



Cancer Evolution: from processes to genomic patterns and back again

International Agency for Research on Cancer
Lyon, France

Nicolas Alcala, PhD
Scientist, Rare Cancers Genetics Team, Genomic Epidemiology Branch
December 1st, 2020

Reminder: cancer evolution shapes observed heterogeneity

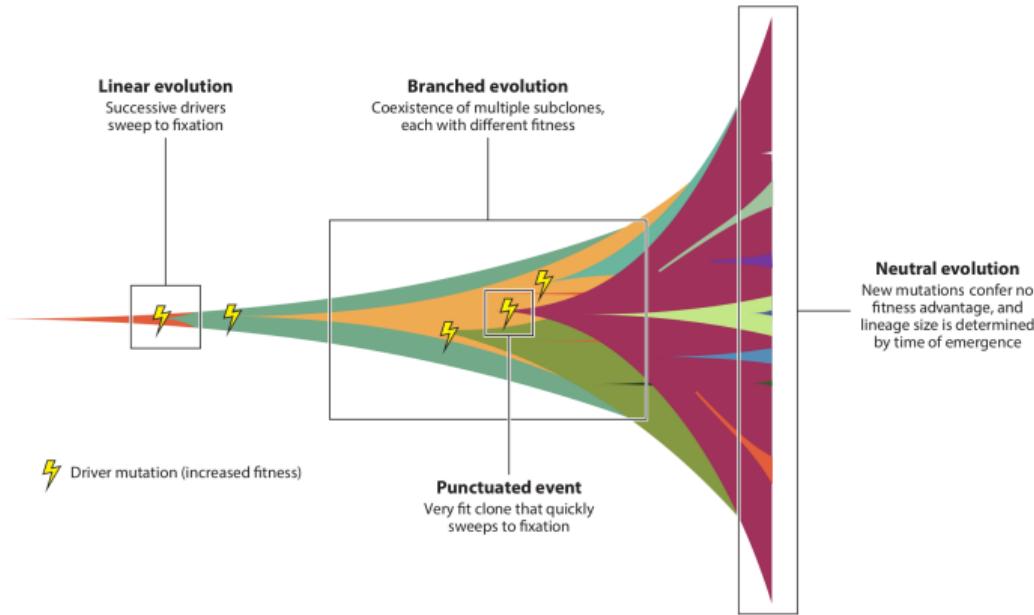
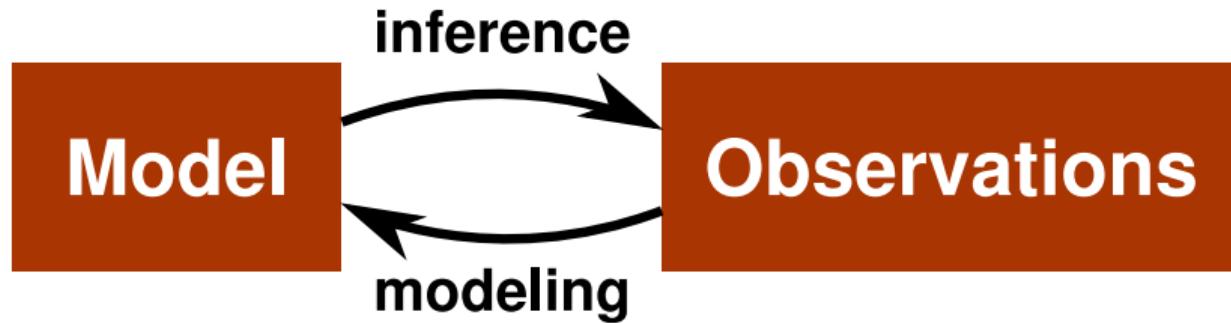


Figure: Schematic view of cancer evolution (Williams *et al.* 2019).

Why do we need models?

Models enable quantitative predictions about genomic patterns that can be tested using statistical inference methods



What are the important parameters?

Evolutionary parameters

- **Cancer cell population size:** gives an idea of tumor progression
- **Cell generation time:** rate of division of self-renewing cells
- **Mutation rate:** how often do alterations appear?
- **Selection coefficients:** do different tumor cells have different growth rates?

Ecological parameters

- **Ecological strategy:** what is the phenotype of the cancer cells? how does it provide an advantage?
- **Carrying capacity/Habitat quality:** what is the extent of vasculature and hypoxic niches? what is the nutrient concentration (e.g., glucose)?
- **Migration rate:** how often do cells migrate/metastasize to other locations?
- **Predation:** how much do immune cells (e.g., tumor-infiltrating lymphocytes) pressure the tumor?

Part I: evolution of the cancer genome

Modeling approaches

Theoretical frameworks are based on **population genetics** (the study of mutations in populations—at micro-evolutionary time-scales) and **molecular phylogenetics** (the study of mutations in species—at macro-evolutionary time-scales)

Modeling approaches

Theoretical frameworks are based on **population genetics** (the study of mutations in populations—at micro-evolutionary time-scales) and **molecular phylogenetics** (the study of mutations in species—at macro-evolutionary time-scales)

Frameworks

- **Mathematical models:** *equations* describing the relationships between processes and observable quantities; allow precise description of the impact of evolutionary processes on genomic patterns
- **Computational models:** *algorithm* describing the effect of processes on observable quantities; allows complex models (e.g., model each cell)

Features

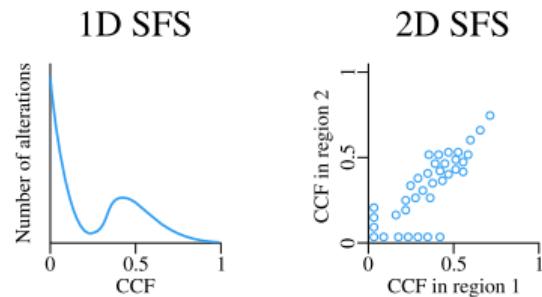
- e.g., spatially implicit or explicit

What is the main quantity of interest?

The cancer cell fraction (CCF)

Definition: proportion of tumor cells that harbor a given somatic alteration (e.g., mutation, indel, structural variant)

Typical representation: the Site Frequency Spectrum (SFS)



How to get the CCF?

Goal: go from variant allelic fraction (VAF, proportion of reads carrying the alteration) to cancer cell fraction (CCF)

Method: call single nucleotide variants (SNVs) and copy number variants (CNVs), and correct VAFs for CN state

- Easiest in diploid regions: then $VAF = \frac{1}{2}CCF \times \phi$, so

$$CCF = \frac{2}{\phi}VAF,$$

where ϕ is the tumor purity. *Example: a heterozygous variant present in CCF = 100% of tumor cells in a tumor with purity $\phi = 80\%$ has $VAF = \frac{1}{2} \times 1 \times 0.8 = 0.4$.*

How to get the CCF?

