# Identifying Somatic Mutations via MuTect (from the Broad Institute)

① 

ref genome
```
                              ☆
        ——————————— A ———————————
```

Normal
```
1. ——————— A
2.    ————————— A—
3.      ——— A———
4.           —A———————
5.
             —A———————
6.           -A    ——
```
read depth 6

← appears homozygous at this site.

Tumor
```
1.        — A——
2.         -T——  ———
3.       ——————T
4.       — A  ——
5.       — T—  ———
6.        -A  ————
7.      ———— T—
```
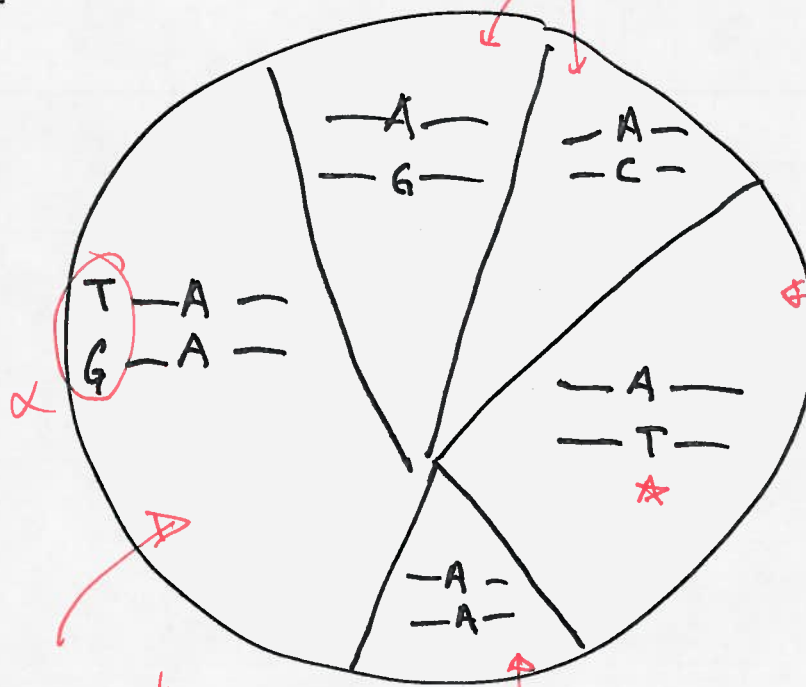read depth 7

← appears heterozygous for A/T.

Somatic mutation.

# Tumors have clonal structure.

Humans have 2 copies of each genomic site.

They have a different mutation. We assume this doesn't occur.

Think of this as a pie chart.

This is the fraction of cells in the tumor sample that have the mutation

These cells are tumor cells but they don't have the $A \rightarrow T$ mutation at site ✶. They might have other somatic mutations though. (e.g. $T-G$ upstream $\alpha$)
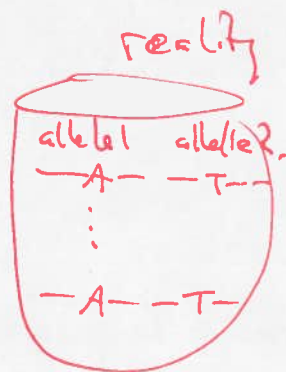
this fraction corresponds to normal cells in the tumor sample. So "contamination". But some normal cells will always be harvested with tumor cells. (What affect does this have on our analysis?)

(pie chart with segments)

— A —
— G —

— A —
— C —

T — A —
G — A —
$\alpha$

— A —
— T —
✶

— A —
— A —

To call somatic mutations in an accurate manner (highly specific & sensitive), there are at least 4 parameters:

1. depth of sequence coverage in tumor & normal.

Imagine depth 6 (only) in tumor. What is the prob. that she is heterozygous at a site where all 6 reads are A?

— A —
— A
— A —
— A —
— A —
— A —

reality

allele1   allele2?
— A —  — T —
⋮
— A —  — T —

6 coin tosses; all 6 are heads.
(the sequencer reaches in the bag and pulls out 1 of 2 copies of the locus at random).

2. Error rate of sequencer.

PHRED measures how reliable a
signal is for a given genomic locus.

Value 0..1.

~~0.0001~~

The idea is that we can measure the
probability the machine makes a mistake

e.g. $\frac{1}{10,000}$ b.p.

