



Bayesian Nonparametric Statistics for Fostering Innovation and Discovery in Biomedical Research

Yanxun Xu

Department of Applied Math and Statistics
Mathematics Institute for Data Science
Division of Biostatistics and Bioinformatics
The Sidney Kimmel Comprehensive Cancer Center
Johns Hopkins University

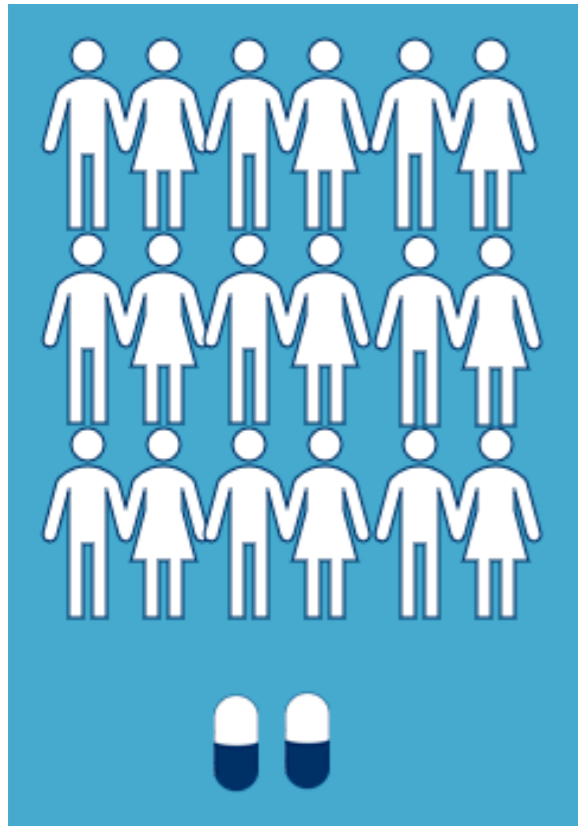
April, BNP 2022

Outline

- **Part 1: *Monday***
 - **Density estimation** for efficient clinical trial designs
 - **Regression** for precision dosing
- **Part 2: *Wednesday***
 - **Clustering** for subgroup finding
 - **Latent feature models** for tumor heterogeneity
- **Part 3: *Friday***
 - Estimating treatment effects from observational data

Treatment Effect

Clinical Trials



Biomarker-based Trials



Observational Data

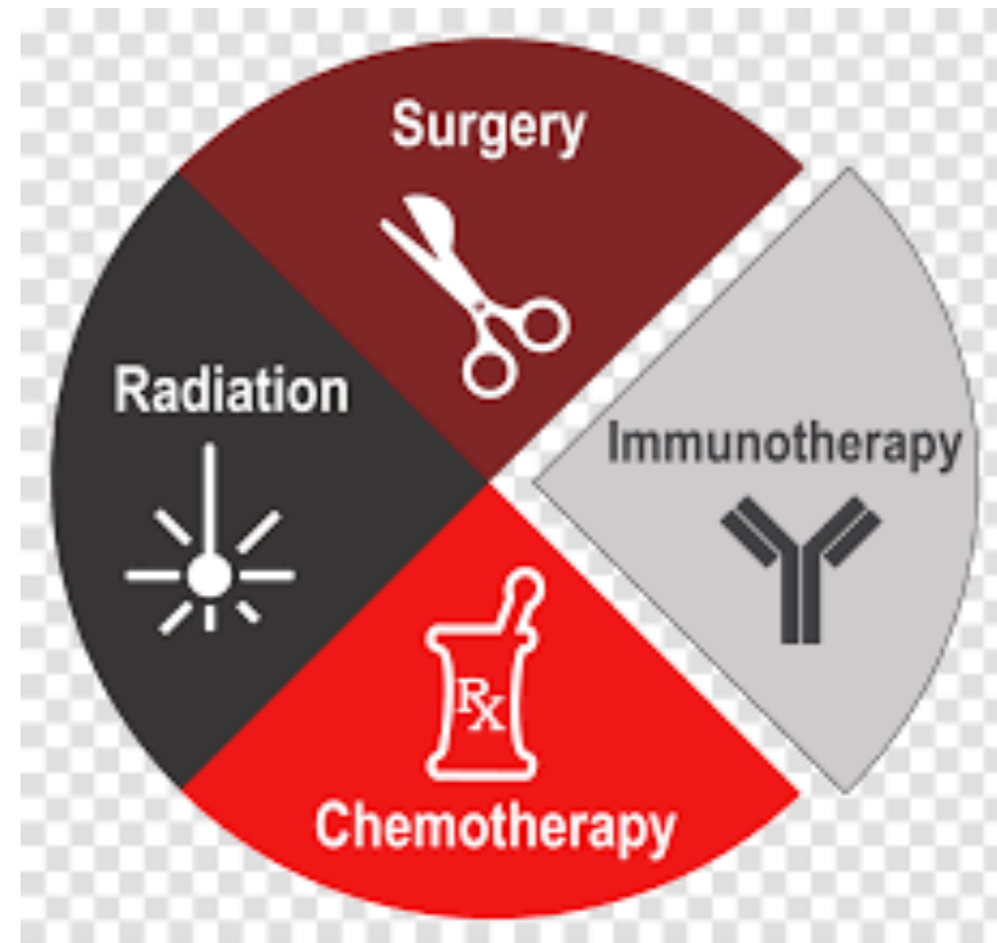


Population

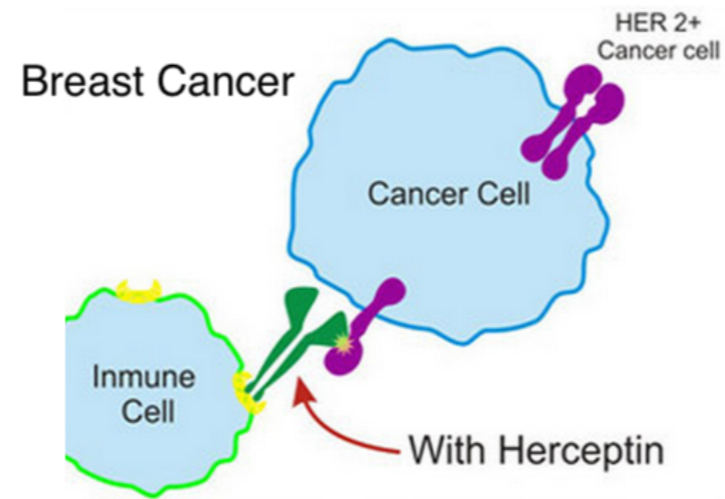
Subgroup

Personalization

One-size Fits All Cancer Treatment

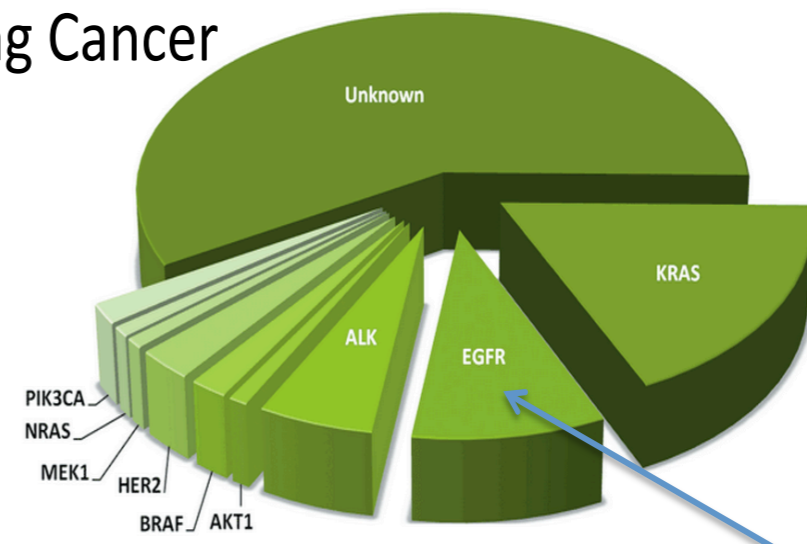


Targeted Therapy



© 2010 Leonor García del Valle. All rights reserved

Lung Cancer



Courtesy of Lung Cancer Foundation of America

Erlotinib

Genomic-driven Cancer Trials

Umbrella Trials

in one **single** cancer type, test the effect of targeted agents on different alterations.