Goal: go from variant allelic fraction (VAF, proportion of reads carrying the alteration) to cancer cell fraction (CCF)

Method: call single nucleotide variants (SNVs) and copy number variants (CNVs), and correct VAFs for CN state

- Easiest in diploid regions: then $VAF = \frac{1}{2}CCF \times \phi$, so

$$CCF = \frac{2}{\phi}VAF,$$

where ϕ is the tumor purity. *Example: a heterozygous variant present in CCF = 100% of tumor cells in a tumor with purity $\phi = 80\%$ has $VAF = \frac{1}{2} \times 1 \times 0.8 = 0.4$.*

- In non-diploid regions:

$$CCF = \frac{2(1 - \phi) + C_{\text{total}}\phi}{C_{\text{mut}}\phi}VAF,$$

where C_{total} is the total CN in the tumor and C_{mut} is the allele-specific CN.

Approach (i): a simple deterministic mathematical model

Williams *et al.* (2016): constant mutation rate + exponential growth (neutral evolution)

Features: mathematical model, spatially implicit

Approach (i): a simple deterministic mathematical model

Williams *et al.* (2016): constant mutation rate + exponential growth (neutral evolution)

Features: mathematical model, spatially implicit

Model:

- Number of mutations M : $\frac{dM}{dt} = \phi\mu\lambda N(t)$ (i), where μ is the mutation rate and λ the growth rate
- Population size: $N(t) = e^{\lambda\beta t}$ (ii), where β is proportion of surviving cells per division
- Frequency of a mutation (under infinite sites model): $f = \frac{1}{\phi N(t)}$ (iii)

Approach (i): a simple deterministic mathematical model

Williams *et al.* (2016): constant mutation rate + exponential growth (neutral evolution)

Features: mathematical model, spatially implicit

Model:

- Number of mutations M : $\frac{dM}{dt} = \phi\mu\lambda N(t)$ (i), where μ is the mutation rate and λ the growth rate
- Population size: $N(t) = e^{\lambda\beta t}$ (ii), where β is proportion of surviving cells per division
- Frequency of a mutation (under infinite sites model): $f = \frac{1}{\phi N(t)}$ (iii)

Observable values: $M(f)$, the cumulative number of mutations at frequency f , $dM(f)/df$, the number of mutations at frequency f

Solving (i)-(iii) for $M(f)$ gives:

$$M(f) = \frac{\mu}{\beta} \left(\frac{1}{f} - \frac{1}{f_{\max}} \right) \quad (1)$$

Approach (i): a simple deterministic mathematical model, inference

Eq. 1 is equivalent to a simple linear model

$$M(x) = ax + b,$$

with $x = 1/f$, $a = \frac{\mu}{\beta}$, and $b = -\frac{\mu}{\beta f_{\max}}$.

Thus, we can estimate a scaled mutation rate μ/β , and a general model fit using **linear regression**

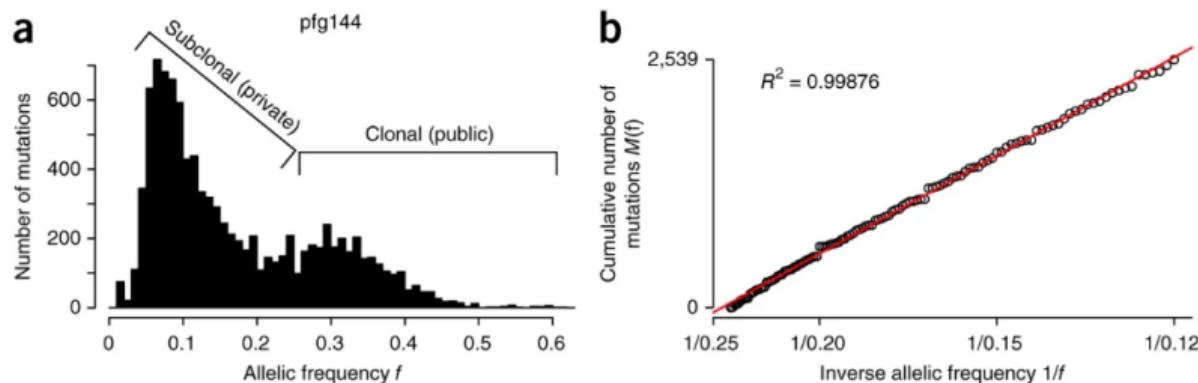


Figure: Gastric cancer whole-genome sequencing data (Williams *et al.* 2016).

Approach (i): a simple deterministic mathematical model, inference

Advantages:

- Mathematically tractable
- Simple, general predictions

Criticisms:

- Other processes can produce similar predictions (Tarabichi *et al.* 2018)
- Deterministic nature does not capture variability of evolutionary processes and does not allow classical hypothesis testing
- Is neutral evolution the null in cancer evolution?

Approach (ii): a more complicated probabilistic mathematical model

Caravagna *et al.* (2020): mixture of neutrally evolving (constant mutation rate + exponential growth) and selected clones (exponential growth at rate higher than the neutral growth)

Features: mathematical model, spatially implicit

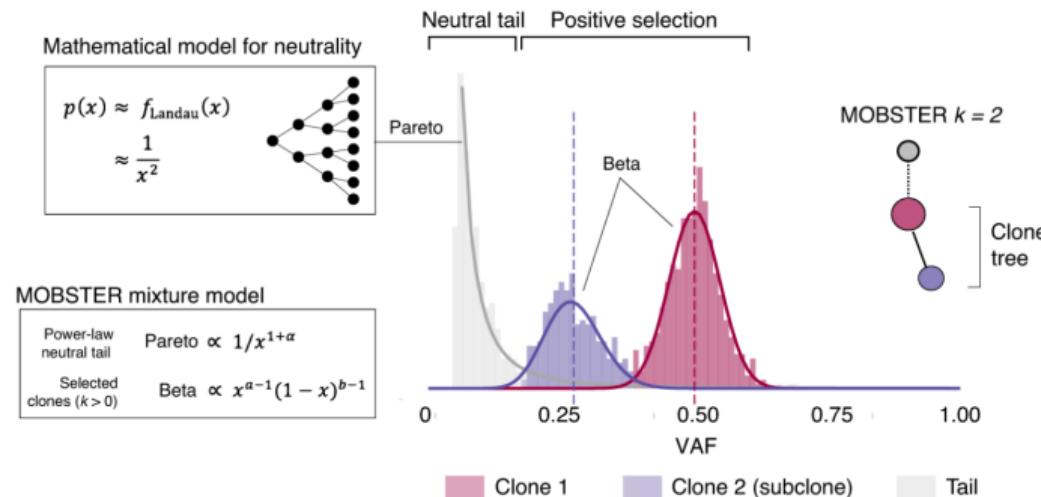


Figure: MOBSTER model (Caravagna *et al.* 2020). The neutral tail is based on the solution to the stochastic Luria–Delbrück model of bacterial growth (Kessler and Levine 2013; Luria and Delbrück 1943).