Clearly it is difficult to distinguish between
somatic mutations & sequencing errors @ low
depth

Reality is

— A —

— T —

— A —

— G — ← sequencing error
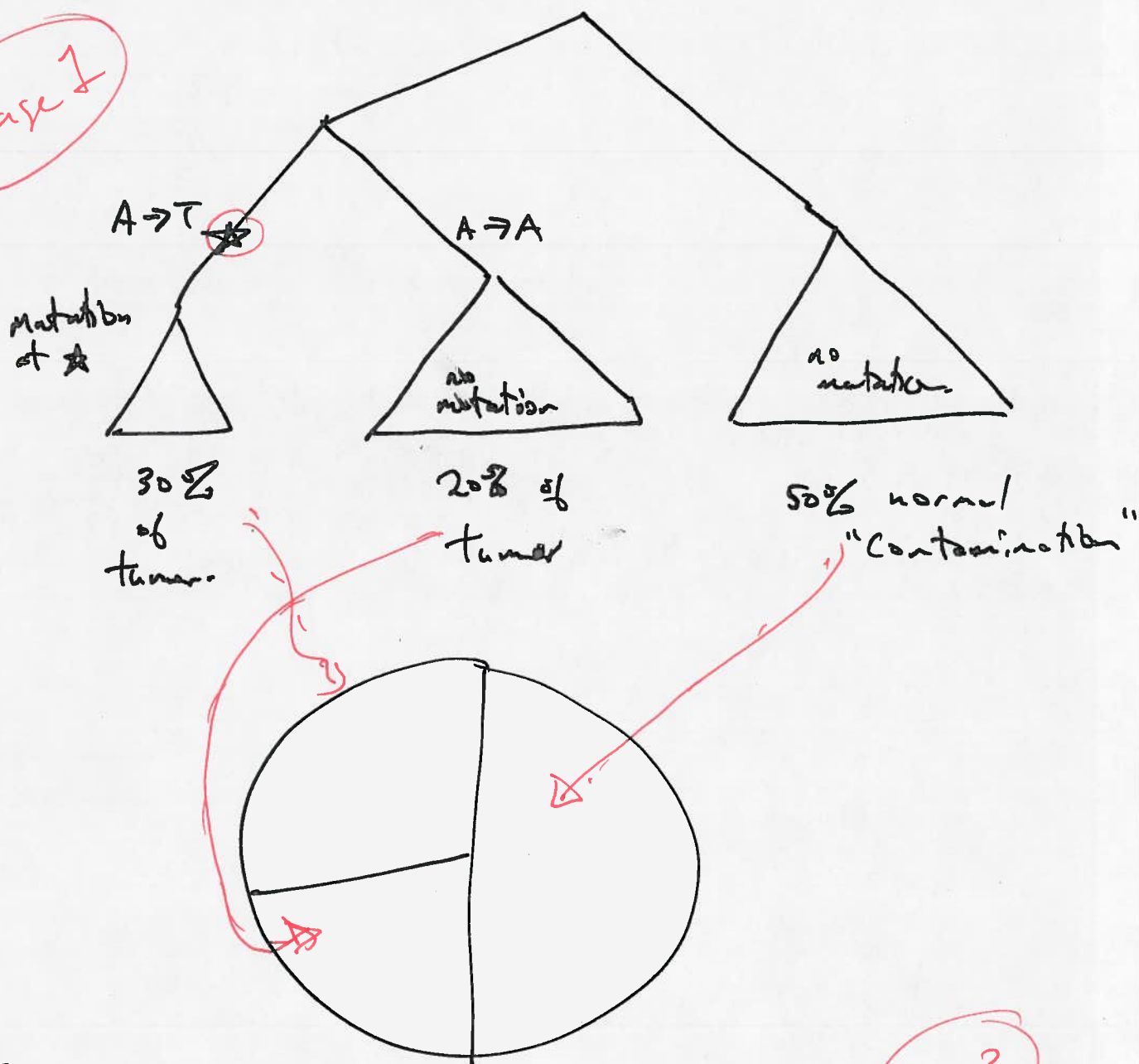
— A —

— T —

A real that covers the alternate
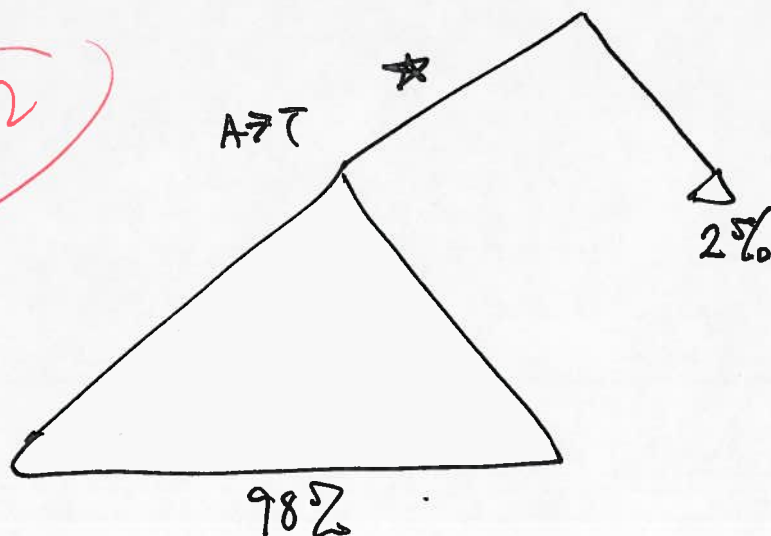allele. m = T.

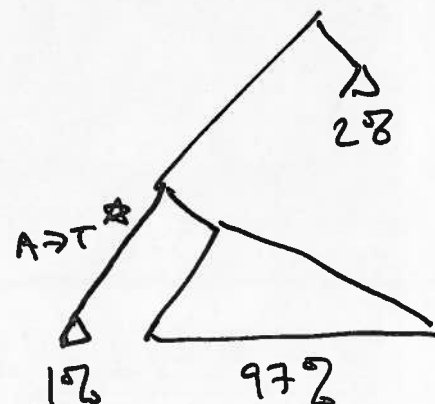# 3. Allele fraction of the mutation

**Case 1**

A→T ☆

Mutation at ☆

30% of tumor.

A→A

no mutation

20% of tumor

no mutation

50% normal "Contamination"

vs.

**Case 2**

A→T ☆

2%

98%

vs. **Case 3**

2%

A→T ☆

1%   97%

4. evidence thresholds.

Normal

$$\begin{array}{c} =\overset{A}{\phantom{}}\ \\ =\overset{A}{\phantom{}}= \\ =A= \\ =A= \end{array}$$