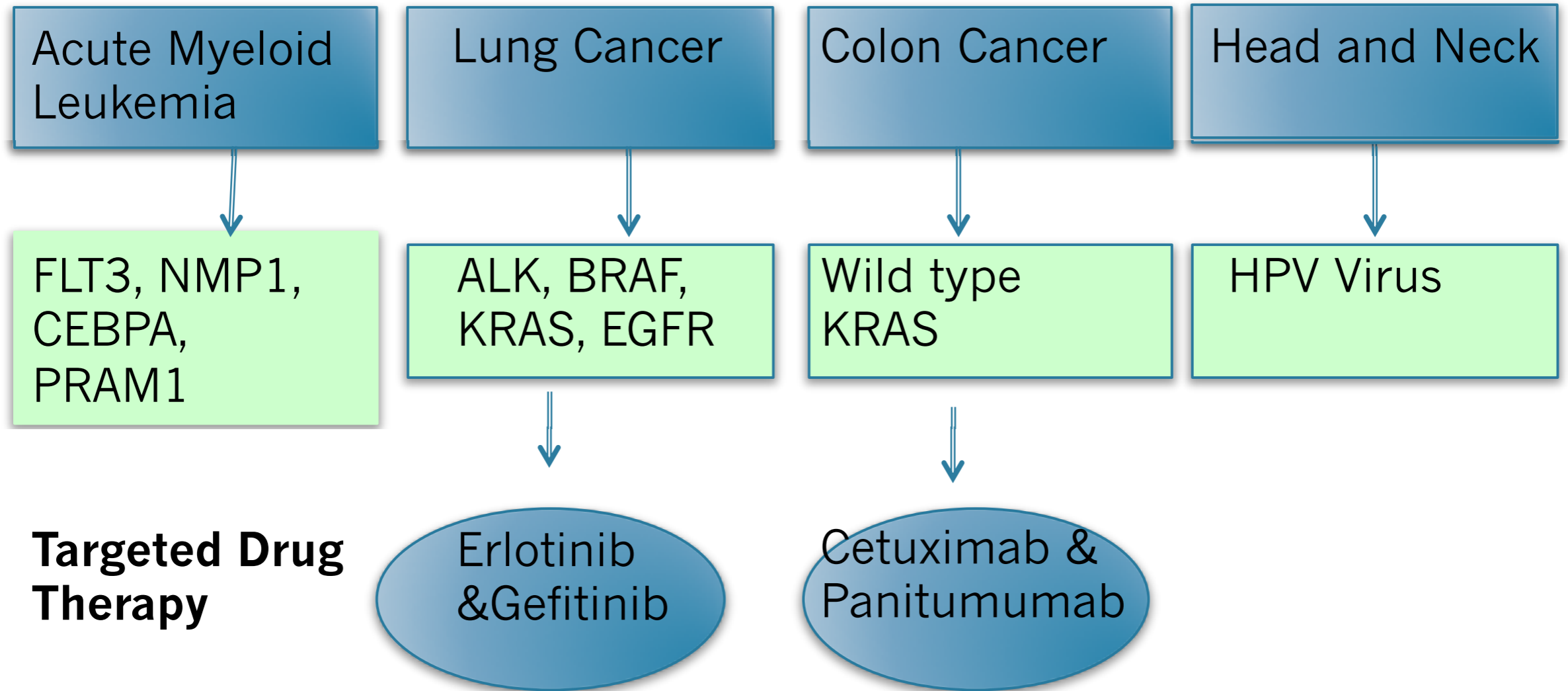


Basket Trials

across **multiple** cancer types, test the effect of targeted agents on the same genomic alterations.



Basket Trial



Motivation Trial: IMPACT II

- **Clinical Trial:** study of targeted agents in metastatic cancers.
- **Patients:** with metastatic cancer (thyroid, ovarian, melanoma, lung, breast, CRC and other)
- **Treatments:** therapy that targets particular molecular aberrations (TT) vs. standard of care (S)
- **Population:** heterogeneous population; different mutations; different cancers; baseline covs . . . Treatment might be effective in a sub-population

Data:

TRT	TUMOR	PFS	CENS	MUTATIONS					
				m1	m2	m3	m4	m5	m6
TT	THYROID	2.6	0	NA	NA	NA	NA	NA	NA
TT	THYROID	3.6	0	NA	0	0	0	NA	0
S	OVARIAN	4.2	1	0	NA	0	0	0	0
S	MELANOMA	5.8	1	NA	0	0	0	NA	0
...						

Motivation Trial: IMPACT II

Objective: determine the subpopulation that achieves the maximum benefit from TT.

	EGFR	KRAS	TP53
Lung Cancer			
Colon Cancer			

We will cast this goal as a **decision problem**.

Subpopulation Finding: Decision Problem

- **Outcome:** progression free survival (PFS) time, $y_i, i = 1, \dots, n$
- **Action:** report a subgroup of patients who might benefit from the TT. A set of **mutation-tumor pairs**,

$$A = \{a : a = (j_a, c_a)\}$$

- $j_a = \{1, \dots, q\}$: Molecular aberration
- $c_a \in \{1, \dots, n_c\}$: tumor type

$\{(KRAS, Lung), (TP53, Breast)\}$

Subpopulation Finding: Decision Problem

- **Action:** report a subgroup of patients who might benefit from the TT. A set of **mutation-tumor pairs**,

$$A = \{a : a = (j_a, c_a)\}$$

Bayes Rule: $A^* = \operatorname{argmax}_A \int u(A, \theta) p(\theta | y, X) d\theta$

Utility: we favor a subpopulation with difference in **log hazards ratio (LR)** and **large size**

Data from IMPACT

- **Outcome:** progression free survival times, y_i
- **Covariates:** $x_i = (c_i, m_i, b_i)$
 - Tumor type c_i (categorical)
 - Molecular aberrations $m_i = (m_{i1}, \dots, m_{iM})$ (binary)
 - Other baseline covariates b_i (age, # prior therapies, etc)

Challenges

Probability model needs to allow for:

- interactions of covariates
- heterogeneous population
- missing data
- Extrapolation with small # observations

BNP!

Random Partition

- $s = (s_1, \dots, s_n)$ be cluster membership indicators,
 $s_i \in \{1, \dots, J\}$
- $S_j = \{i : s_i = j\}$

Product partition model: $p(s) \propto \prod_{j=1}^J c(S_j)$

For DP, $c(S_j) = \alpha(|S_j| - 1)!$

Random Partition

- $s = (s_1, \dots, s_n)$ be cluster membership indicators,
 $s_i \in \{1, \dots, J\}$
- $S_j = (i : s_i = j)$
- x_j^* by cluster

Product partition model with covariates (PPMx):

$$p(s \mid x) \propto \prod_{j=1}^J c(S_j) g(x_j^*)$$

Favors clusters homogeneous in x_i with $g(x_j^*)$ scoring similarity of $x_j^* = \{i : s_i = j\}$.

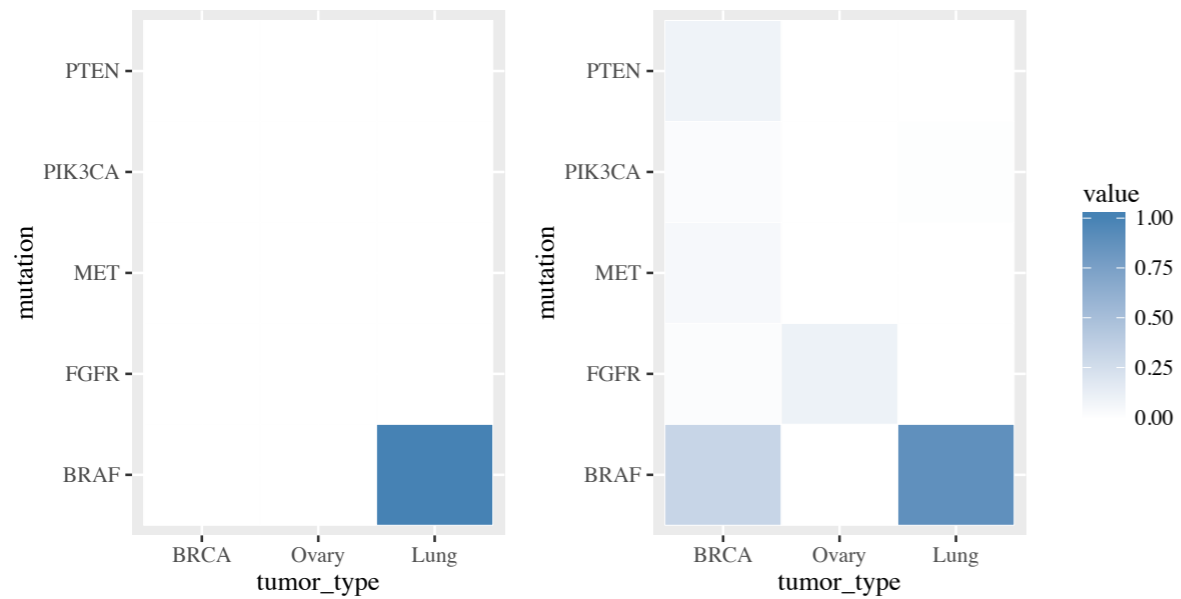
Similarity function: over observed covariates only

$$g(x_j^*) = \prod_{l=1}^p g_l(\{x_{il}, i \in S_j \text{ and } x_{il} \text{ observed}\})$$

