Between-region genetic divergence reflects mode and tempo of tumor evolution

By Ruping Sun et al.

General Outline

- Problem Statement
- Methods
- Results

General Outline

- Problem Statement
 - \circ Some Context
 - Key Biological Problem
- Methods
- Results

• In order to make precision (specialized medicine) work, there is a need to understand the specifics of how cells evolve across the spectrum of solid tumor types.

- In order to make precision (specialized medicine) work, there is a need to understand the specifics of how cells evolve across the spectrum of solid tumor types.
- Methods to infer the role of natural selection within established tumor types is lacking.

- In order to make precision (specialized medicine) work, there is a need to understand the specifics of how cells evolve across the spectrum of solid tumor types.
- Methods to infer the role of natural selection within established tumor types is lacking.
 - Tumor progression takes years,

- In order to make precision (specialized medicine) work, there is a need to understand the specifics of how cells evolve across the spectrum of solid tumor types.
- Methods to infer the role of natural selection within established tumor types is lacking.
 - Tumor progression takes years,
 - $\circ \qquad {\rm Often \ only \ detected \ in \ late \ stages}$

- In order to make precision (specialized medicine) work, there is a need to understand the specifics of how cells evolve across the spectrum of solid tumor types.
- Methods to infer the role of natural selection within established tumor types is lacking.
 - Tumor progression takes years,
 - Often only detected in late stages
 - This contributes to poor understanding of tumorigenesis and growth immediately after transformation

• Evolution is the product of 3 major underlying processes

- Evolution is the product of 3 major underlying processes
 - Mutation
 - Readily measured in human tumours SSNV calling, CNA estimation etc.

- Evolution is the product of 3 major underlying processes
 - Mutation
 - Readily measured in human tumours
 - Selection
 - Assumed to govern the growth of an established tumor after tumorigenesis (transformation)
 - Acquisition of additional "driver" mutations results in multiple selective sweeps.
 - Due to the selection of driver mutations, hitchhiking passenger mutations can also attain high frequency and manifest as subclonal clusters in bulk sequencing data

- Evolution is the product of 3 major underlying processes
 - Mutation
 - Readily measured in human tumours
 - \circ Selection
 - Assumed to govern the growth of an established tumor after tumorigenesis (transformation)
 - Acquisition of additional "driver" mutations results in multiple selective sweeps.
 - Due to the selection of driver mutations, hitchhiking passenger mutations can also attan high frequency and manifest as subclonal clusters in bulk sequencing data
 - Drift
 - Difficult to distinguish from selection but may also cause extensive ITH.





Goal

• Test hypotheses about underlying evolutionary processes using spatiotemporal patterns of genetic variation among cell populations.

Goal

- Test hypotheses about underlying evolutionary processes using spatiotemporal patterns of genetic variation among cell populations.
- Use simulations to make up lack of high quality early stage data

Goal

- Test hypotheses about underlying evolutionary processes using spatiotemporal patterns of genetic variation among cell populations.
- Use simulations to make up lack of high quality early stage data
- Leverage powerful statistical methods to reveal something about

Key Questions

• What kind of role does selection play in the life of a tumor?

Key Questions

- What kind of role does selection play in the life of a tumor?
- Can you distinguish strong positive selection from weak selection or neutral evolution during tumor progression?

General Outline

- Problem Statement
- Methods
- Results

General Outline

- Problem Statement
- Methods
- Results

Methods





• Simulate various modes of evolution

- Simulate various modes of evolution
 - Neutral
 - Neutral CSC
 - Positive Selection (s = 0.01 -> 0.1 etc.) Acquisition of advantageous mutations adjusts birth-death rate according to s.



• Simulate various modes of evolution

- Simulate various modes of evolution.
- At each time step, a random deme near the periphery of the tumor is chosen for division.
 - This peripheral growth model is consistent with recent studies that show that cells at the periphery are more proliferative than those at the core of a tumor.
 - A random empty lattice site is chosen for the replicate deme.
 - Deme size is chosen conservatively as a maximum of 10k, as large demes hinder selection structurally
 - When a deme reaches maximum size it splits into two offspring demes. The split is modelled by a binomial distribution (N_c, p=0.5) where N_c is the current size of the deme



- Simulate various modes of evolution.
- At each time step, a random deme near the periphery of the tumor is chosen for division



- Simulate various modes of evolution.
- At each time step, a random deme near the periphery of the tumor is chosen for division
- Cell division is a continuous time Markov Process known as a Birth-Death process
 - At any time step there a *p* probability that a cell will divide and a *q* = 1-*p* probability that it will die.
 - Death/Birth rate ratio is important and is characteristic of cancers. E.g. h = q/p = 0.99 in early tumors but
 0.72 in metastatic colorectal cancer.



