Thyroid cancer. The curve plots show the sensitivity of the MSCE-T model hazard to parameters: i) Normal cell mutation rate  $\mu_0$ ; ii) Aggressiveness of the first-stage mutated cells (FSMCs)  $\mu_1 \times (\alpha_1 - \beta_1)$ ; and iii) Aggressiveness of the second-stage mutated cells (SSMCs)  $\mu_2 \times (\alpha_2 - \beta_2)$ . The box plots show the distribution of these parameters for three cohorts. Parameters  $\mu_i$ ,  $\alpha_i$  and  $\beta_i$  all have the unit  $\frac{1}{\text{time}}$ .



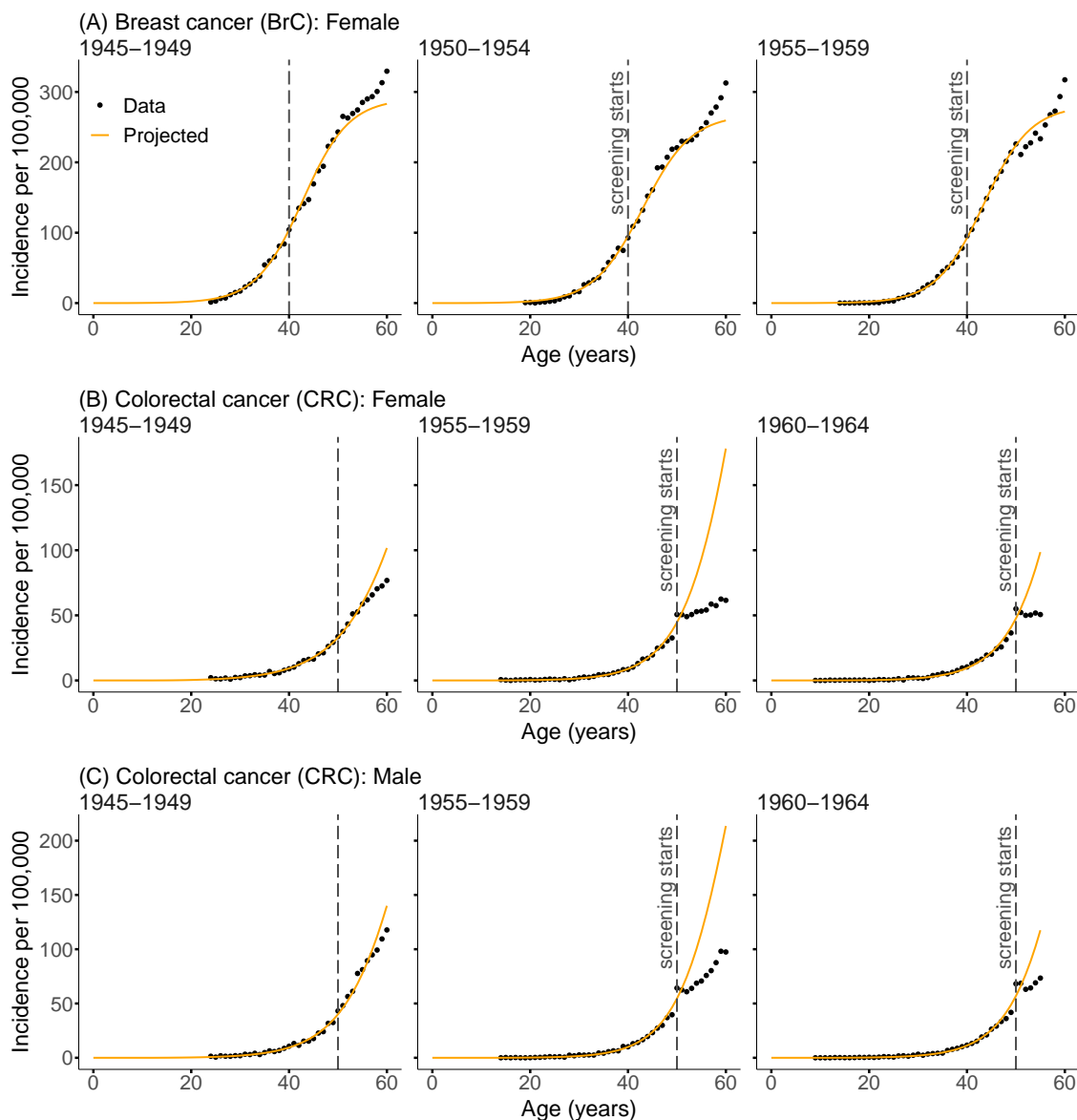


Figure 5: Counterfactual modeling to examine the impact of population-level screening for BrC and CRC. The MSCE-T model is fitted to female BrC patients born in 1945-1949, 1950-1954, and 1955-1959 (sub-figure A). For CRC, the model is fitted to patients born in 1945-1949, 1955-1959, and 1960-1964 (sub-figures B for female and C for male patients). The dashed lines mark the screening initiation. The model is first fitted to incidence before the recommended screening age (dots to the left of the dashed lines) and then used to generate projections (orange curves; i.e., counterfactuals with no screening) for older ages, comparing to the actual incidence (dots to the right of the dashed lines). The dashed lines for the cohort of 1945-1949 are not annotated since screening was not yet as ubiquitous for this cohort.