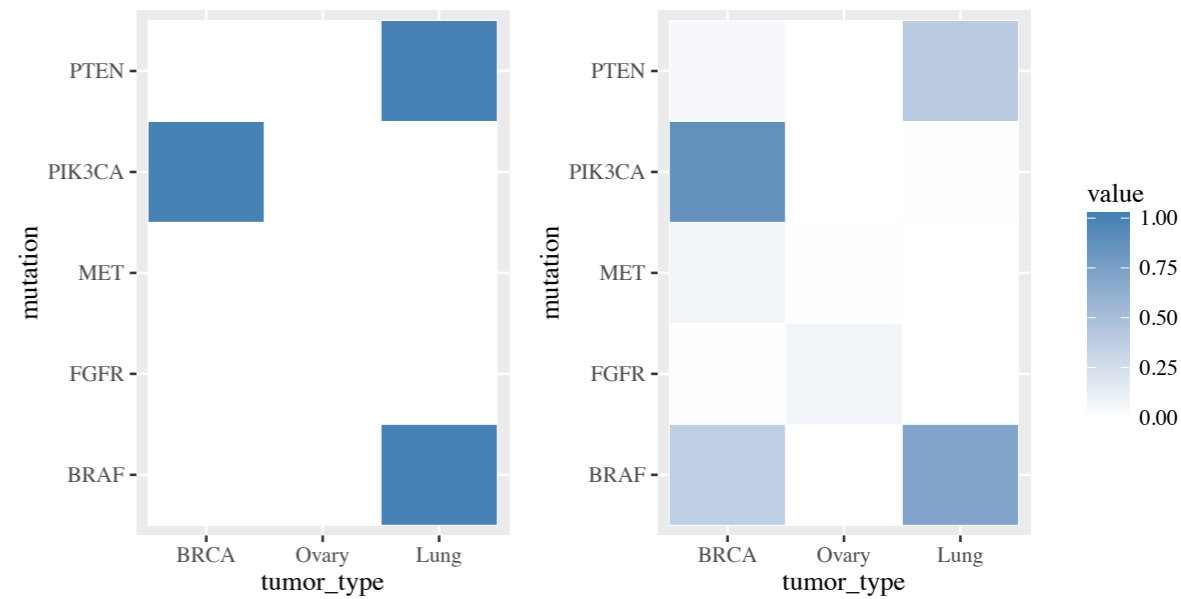
Sampling model: exchangeable within clusters (e.g., lognormal regression model)

$$p(y \mid s, x, \eta) = \prod_{j=1}^J \prod_{i \in S_j} p(y_i \mid \eta_j)$$

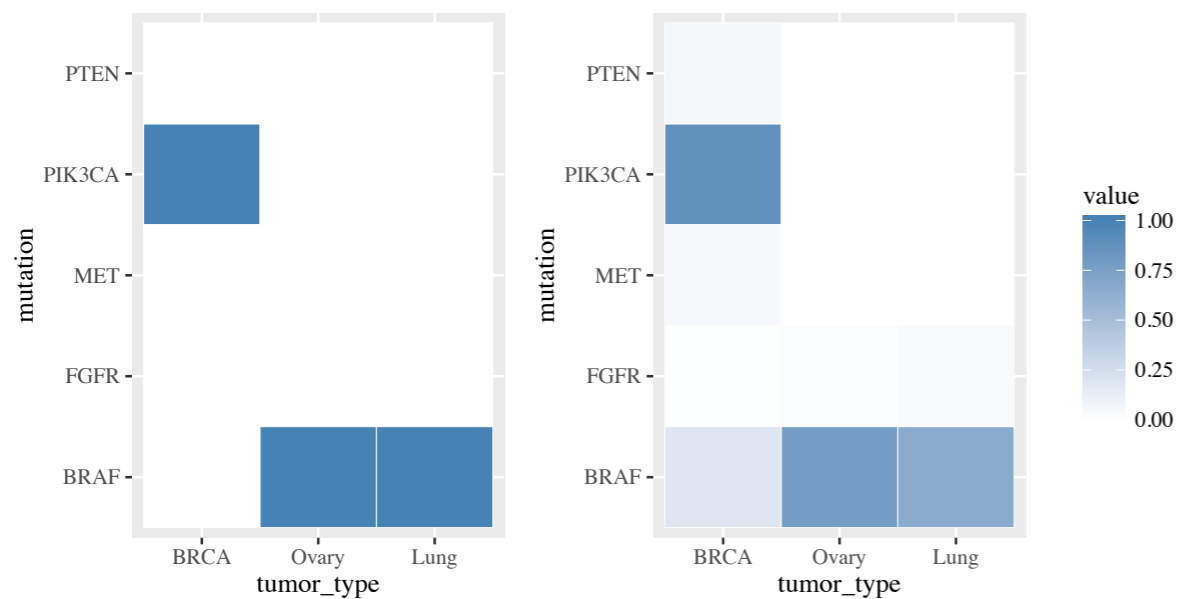
Results



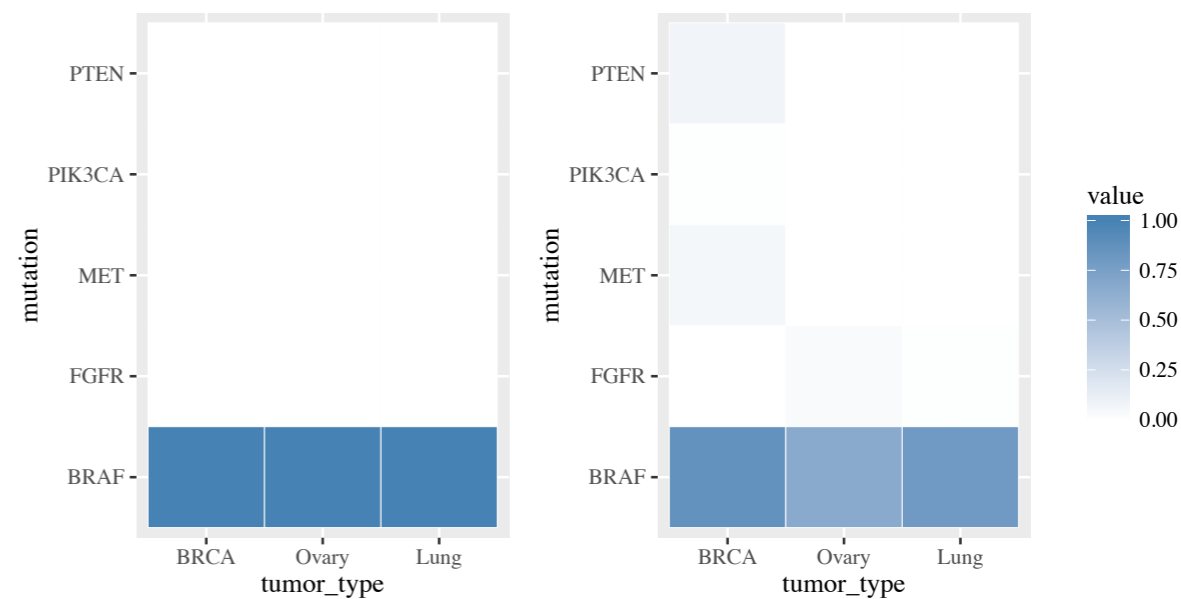
Scenario 3



Scenario 4



Scenario 5



Scenario 6

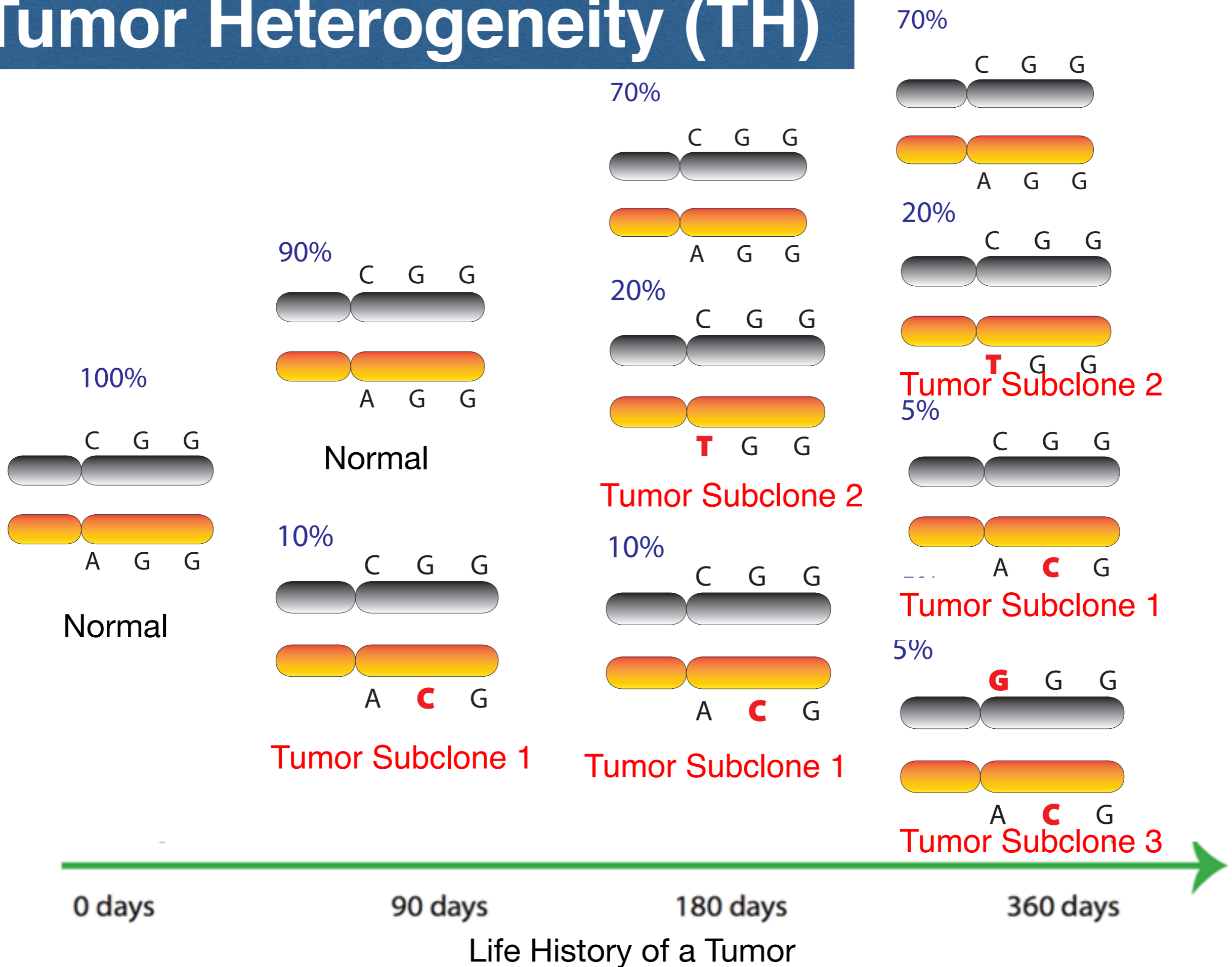
Take-home Messages

- A general class of probability models that allow for interactions and missing data
- Subgroup finding can be casted as a decision problem.
- Separate the decision problem with probability model
- Can be used in clinical trial designs to adaptively assign patients