- Simulate various modes of evolution.
- At each time step, a random deme near the periphery of the tumor is chosen for division
- Cell division is a continuous time Markov Process known as a Birth-Death process
Methods - Spatial Computational Modelling

- Simulate various modes of evolution.
- At each time step, a random deme near the periphery of the tumor is chosen for division
- Cell division is a continuous time Markov Process known as a Birth-Death process
- Poisson process random point mutations *assuming infinite sites model*, at deme division time.
 - Under null models, all mutations are neutral and don't confer fitness advantage
 - Under selection models, beneficial mutations occur as a poisson process with a different mean, and increase the birth rate of the mutated cell by *s*.
 - All beneficial mutations are given a unique index, and other statistics like host cells are recorded.

Methods - Spatial Computational Modelling



Methods - Spatial Computational Modelling

- Simulate various modes of evolution.
- At each time step, a random deme near the periphery of the tumor is chosen for division
- Cell division is a continuous time Markov Process known as a Birth-Death process
- Poisson process random point mutations *assuming infinite sites model*, at deme division time.

Methods - SSNV Calling, SCNA detection and VAF adjustment

- For each raw SNV call from MuTect, read alignmet features from all samples were reinspected in an automated fashion to assess confidence.
- TitanCNA was used to estimate somatic CNAs, and the observed VAF for each detected SSNV was adjusted on the basis of cancer cell fraction, tumor purities, as well as local copy number estimates.

$$VAF_{a} = \frac{CCF_{est}}{2}, \qquad CCF_{est} = \begin{cases} N_{c} \times \frac{VAF_{o}}{p'} - P_{CNA} \times (N_{t} - N_{b} - 1) & Early Major \\ N_{c} \times \frac{VAF_{o}}{p'} - P_{CNA} \times (N_{b} - 1) & Early Minor \\ N_{c} \times \frac{VAF_{o}}{p'} & Late/Independent \end{cases}$$

- An SSNV m is *subclonal* if all of the following criteria are met
 - \circ Total Probability P_m < 0.05

$$P_m = \prod_{i=1}^k P_{mi} \left(X_{mi} \le S_{mi}, N_{mi}, f.pub_{mi} \right) < 0.05$$

$$f.pub_{mi} = \begin{cases} pu_i \times \frac{nb_{mi}}{nc_{mi}} & \text{if } nb_{mi} \ge 1, nt_{mi} \ge 2\\ (pu_i \times (nt_{mi} - nb_{mi})) / nc_{mi} & \text{otherwise} \end{cases}$$

- P_{mi} is the probability for region *i* of observing less than S_{mi} reads having the mutant allele out of N_{mi} reads.
- This is provided the lower bound of *f* described above which is the expected frequency if the SSNV is public.

- An SSNV m is *subclonal* if all of the following criteria are met
 - \circ Total Probability P_m < 0.05

- An SSNV m is *subclonal* if all of the following criteria are met
 - Total Probability $P_m < 0.05$
 - At least one region *i* with $CCF_m \pm 95\% CI_{mi} < 1$

- An SSNV m is *subclonal* if all of the following criteria are met
 - Total Probability $P_m < 0.05$
 - At least one region *i* with $CCF_m \pm 95\% CI_{mi} < 1$
 - At least one region *i* with adjusted VAF_{mi} < 0.25, which was chosen experimentally.

- An SSNV m is *subclonal* if all of the following criteria are met
 - Total Probability $P_m < 0.05$
 - At least one region *i* with $CCF_m \pm 95\% CI_{mi} < 1$
 - At least one region *i* with adjusted VAF_{mi} < 0.25, which was chosen experimentally.
- Any SSNV that fails to meet one of the above criteria is considered public!