Approach (ii): a more complicated probabilistic mathematical model, inference

Fit a $k + 1$ component mixture model, where each alterations comes either from the neutral tail or from one of k selected clones

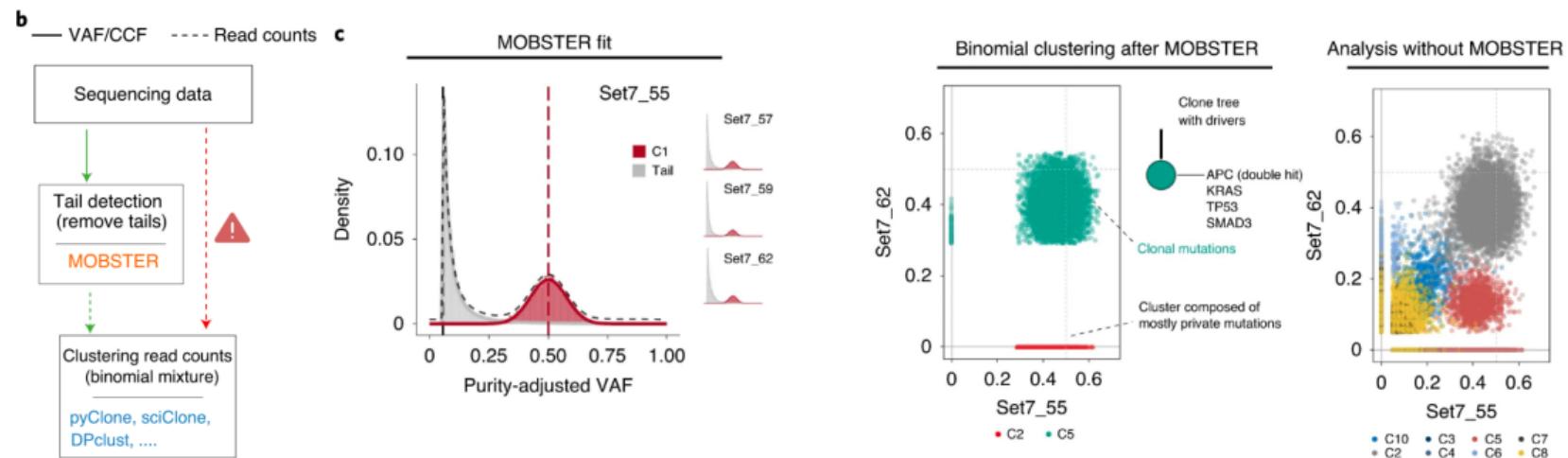


Figure: MOBSTER inference scheme (Caravagna *et al.* 2020).

Approach (ii): a more complicated probabilistic mathematical model, inference

Advantages:

- Encompasses many scenarios
- Allows inference of population parameters

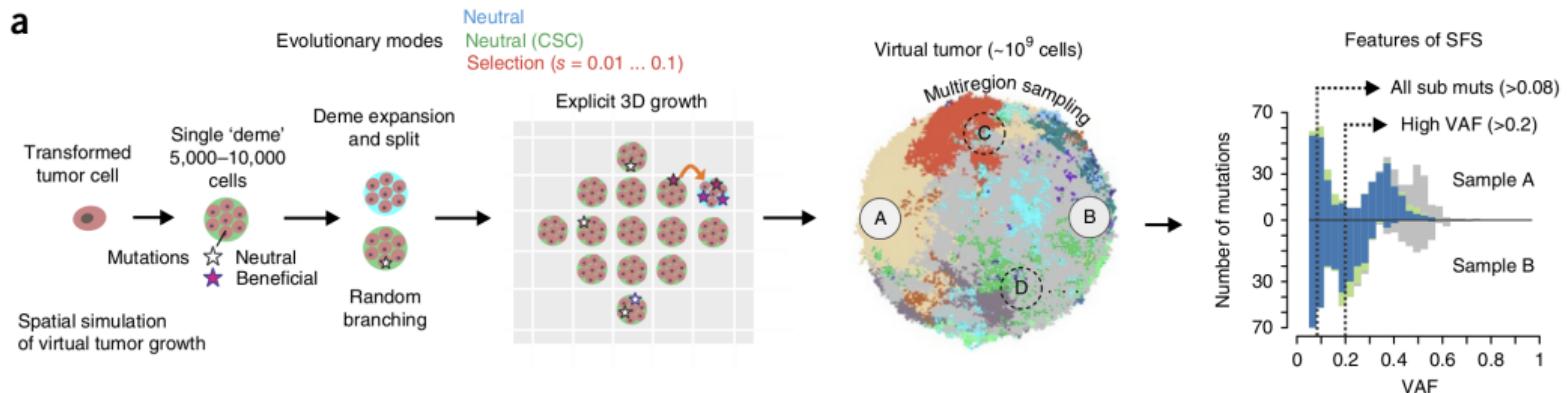
Criticisms:

- Spatially implicit \Rightarrow does not account for spatial heterogeneity
- Only considers exponential growth \Rightarrow mostly suited for initial progression phase, cannot detect slowing down or tumor regression
- Only considers competition between clones \Rightarrow not commensalism or cooperation that could help maintain stable proportions of different clones through time

Approach (iii): a computational model

Sun *et al.* (2017): constant mutation rate + spatially explicit 3-D growth + selection