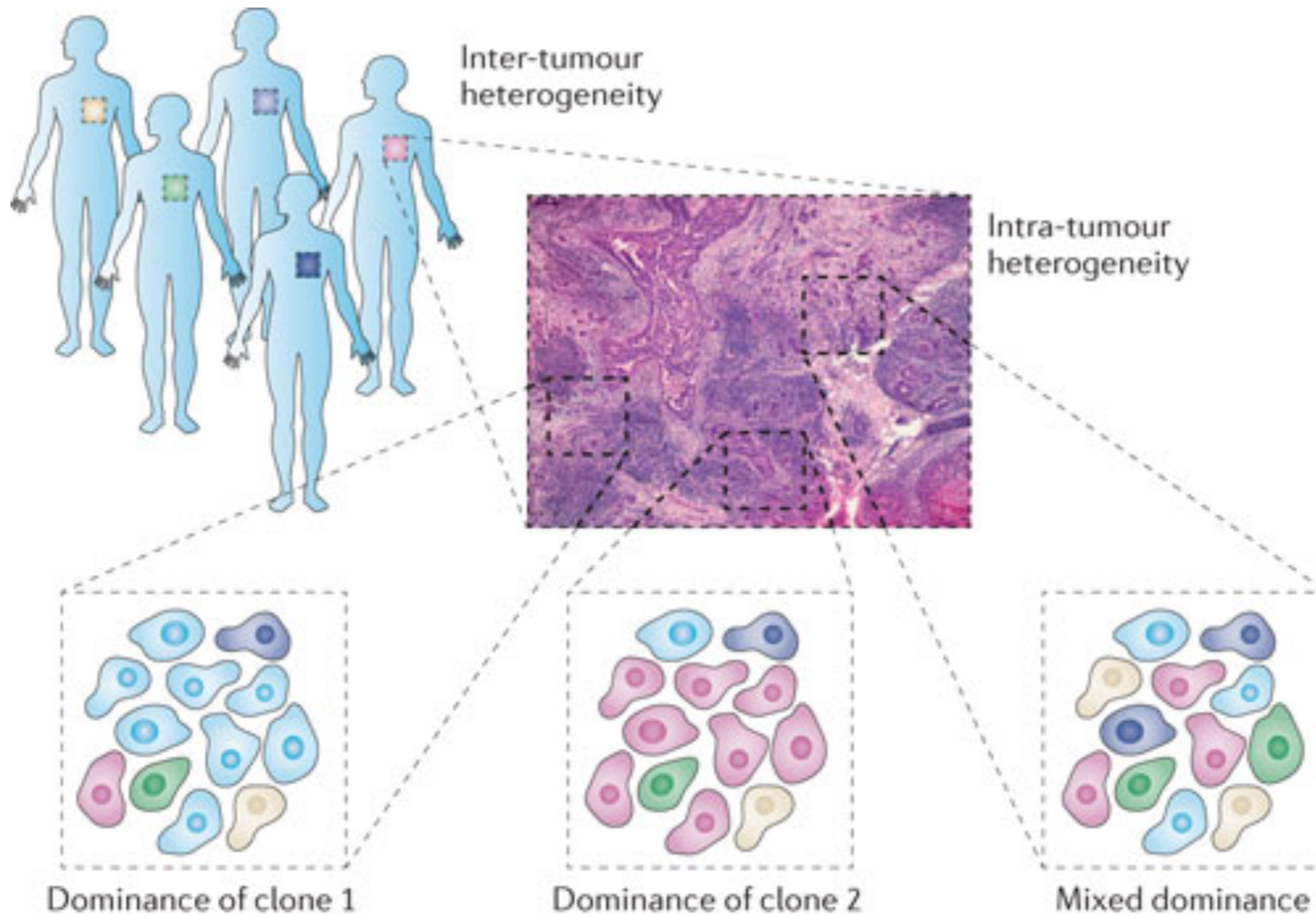
Outline

- **Part 1: *Monday***
 - **Density estimation** for efficient clinical trial designs
 - **Regression** for precision dosing
- **Part 2: *Wednesday***
 - **Clustering** for subgroup finding
 - **Latent feature models** for tumor heterogeneity
- **Part 3: *Friday***
 - **Estimating treatment effects** from observational data

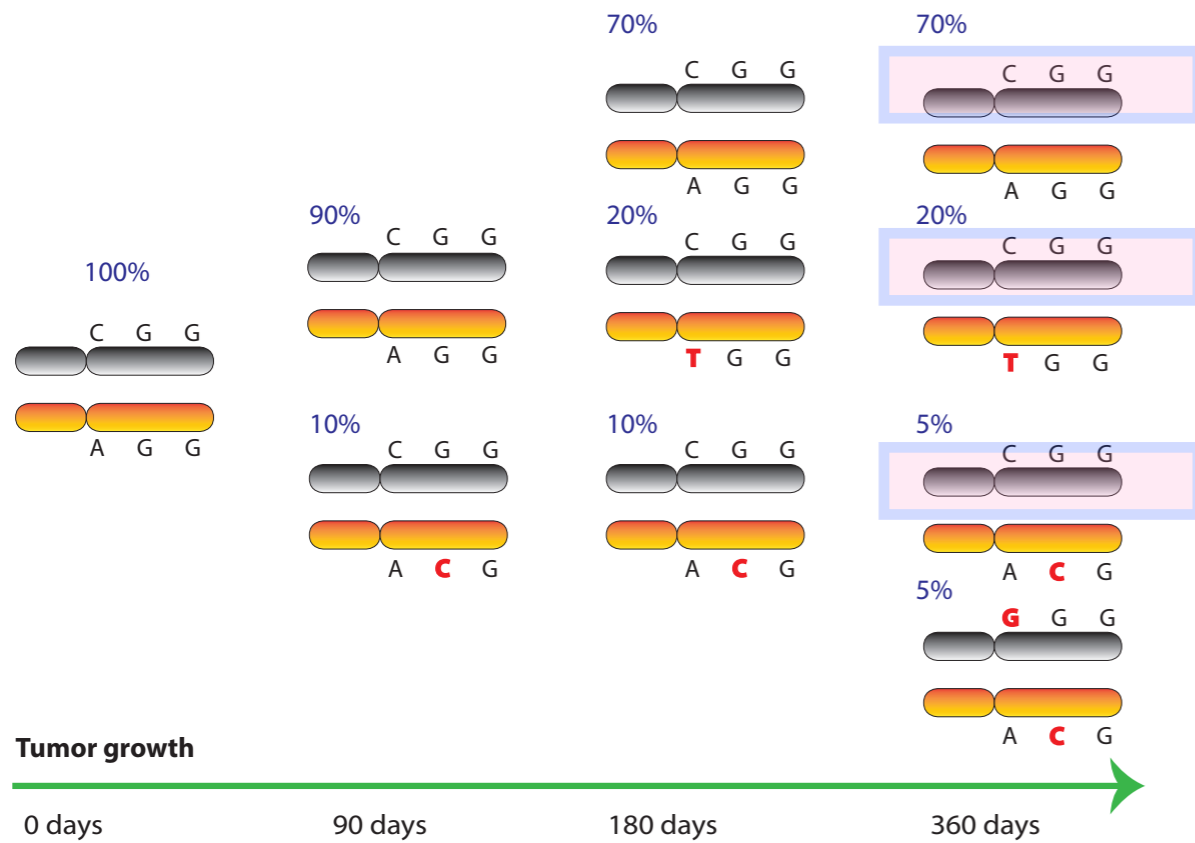
Tumor Heterogeneity (TH)