• The first is fHrs, the fraction of high frequency region-specific subclonal SSNVs out of all region-specific subclonal SSNVs

2)
$$\text{fHrs} = \frac{1}{r} \times \sum_{j=1}^{r} \left(\frac{\text{RSM}_{ja}^{\text{high}}}{2 \times \text{RSM}_{ja}^{\text{all}}} + \frac{\text{RSM}_{jb}^{\text{high}}}{2 \times \text{RSM}_{jb}^{\text{all}}} \right)$$

where $\text{RSM}_{ja}^{\text{high}}$ and $\text{RSM}_{ja}^{\text{all}}$ represent the number of high-frequency (VAF > 0.2) region-specific SSNVs and the number of all region-specific SSNVs with VAF >0.08 for region *a*, in a pairwise comparison *j* between regions *a* and *b*.

- The first is **fHrs**, the fraction of high frequency region-specific subclonal SSNVs out of all region-specific subclonal SSNVs.
- Next is **fHsub**, which is the fraction of subclonal SSNVs with high frequency (VAF > 0.2).

1)
$$fHsub = \frac{1}{k} \times \sum_{i=1}^{k} \frac{SM_i^{high}}{SM_i^{all}}$$

where SM_i^{high} and SM_i^{all} are the number of high-frequency subclonal SSNVs (adjusted VAF > 0.2, hereafter referred to as VAF) and the number of all subclonal SSNVs with VAF >0.08 for region *i*. The cutoff was set to 0.2 because above this value fHsub tends to plateau in its sensitivity to distinguishing the neutral and selection models (**Supplementary Fig. 36**). A lower cutoff of 0.08 was chosen empirically to satisfy the tradeoff between the number of subclonal SSNVs and variant calling errors.

- The first is **fHrs**, the fraction of high frequency region-specific subclonal SSNVs out of all region-specific subclonal SSNVs.
- Next is **fHsub**, which is the fraction of subclonal SSNVs with high frequency (VAF > 0.2).
- KSD (Kolmogorov-Smirnov distance) estimates the dissimilarity of the SFSs between the two regions

$$\text{KSD} = \frac{1}{r} \times \sum_{j=1}^{r} \text{KSD}_{j}$$

and

4)

$$\mathrm{KSD}_j = \max \left| F_a - F_b \right|$$

where F_a is the cumulative SFS of region a, in a pairwise comparison j between regions a and b.

- The first is **fHrs**, the fraction of high frequency region-specific subclonal SSNVs out of all region-specific subclonal SSNVs.
- Next is **fHsub**, which is the fraction of subclonal SSNVs with high frequency (VAF > 0.2).
- KSD (Kolmogorov-Smirnov distance) estimates the dissimilarity of the SFSs between the two regions.
- **FST** (fixation index) a measure of genetic divergence between the two regions.

3)

 $FST = \frac{1}{r} \times \sum_{i=1}^{r} FST_j^{Hudson}$ and $FST_{j}^{Hudson} = \frac{\sum_{m=1}^{m_{t}} (f_{a}^{m} - f_{b}^{m})^{2} - \frac{f_{a}^{m} \times (1 - f_{a}^{m})}{d_{a}^{m} - 1} - \frac{f_{b}^{m} \times (1 - f_{b}^{m})}{d_{b}^{m} - 1}}{\sum_{m=1}^{m_{t}} f_{a}^{m} \times (1 - f_{b}^{m}) + f_{b}^{m} \times (1 - f_{a}^{m})}}$ and

where f_a^m is the VAF for SSNV *m* and d_a^m is the sequencing depth for SSNV *m* in region *a*. The genetic variance components (numerator and denominator) are averaged separately to obtain a ratio combining the Hudson FST estimates across all m_t SSNVs⁵⁴.

- The first is **fHrs**, the fraction of high frequency region-specific subclonal SSNVs out of all region-specific subclonal SSNVs.
- Next is **fHsub**, which is the fraction of subclonal SSNVs with high frequency (VAF > 0.2).
- KSD (Kolmogorov-Smirnov distance) estimates the dissimilarity of the SFSs between the two regions.
- **FST** (fixation index) a measure of genetic divergence between the two regions.
- **rAUC** The ratio of the area under the pooled cumulative SFS to the area under the theoretical cumulative SFS assuming neutral growth in a well-mixed population. **For MRS**, *pooled VAF = #alternative alleles/total read depth*.

5)

$$rAUC = \frac{AUC_{merged}}{AUC_{theoretical}}$$

• A support vector machine classifier with a radial basis function (RBF) kernel was trained on 1400 simulated tumors derived from 7 different growth models first 3 models labelled as "effectively neutral" and the other 4 as "selection", using 10-fold cross-validation with the 5 ITH metrics as features.