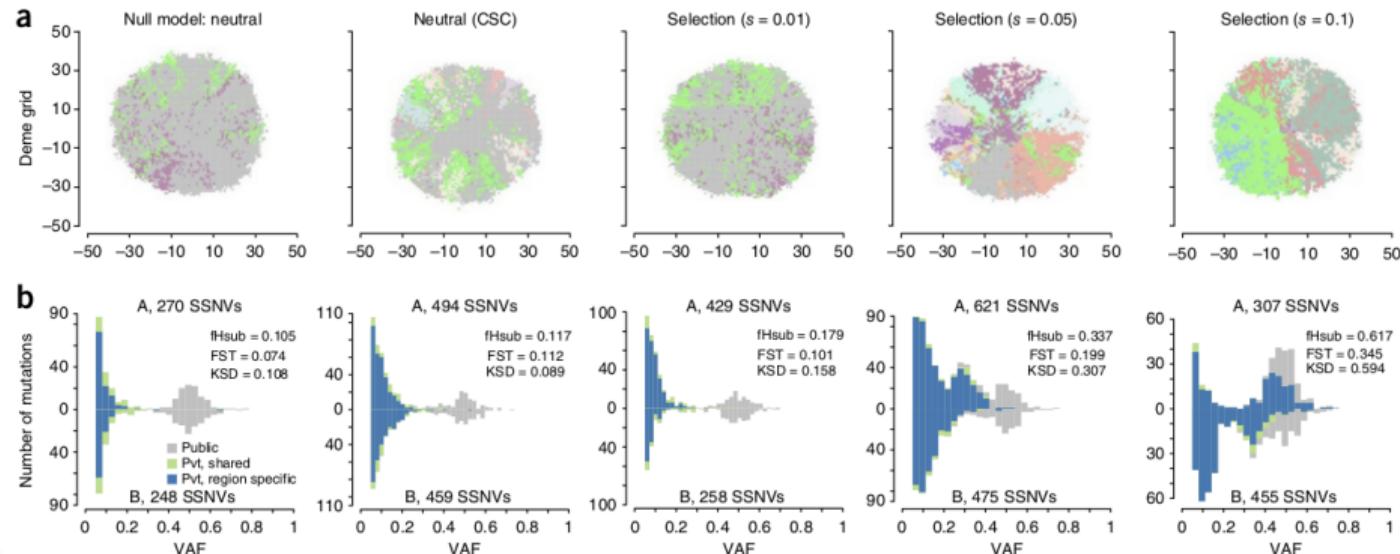
Model:



Approach (iii): a computational model

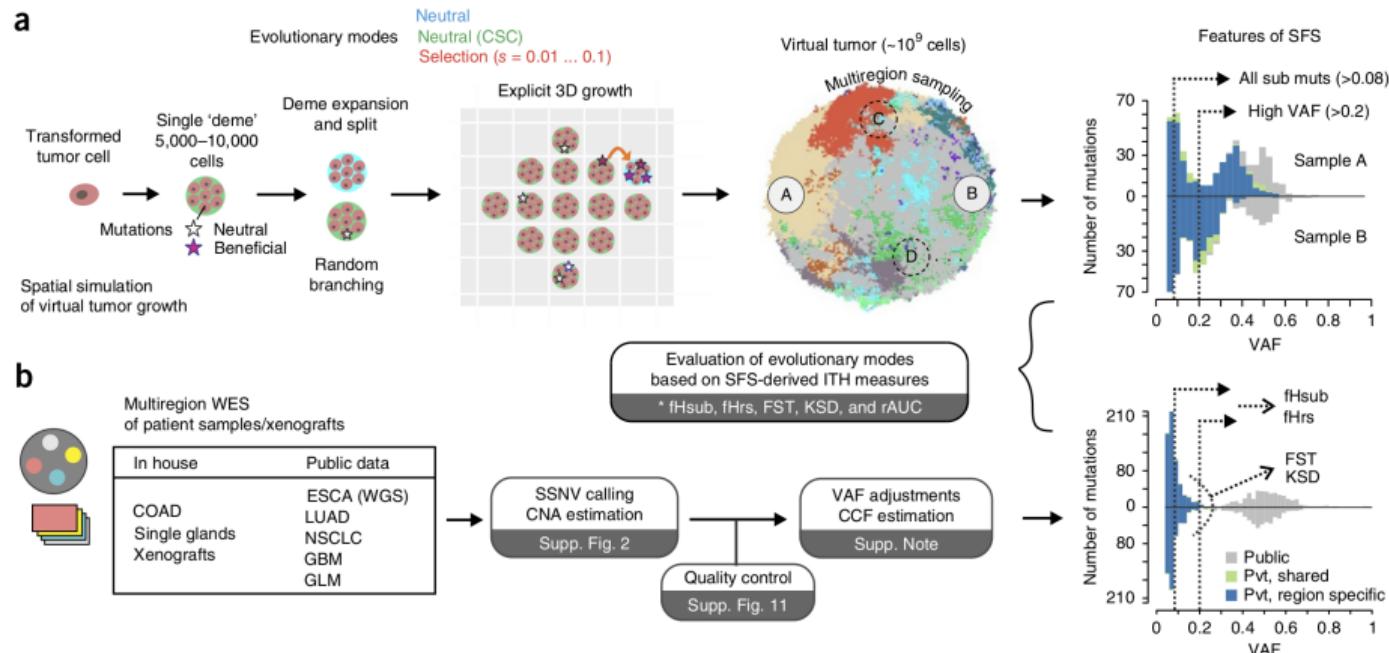
The model allows to make **quantitative predictions** in different scenarios, and distinguish **qualitative behaviors**:

- neutral or effectively neutral models lead to weak spatial correlation, little differentiation between regions, and no private “bumps” in the SFS



Approach (iii): a computational model, inference

Compute summary statistics, compare their values with that of many simulations ($\sim 10^6$); Retain the parameters of simulations with stats closest to the observations (Approximate Bayesian Computation)



Approach (iii): a computational model, inference

Advantages:

- Allows arbitrarily complex scenarios that fit the biology

Criticisms:

- Computationally intensive (large parameter space to explore)
- Identifiability is unknown (worked in this case but might not for other applications)

Part II: evolution of the cancer transcriptome

Modeling approaches

Theoretical frameworks are based on **quantitative genetics** (the study of the evolution of phenotypes in populations) and **phylogenetics of traits** (the study of the evolution of phenotypes in species)

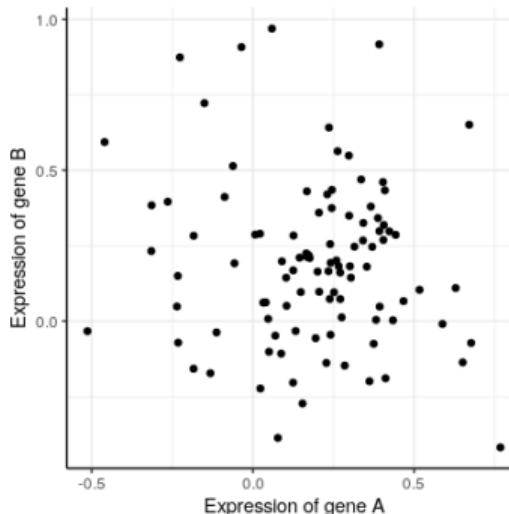
Most traits have **complex genetic architectures**, that are often unknown

What is the main quantity of interest?