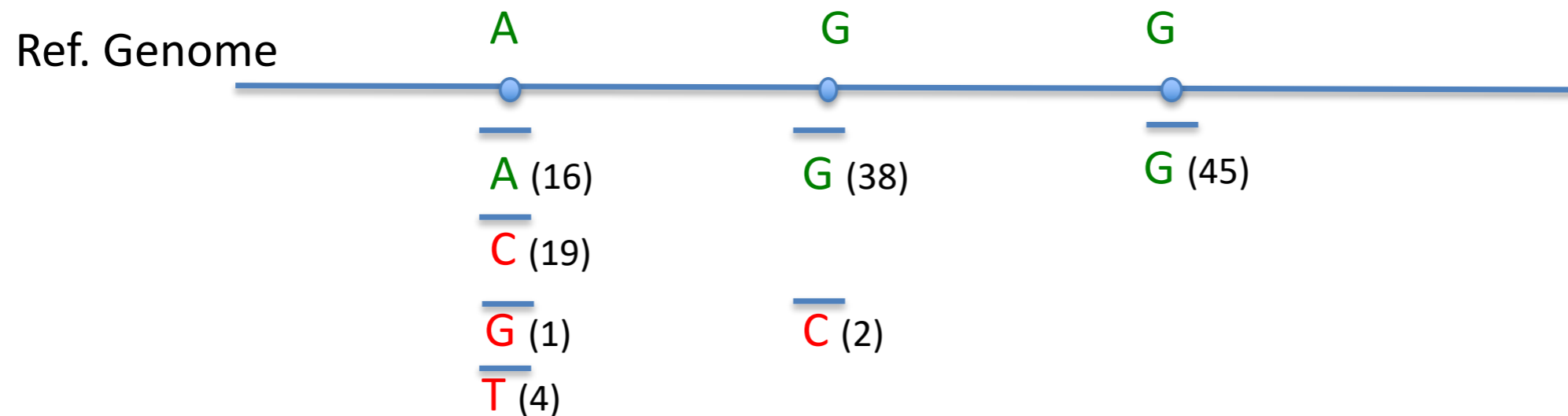
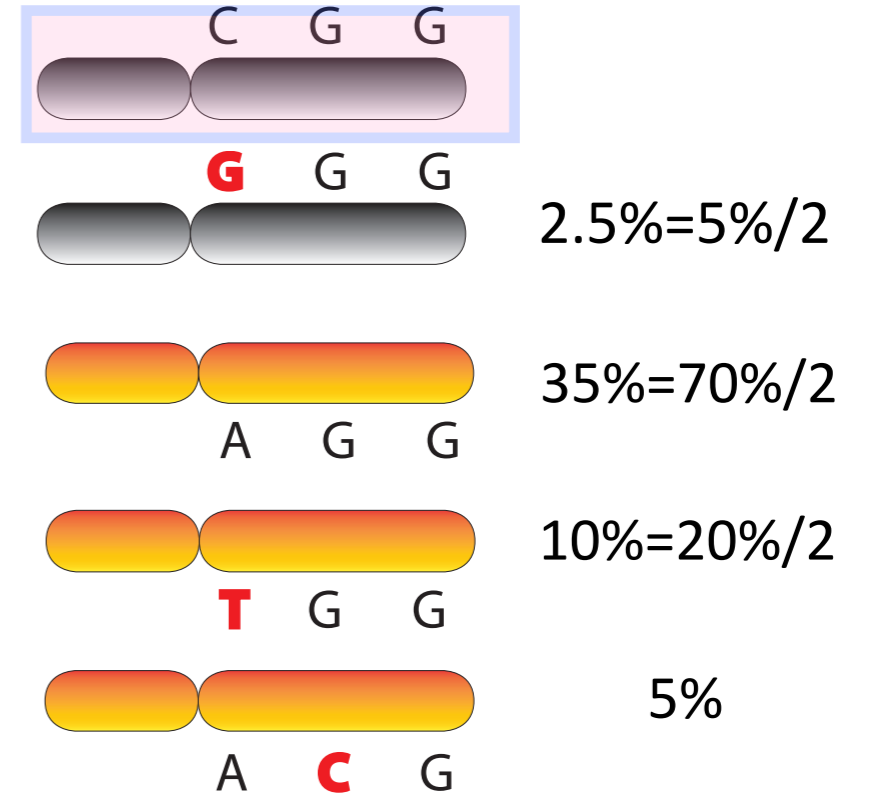
Clinical Utility of TH



Haplotype

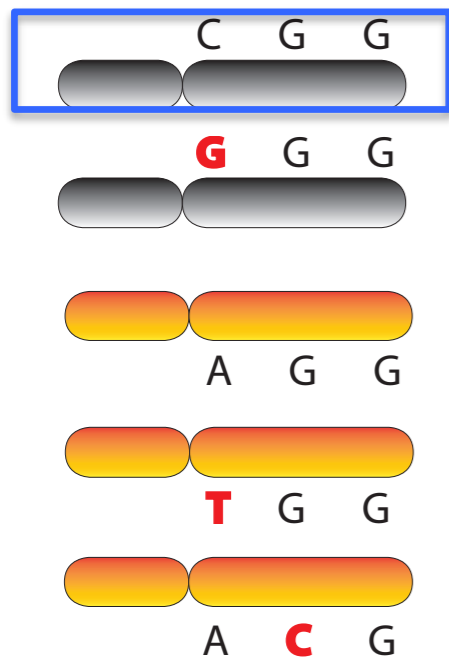


$$47.5\% = (70\% + 20\% + 5\%) / 2$$



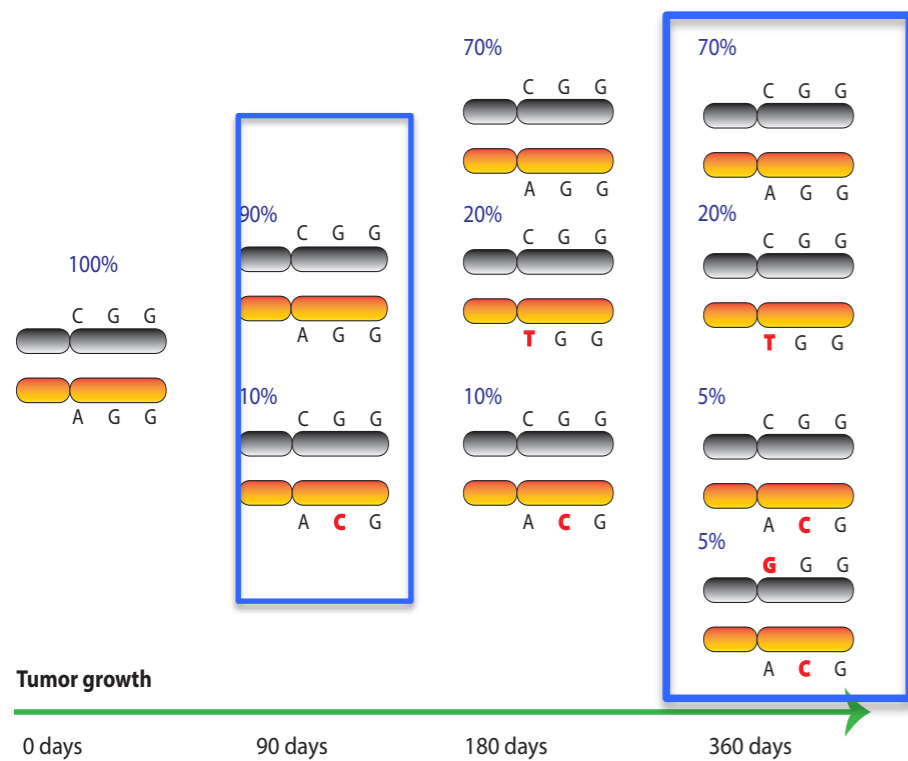
$$47.5\% \text{ (CGG)} + 2.5 \text{ (GGG)} + 35\% \text{ (AGG)} + 10\% \text{ (TGG)} + 5\% \text{ (ACG)}$$

Tumor Heterogeneity in Terms of Haplotype Genome (Z) and Cellular Fractions (W)



	Hap1	Hap2	Hap3	Hap4	Hap5
SNV1	1	1	0	1	0
SNV2	0	0	?	0	1
SNV3	0	0	?	0	0

The Z Matrix

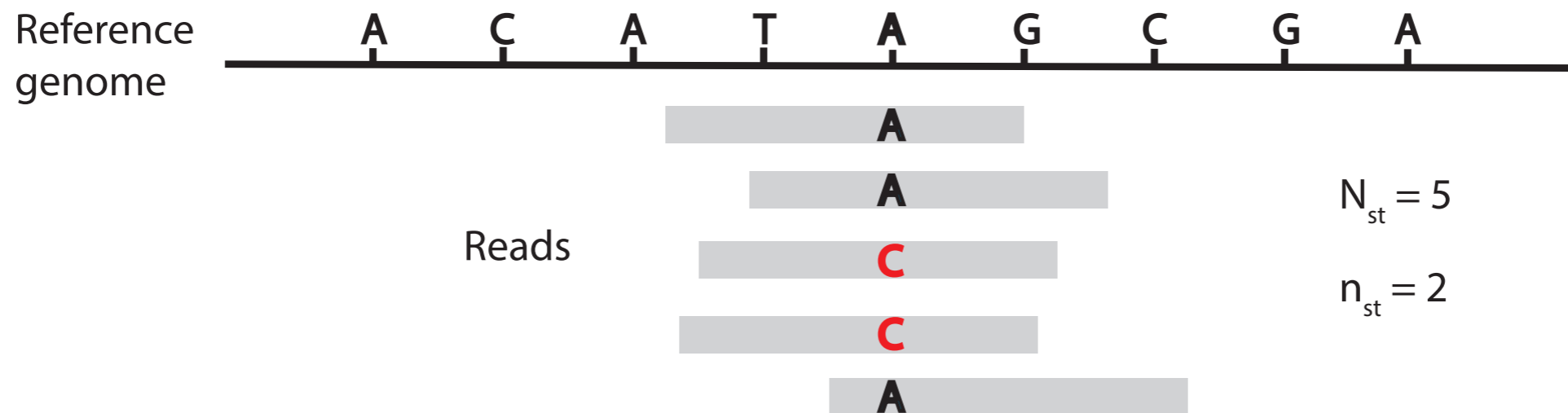


	Hap1	Hap2	Hap3	Hap4	Hap5
Sample 1	47.5%	2.5%	?	10%	5%
Sample 2	50%	0	35	10%	5%
Sample 3	50%	0	?	0	5%
Sample 4	50%	0	?	0	0