- A support vector machine classifier with a radial basis function (RBF) kernel was trained on 1400 simulated tumors derived from 7 different growth models first 3 models labelled as "effectively neutral" and the other 4 as "selection", using 10-fold cross-validation with the 5 ITH metrics as features.
- Hyperparameter optimization/Grid Search for optimizing RBF parameters.

- A support vector machine classifier with a radial basis function (RBF) kernel was trained on 1400 simulated tumors derived from 7 different growth models first 3 models labelled as "effectively neutral" and the other 4 as "selection", using 10-fold cross-validation with the 5 ITH metrics as features.
- Hyperparameter optimization/Grid Search for optimizing RBF parameters.
- Another SVM was trained on the ICs derived from ICA of the 5 ITH metrics above.

General Outline

- Problem Statement
- Solution and Methods
- Results



• Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.

Results of Spatial Simulation



Results of Spatial Simulation



Results of Spatial Simulation



Theoretical Neutral AFS with exponential growth (1/f) - - Mean AFS for 100 simulations

• Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)
 - Consists of 8 tumors W,M,G,N,U,O, and S.
 - These samples were passed through the VAP :
 - M, O, U showed ITH metrics comparable to virtual tumors under neutral growth
 - G,N,W, and S exhibited only slightly higher values, but consistent with weak selection.



- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)
 - Consists of 8 tumors W,M,G,N,U,O, and S.
 - These samples were passed through the VAP :
 - M, O, U showed ITH metrics comparable to virtual tumors under neutral growth
 - G,N,W, and S exhibited only slightly higher values, but consistent with weak selection.
- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)
 - In vivo COAD xenograft
 - Generated a type of cancerous COAD cell line and xenotransplanted it into immune-compromised mice.
 - The resulting tumors were sequenced.



HCT116 xenografts Single Cell Tissue Culture (TC) Transplantation

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)
 - In vivo COAD xenograft
 - Generated a type of cancerous COAD cell line and xenotransplanted it into immune-compromised mice.
 - The resulting tumors were sequenced.

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)
 - In vivo COAD xenograft

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)
 - In vivo COAD xenograft
 - Samples from 4 types of solid tumors.
 - Additionally, for use as positive controls:
 - 2 tumors treated with a mutating agent known to impose positive selective pressure
 - Premalignant lesions (Barrett's esophagus lesions) were also used.







0.5

VAF

0.7

0.9 1

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
 - COAD colorectal adenocarcinoma dataset (taken > 3cm apart)
 - In vivo COAD xenograft
 - Samples from 4 types of solid tumors + 4 positive controls

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.

- Tumors evolving neutrally and through strong positive selection show fundamentally different patterns of intra-tumor heterogeneity, and these can be distinguished through multi-region sequencing.
- Patterns of ITH evaluated in several publicly available MRS datasets spanning multiple tumor types.
- ITH metrics were calculated for all datasets, classified, and projected into model space.
 - ICA (independent component analysis) using the five ITH metrics identified two distinct clusters, corresponding to selection with s>0.02.



Results b Neutral 0 e-neutral Neutral (CSC) • s=0.01 • s=0.02 ⁻⁻ selecton Independent Component (IC) 2 s=0.03 0.8 • s=0.05 • s=0.1 GLM-27PF ESCA_BE-14-exon 0.6 CRA-SO ESCA-14-exon BE-4-exon Xeno-HCT116-S1 ESCA-4-exon OGBM-0221P GLM-18 9-HCT116-S3 OESCA-14 COAD-M ESCA_BE-14 0.4 COAD-GO 116-S3 COAD-D Xeno-LOVO-D COAD-U OLUAD-270 OCOAD-N **OGBM-0125P** COAD-W GBM-0221PF ESCA-4 O ESCA-8-exoro LUAD-292^O OESCA-LUAD-001 GLM-26P Xeno-LOVO-00 0.2 OESCA-8 NSCLC-004 Decision Boundary for . . Virtual Tumors ESCA BE-4-exon Primary Tumors LUAD-4990 • BE BE vs ESCA ESCA_BE-4 O - O TX GBM-0125PR 0.2 0.4 0.6 0.8 0

Independent Component (IC) 1

1