Molecular traits

Definition: quantitative molecular variables (e.g., gene expression) that have an important role in cancer initiation and progression

Typical representation: a scatterplot, where each dimension represents a trait of interest



How to get the right molecular traits?

Goal: identify molecular the variables that are important for cancer evolution

Methods:

- (Supervised) use biological knowledge of a cancer type to choose relevant pathways
- (Unsupervised) leverage the between-patient molecular variation within a cohort to identify major sources of variation

An equilibrium deterministic mathematical model

Hausser and Alon (2020): directional selection + constraints (trade-off)

Features: mathematical model, unsupervised, equilibrium model (selection had enough time to act)

An equilibrium deterministic mathematical model

Hausser and Alon (2020): directional selection + constraints (trade-off)

Features: mathematical model, unsupervised, equilibrium model (selection had enough time to act)

Model:

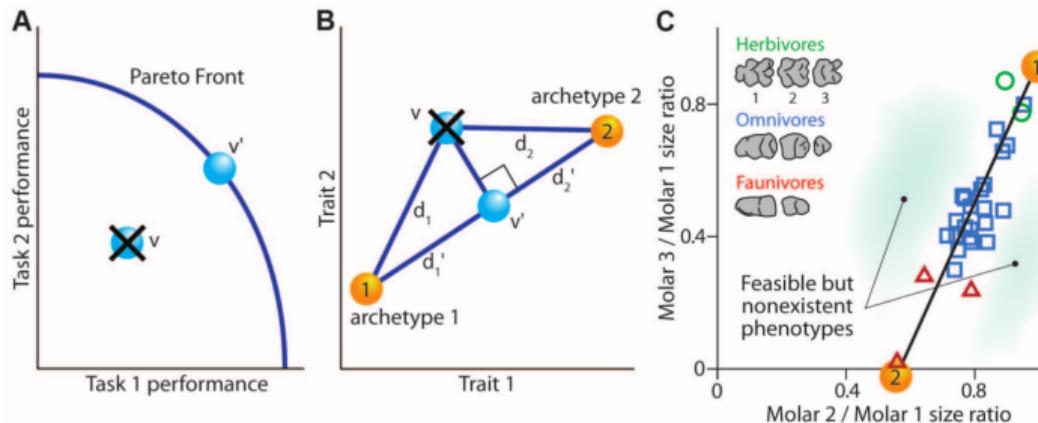
- cancer cells have to perform **multiple tasks** (e.g., cell division, energy production, interaction with immune system, invasion and tissue remodeling)
- cells **cannot excel at everything** (e.g., the optimal phenotype for cell division will not be optimal in terms of interaction with immune system) \Rightarrow **trade-offs**

An equilibrium deterministic mathematical model

Hausser and Alon (2020): directional selection + constraints (trade-off)

Features: mathematical model, unsupervised, equilibrium model (selection had enough time to act)

- **Pareto front:** set of phenotypes that cannot be improved at all tasks at once (increasing some performances requires reducing performance at other tasks)
- **Archetype:** optimal phenotype for a task



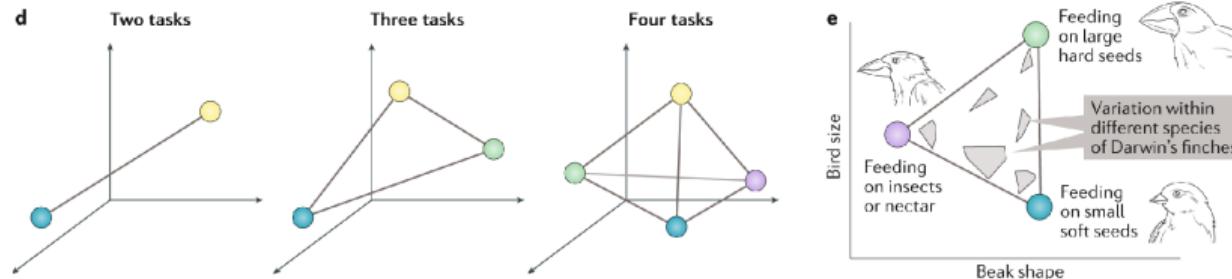
An equilibrium deterministic mathematical model

Hausser and Alon (2020): directional selection + constraints (trade-off)

Features: mathematical model, unsupervised, equilibrium model (selection had enough time to act)

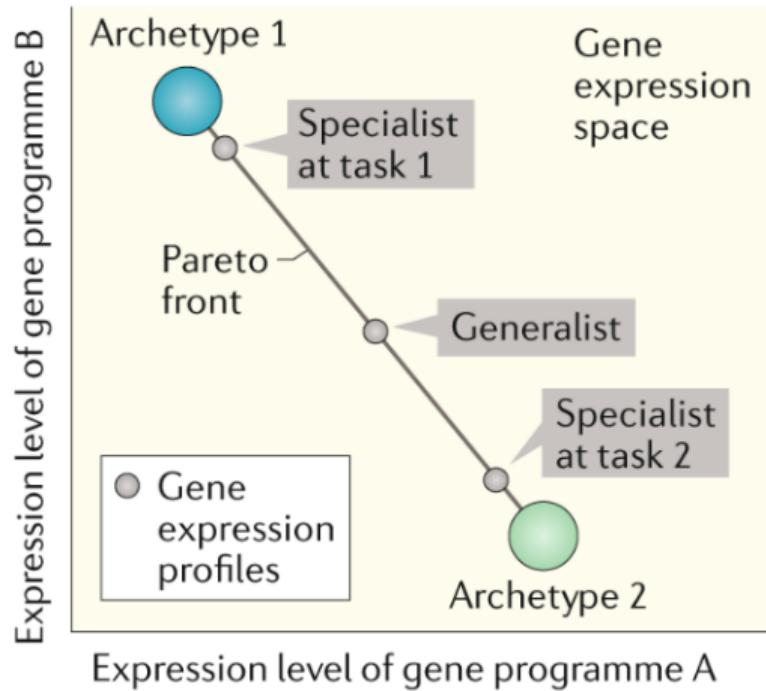
Results:

- the Pareto front has a **simple geometric shape** that only depends on the number of tasks (2 tasks: line, 3 tasks: triangle, four tasks: tetrahedron, etc)
- position within the front depends on the strength of selection for each task (strong selection for a task: specialist, equal selection for each task: generalist)



An equilibrium deterministic mathematical model

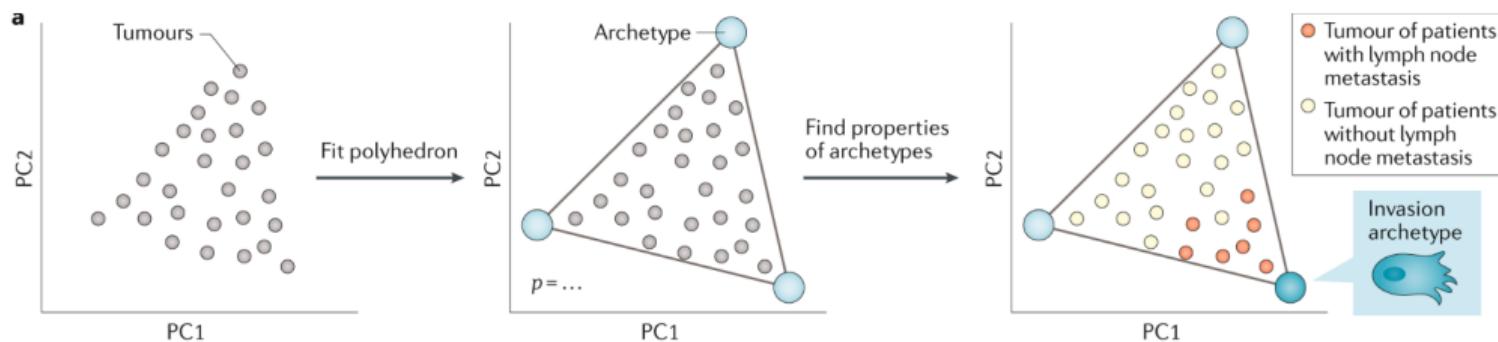
Hausser *et al.* (2019): we expect cancer transcriptomes to be a proxy for the molecular phenotypes facing trade-offs between tasks



An equilibrium deterministic mathematical model

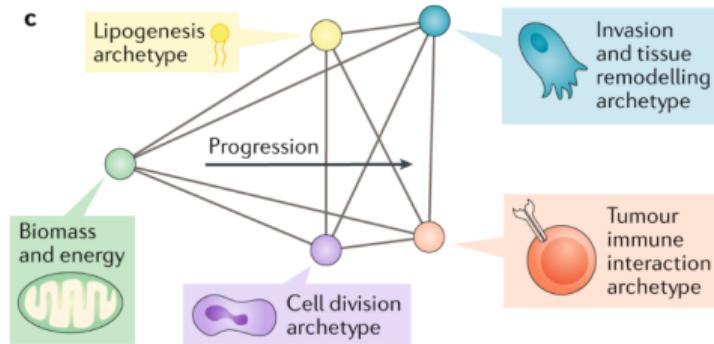
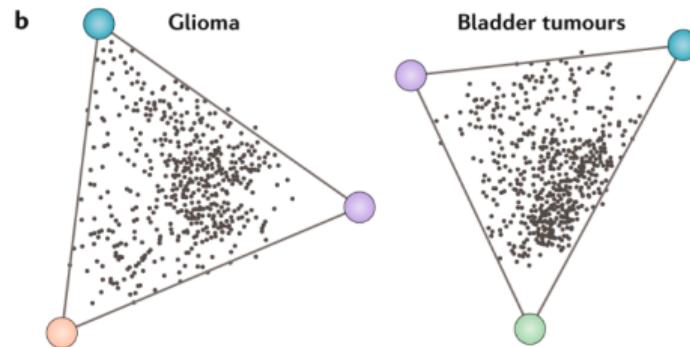
Hausser *et al.* (2019): Method

- Create a low-dimensional representation of the data (PCA)
- Fit a simplex and identifying vertices (archetypes)
- Interpret archetypes through gene-set enrichment analyses