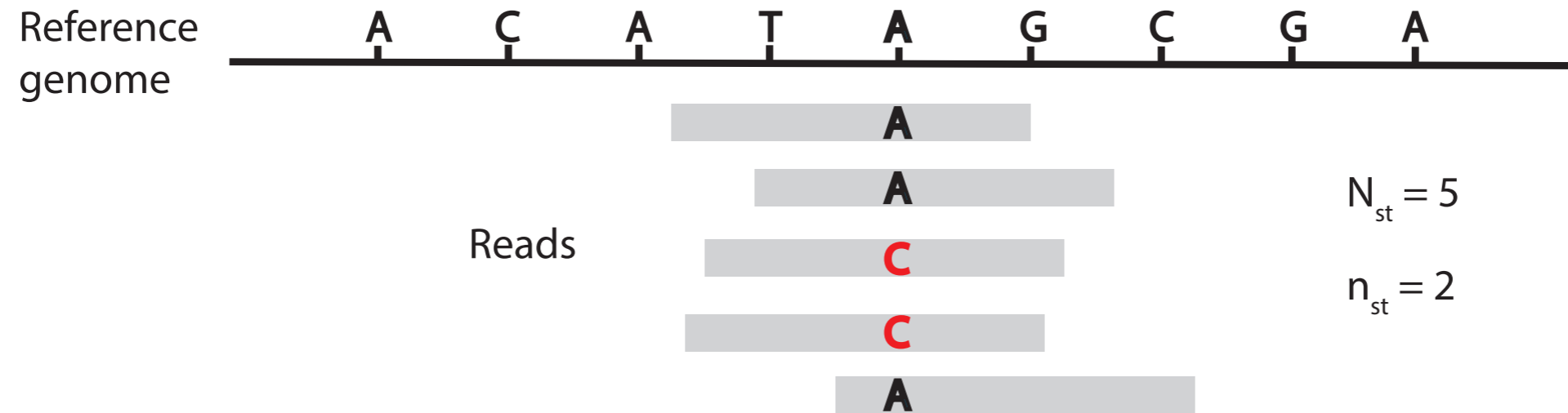
The W Matrix

Notations

- **SNV:** point mutations, $s = 1, \dots, S$
- **Sample:** $t = 1, \dots, T$
- **Data:** $N_{st} = \#$ reads mapped to locus of SNV s in sample t
 $n_{st} = \#$ of them with SNV.



Sampling Model



$$n_{st} \sim \text{Binomial}(N_{st}, p_{st})$$

VAF: variant allele fraction

Observed VAF: n_{st}/N_{st}

Expected VAF: $p_{st} = E(n_{st}/N_{st})$

Link VAFs with Haplotypes



Observed VAF: n_{st}/N_{st}

Expected VAF: $p_{st} = E(n_{st}/N_{st})$

Key Idea: A **variant** read must be from a haplotype with variant.

Link VAFs with Haplotypes

Key Idea: A **variant** read must be from a haplotype with variant.

s : SNV; c : haplotype (latent); t : sample

$z_{sc} = 1$: haplotype c has a variant on SNV s .

$z_{sc} = 0$: haplotype c has no variant on SNV s .

w_{tc} : fraction of haplotype c in sample t .

Linking Equation:

$$p_{st} = \sum_c w_{tc} z_{sc}$$

Haplotype Genotype Z

Haplotypes

	1	2	3	...	C
1	1	1	0	0	0
2	1	1	1	0	0
3	0	1	1	0	0
4	1	0	0	0	0
...	1	0	0	1	1
...	0	1	0	1	1
S	1	1	0	0	0

$p(Z)$ on $(S \times C)$ binary matrix

Indian Buffet Process (IBP)

SNVs

Indian Buffet Process (IBP)

		Dishes				
		1	2	3	...	C
Customers	1	■	■			
	2	■	■	■		
	3		■	■		
	4	■				
	...	■			■	■
	...				■	■
	...		■		■	■
	S	■	■			

For $s = 1, \dots, S$

- Customer s chose dish c that has been already chosen m_k time with probability m_k/s
- Number of new dishes:
 $K_s \sim \text{Poisson}(\gamma/s)$

Indian Buffet Process (IBP)

		Dishes				
		1	2	3	...	C
Customers	1					
	2					
	3					
	4					
	⋮					
	S					

For $s = 1, \dots, S$

- Customer s chose dish c that has been already chosen m_k time with probability m_k/s
- Number of new dishes:
 $K_s \sim \text{Poisson}(\gamma/s)$

Indian Buffet Process (IBP)

		Dishes				
		1	2	3	...	C
Customers	1	■	■			
	2					
	3					
	4					
	⋮					
	S					

For $s = 1, \dots, S$

- Customer s chose dish c that has been already chosen m_k time with probability m_k/s
- Number of new dishes:
 $K_s \sim \text{Poisson}(\gamma/s)$

Indian Buffet Process (IBP)

		Dishes				
		1	2	3	...	C
Customers	1	■	■			
	2	■	■	■		
	3					
	4					
	⋮					
	S					

For $s = 1, \dots, S$

- Customer s chose dish c that has been already chosen m_k time with probability m_k/s
- Number of new dishes:
 $K_s \sim \text{Poisson}(\gamma/s)$

Indian Buffet Process (IBP)

		Dishes				
		1	2	3	...	C
Customers	1	■	■			
	2	■	■	■		
	3		■	■		
	4					
	⋮					
	S					

For $s = 1, \dots, S$

- Customer s chose dish c that has been already chosen m_k time with probability m_k/s
- Number of new dishes:
 $K_s \sim \text{Poisson}(\gamma/s)$

Indian Buffet Process (IBP)

		Dishes				
		1	2	3	...	C
Customers	1	■	■			
	2	■	■	■		
	3		■	■		
	4	■				
	...	■			■	■
	...		■		■	■
	...				■	■
	S	■	■			

For $s = 1, \dots, S$

- Customer s chose dish c that has been already chosen m_k time with probability m_k/s
- Number of new dishes:
 $K_s \sim \text{Poisson}(\gamma/s)$

IBP Prior

$p(\mathbf{Z} | \gamma) =$

	1	2	3	...	C
1					
2					
3					
4					
...					
S					

customers (SNVs) who chose dish c (haplotype)

Number of customers (SNVs)

Number of dishes (haplotypes)

$$p(\mathbf{Z} | \gamma) = \frac{\gamma^C e^{-\gamma H_S}}{C!} \prod_{c=1}^C \frac{(S - m_c)! (m_c - 1)!}{S!}$$

Model Summary

$$p(Z, w, n | N) = \underbrace{p(Z)}_{\text{IBP}} p(w | Z) \underbrace{p(n | Z, w, N)}_{\text{Binomial}}.$$