An equilibrium deterministic mathematical model

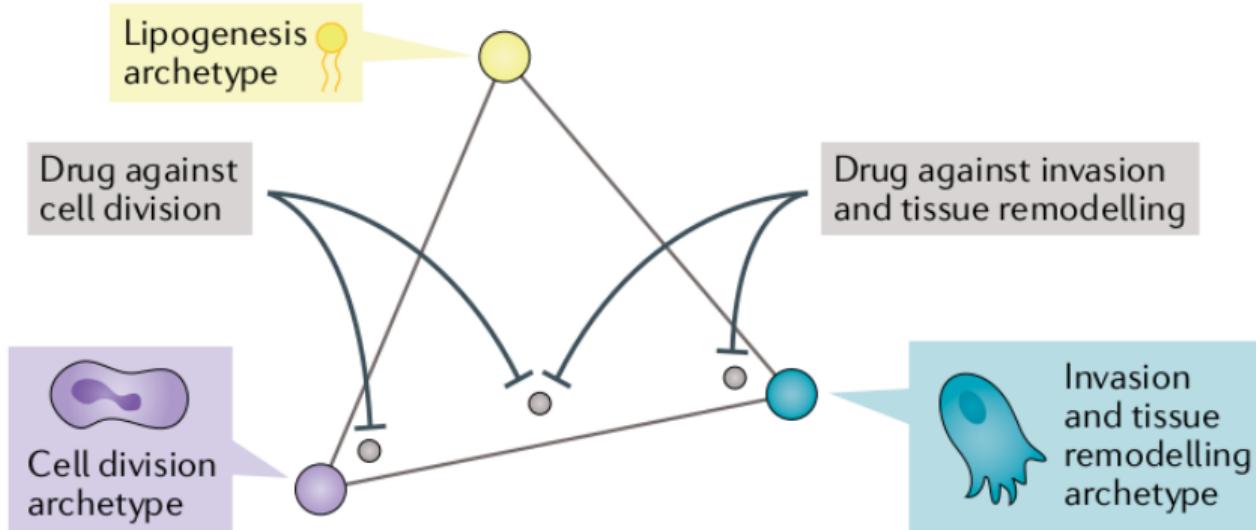
Hausser *et al.* (2019): Results



An equilibrium deterministic mathematical model

Hausser *et al.* (2019): Clinical implications

f



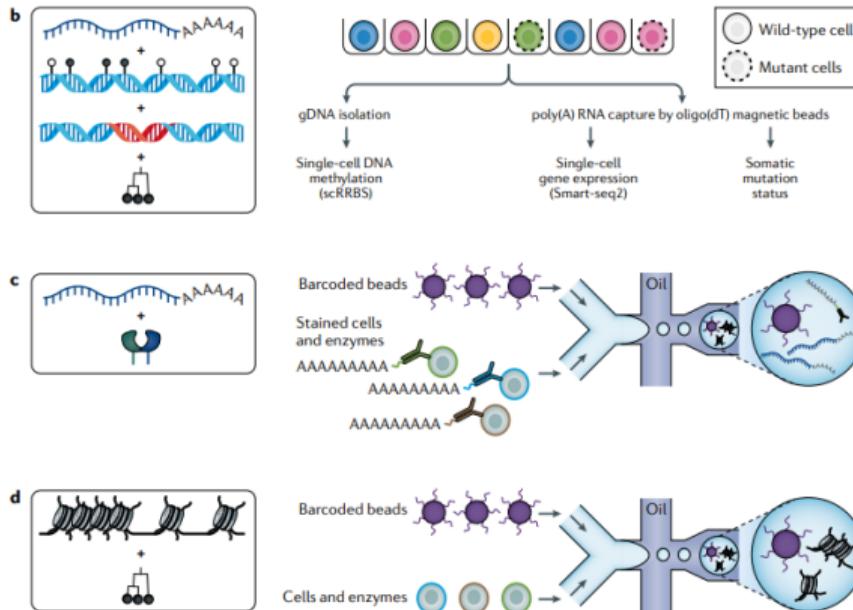
Take-home message

Interpret the geometry of the data

- 1 **Look at the molecular variation** (preferably after applying some linear dimensionality reduction)
- 2 Think about the **constraints in the data**: which genes are never overexpressed simultaneously? what does it say about the trade-offs between biological functions?
- 3 **Interpret the vertices and edges** of the distribution: do they have clear molecular phenotypes corresponding to pure tumoral strategies?

Part III: cancer evolution at the single cell resolution

- Evolutionary theory describes how individuals (here cells) compete and cooperate within populations (here tumor)
- Most evolutionary theory methods require individual-level genetic information, which is finally available for cell populations

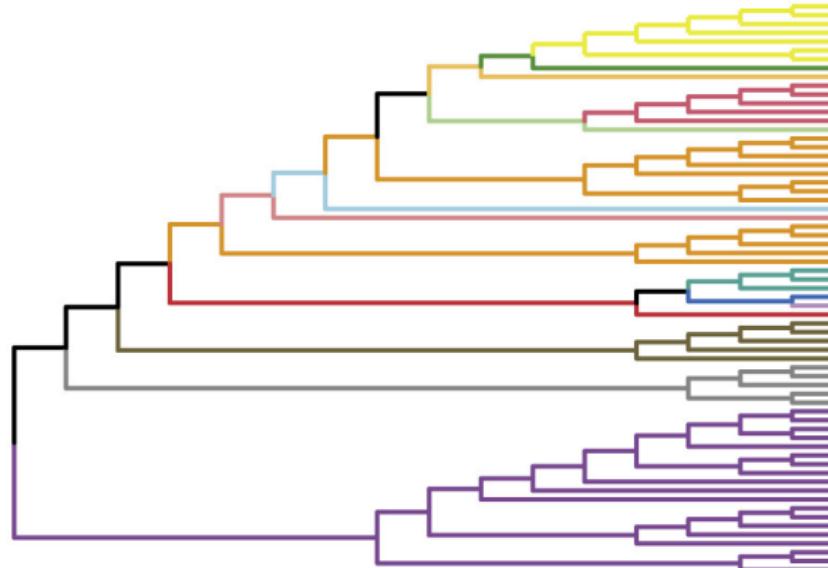


What is the main quantity of interest?

The cancer cell genealogy

Definition: ancestral relationships between tumor cells

Typical representation: genealogical tree



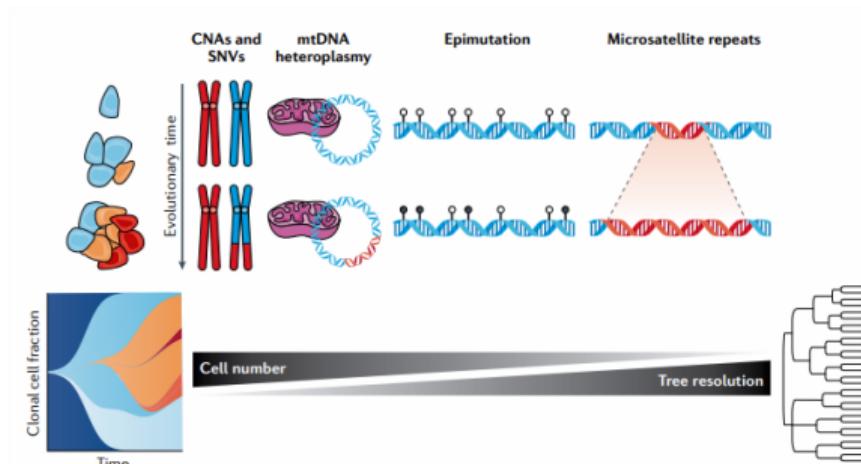
How to reconstruct cell genealogies?

Using genetic markers as natural barcodes (Nam *et al.* 2021)

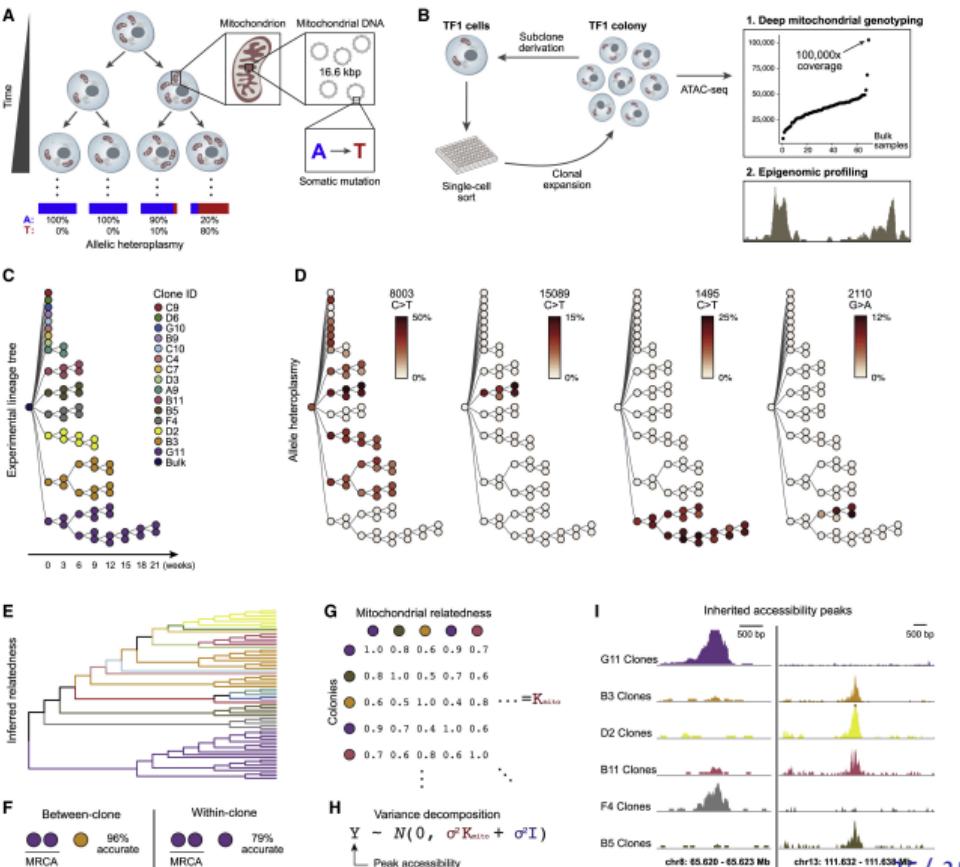
scWGS: use somatic CNVs, indels in microsatellite repeats, and a finite sites mutation model to resolve phylogeny (SNVs not reliable enough due to sparsity and sequencing artefacts)

scRNA-seq/ATAC-seq: use somatic mtDNA SNVs to resolve phylogeny

scDNAm: use DNA methylation as molecular clock to infer phylogeny



How to reconstruct cell genealogies?



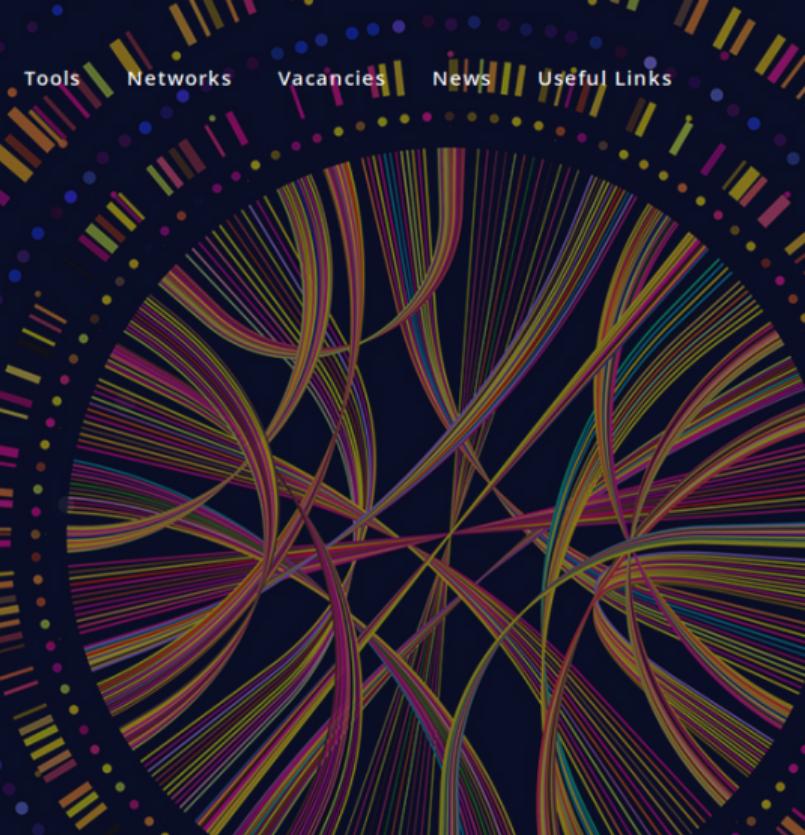
Using genetic markers as natural barcodes

scRNA-seq/ATAC-seq: use somatic mtDNA SNVs to resolve phylogeny (Ludwig *et al.* 2019)

Contact Us

Rare Cancers Genomics

Multidisciplinary and multi-omics molecular
characterisation of rare cancers



References

- Caravagna, G., T. Heide, M. J. Williams, L. Zapata, D. Nichol, *et al.*, 2020 Subclonal reconstruction of tumors by using machine learning and population genetics. *Nature Genetics* **52**: 898–907.
- Hausser, J., and U. Alon, 2020 Tumour heterogeneity and the evolutionary trade-offs of cancer. *Nature Reviews Cancer* : 1–11.
- Hausser, J., P. Szekely, N. Bar, A. Zimmer, H. Sheftel, *et al.*, 2019 Tumor diversity and the trade-off between universal cancer tasks. *Nature communications* **10**: 1–13.
- Kessler, D. A., and H. Levine, 2013 Large population solution of the stochastic luria–delbrück evolution model. *Proceedings of the National Academy of Sciences* **110**: 11682–11687.
- Ludwig, L. S., C. A. Lareau, J. C. Ulirsch, E. Christian, C. Muus, *et al.*, 2019 Lineage tracing in humans enabled by mitochondrial mutations and single-cell genomics. *Cell* **176**: 1325–1339.
- Luria, S. E., and M. Delbrück, 1943 Mutations of bacteria from virus sensitivity to virus resistance. *Genetics* **28**: 491.
- Nam, A. S., R. Chaligne, and D. A. Landau, 2021 Integrating genetic and non-genetic determinants of cancer evolution by single-cell multi-omics. *Nature Reviews Genetics* **22**: 3–18.
- Sun, R., Z. Hu, A. Sottoriva, T. A. Graham, A. Harpak, *et al.*, 2017 Between-region genetic divergence reflects the mode and tempo of tumor evolution. *Nature genetics* **49**: 1015.
- Tarabichi, M., I. Martincorena, M. Gerstung, A. M. Leroi, F. Markowetz, *et al.*, 2018 Neutral tumor evolution? *Nature genetics* **50**: 1630.
- Williams, M. J., A. Sottoriva, and T. A. Graham, 2019 Measuring clonal evolution in cancer with genomics. *Annual review of genomics and human genetics* **20**.
- Williams, M. J., B. Werner, C. P. Barnes, T. n. Graham, and A. Sottoriva, 2016 Identification of neutral tumor evolution across cancer types. *Nature genetics* **48**: 238.