$$p_{st} = \sum_c w_{tc} z_{sc}$$

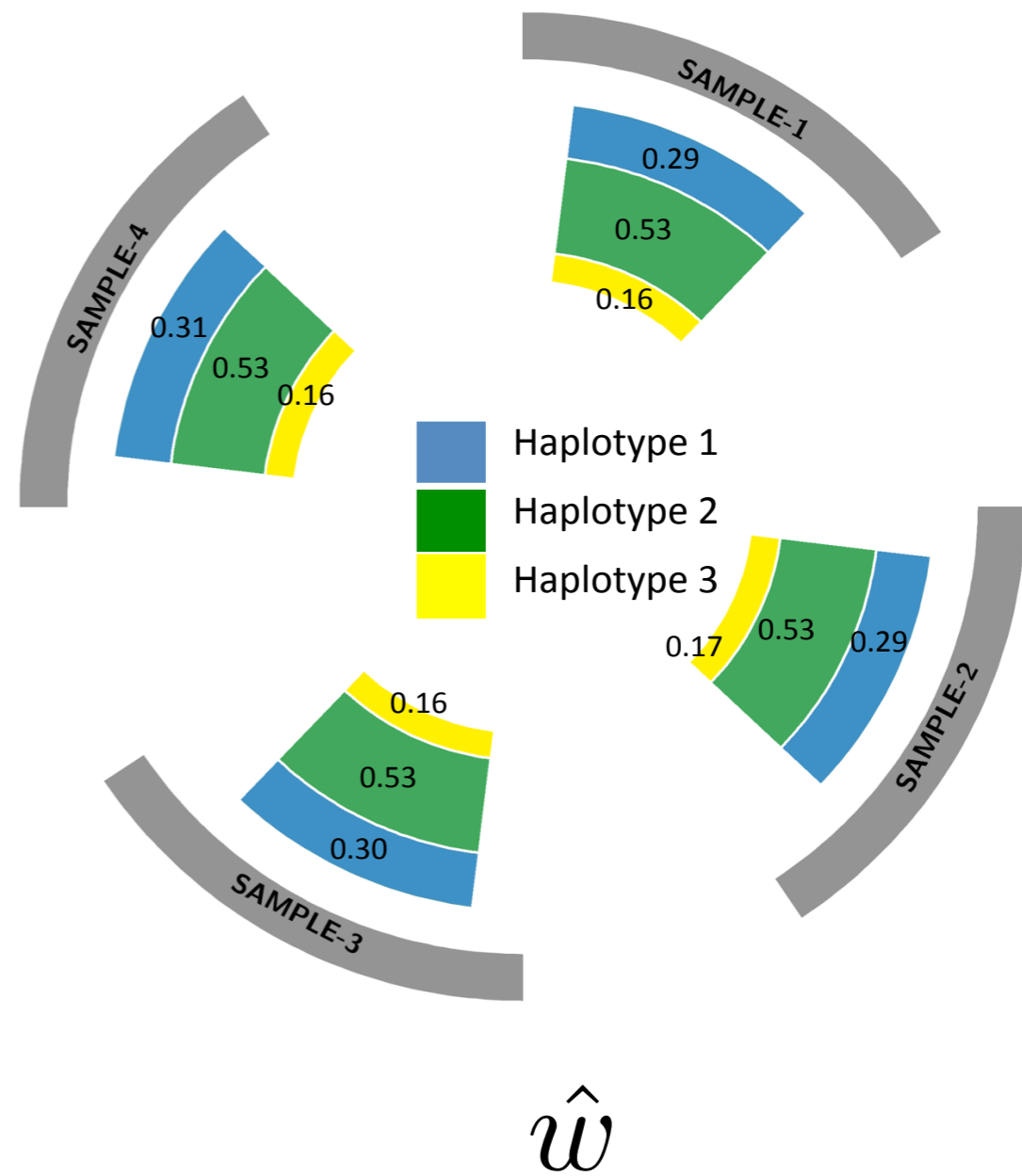
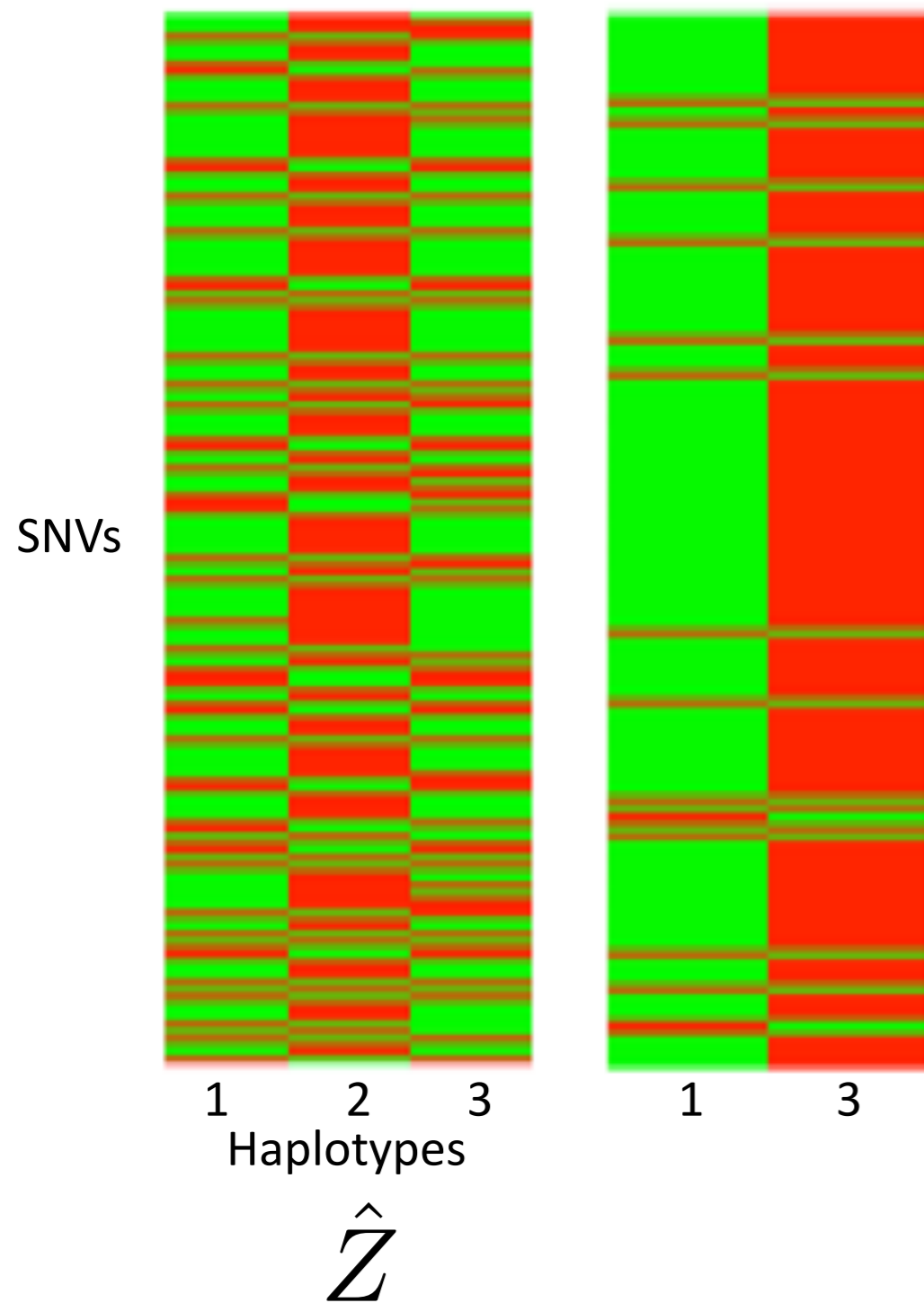
$$p(w_t) \sim \text{Dir}(a_1, \dots, a_C), t = 1, \dots, T.$$

$$p(Z, w | N, n)$$

Application: Intra-Tumor Heterogeneity

- One tumor from lung cancer; 4 samples surgically dissected
- Each sample generates a whole-genome sequencing data set
- Bio-X pipeline (BWA, Samtools, GATK) for data preprocessing: coverage ~ 100X.
- Selected $S=17,160$ SNVs

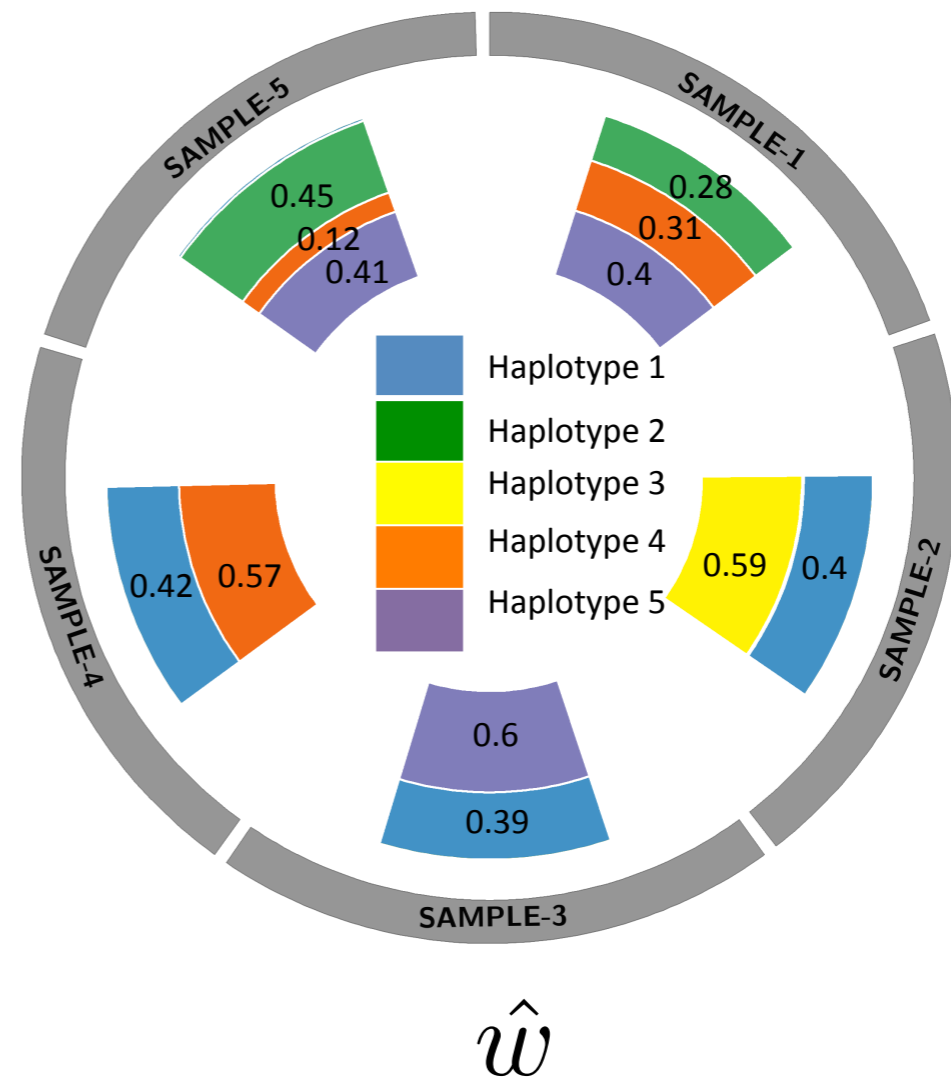
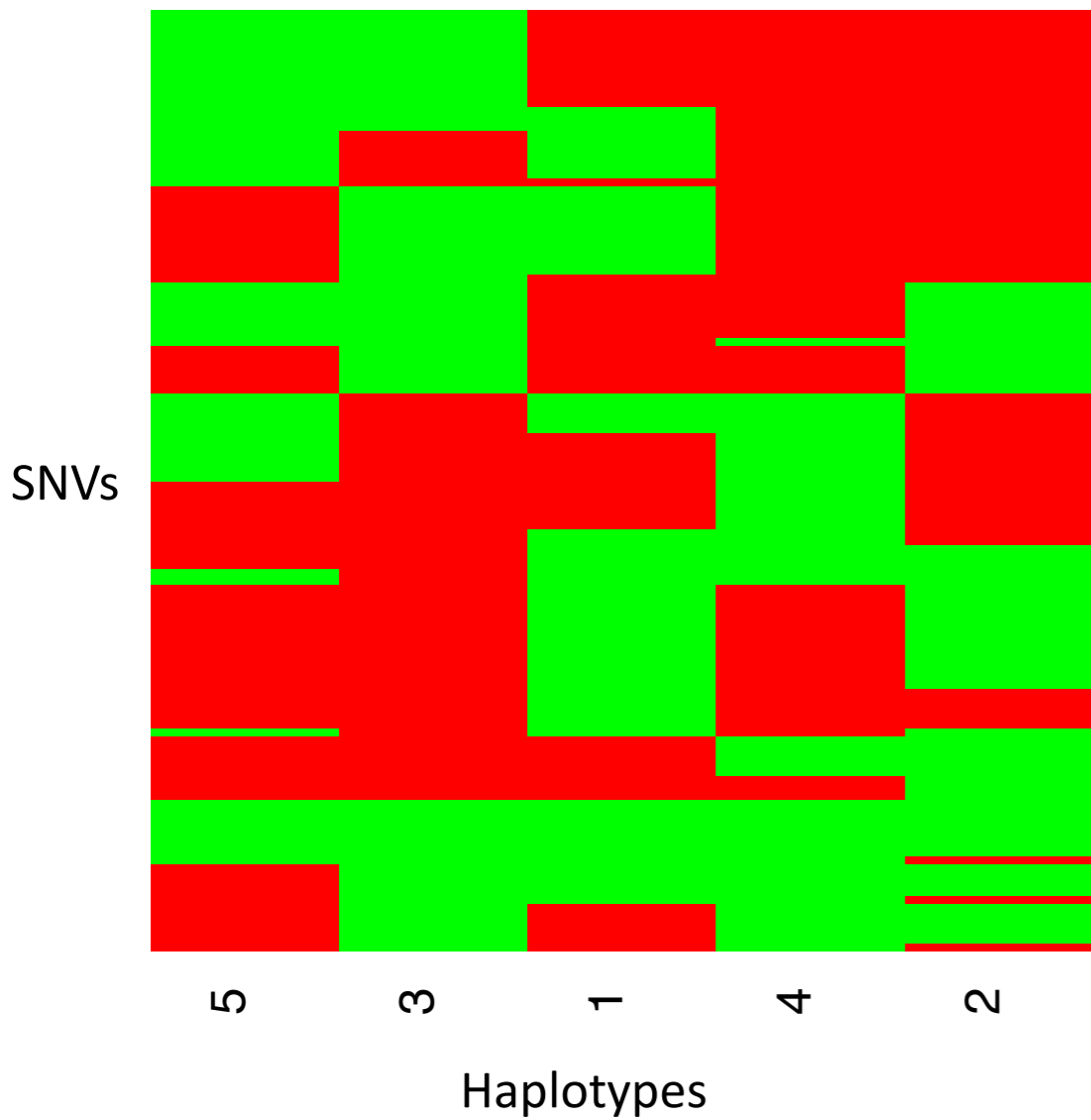
Application: Intra-Tumor Heterogeneity



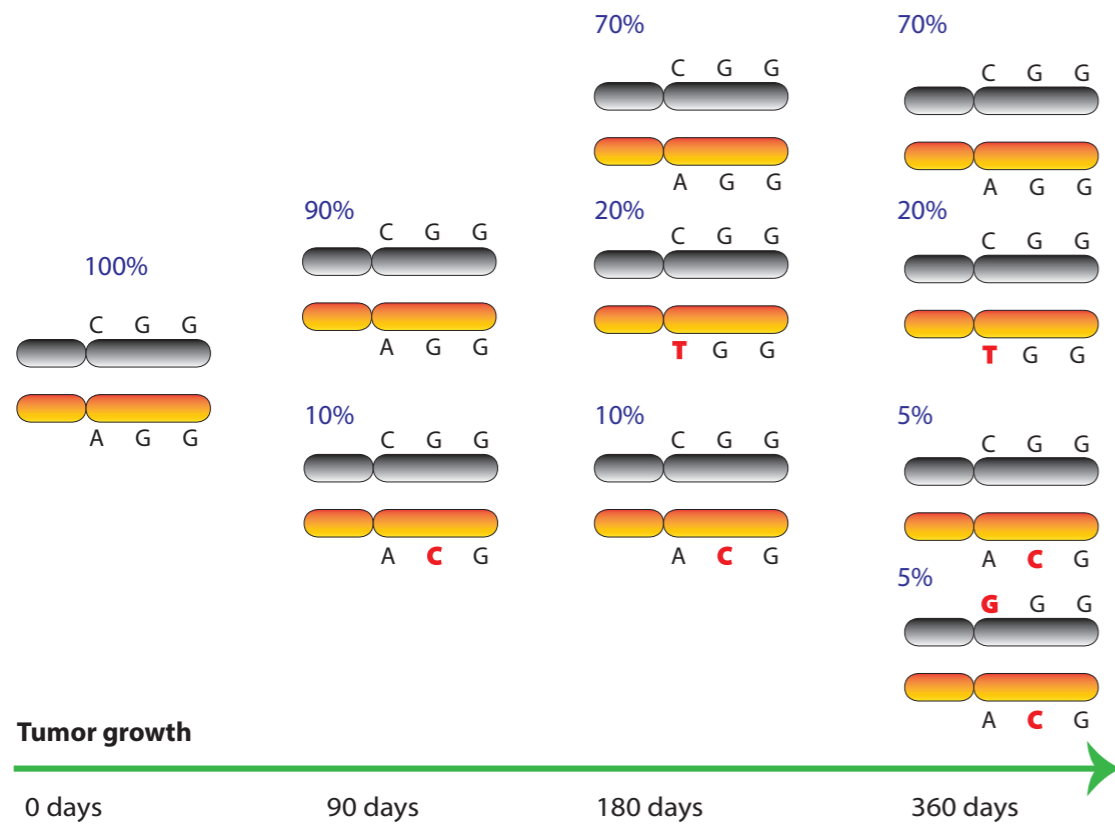
Application: Inter-Tumor Heterogeneity

- Exome-sequencing data for **five** tumor samples from four different pancreatic ductal adenocarcinoma (PDAC) patients
- Bio-X pipeline (BWA, Samtools, GATK) for data preprocessing: coverage $\sim 70X$.
- Selected **118** SNVs: 1) significant coverage in all samples; 2) related to PDAC in the KEGG pathway database; 3) are nonsynonymous

Application: Inter-Tumor Heterogeneity



Extension: Categorical IBP



Subclone

	1	2	3	4	5
1	0.5	1	0	1	0
2	1	0.5	1	1	1
3	0.5	0	0	0	0.5
4	0.5	0	0.5	0	0.5
5	1	1	0.5	0.5	0.5
6	1	0	0.5	0	0
7	1	0	0	0	0
8	1	0.5	0	0.5	1
9	1	0.5	1	1	1
10	0.5	0	0	0	1

Clinical Trial Based on TH

The PANGEA-IMBBP Trial

Personalized ANTibodies for Gastro-Esophageal Adenocarcinoma:
A Pilot 1st Metastatic Trial of Biology Beyond Progression

Historical Control (Arm A) → FOLFOX $\xrightarrow{\text{PFS}_1 \text{ 6m}}$ FOLFIRI 60% $\xrightarrow{\text{PFS}_2 \text{ 4m}}$ FOLTAX 30% $\xrightarrow{\text{PFS}_3 \text{ 2m}}$

Biomarker Evaluation in all samples to allow for treatment assignment

i) Primary Endpoints:
Feasibility:
 Time to treatment assignment
Safety: toxicity
ii) mOS (HR 0.67) in HER2+/MET+ Historic Arm A v Arm B
 (Improve to 18 months from 12 months)
 HR: 0.67 of combined HER2+ & MET+
 {12 month survival rate ~63%}
 N= 68 (80% power)

•Secondary Endpoints:
 PFS_{1,2,3}, PFS₁₊₂₊₃ for HER2+/MET+
 mOS for all 5 groups
 2nd/3rd line treatment rates, RR
 -Arm A₁ v A₂/ B₁ v B₂
 -Compared to Historic Controls
 Tissue correlates
 - primary tumor to metastatic lesion
 - baseline v PFS_{1,2,3}

Diagnosis: metastatic cancer