——A——

Tumor

— A ——

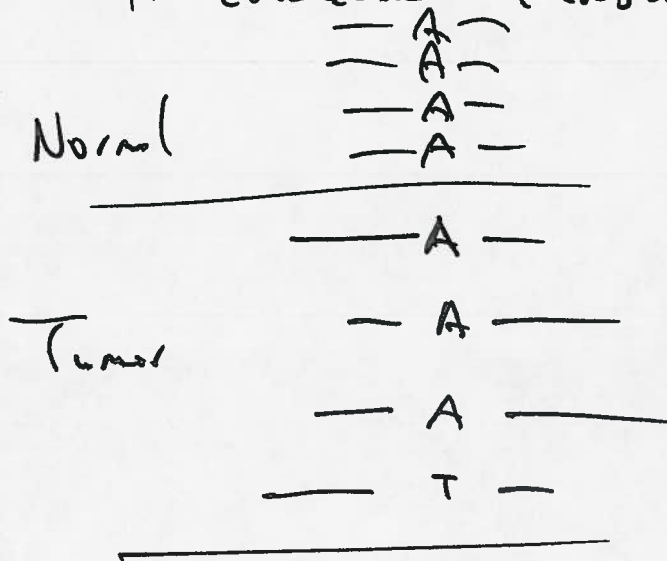— A ——

— T —

With low read depth, difficult to decide if 1 T in 4 reads is enough.

Heterozygous or sequencing mistake?

How certain do you want to be?

Can you tolerate a few mistakes?

Model $M_0$ : no variant at the site.
all non-ref bases are random
sequencing errors.

Alt. model $M_f^m$ : the site contains a true somatic
mutation $m$ with allelic freq $f$.

We could estimate this
by counting the fraction
of reads with the
variant $m$ versus the
total # of reads

—A—
—T—
—A—
—T—
—T—
—A—

50/50
A/T.
$\hat{f} = .50$

log. odds

$$\log_{10}\left(\frac{\mathcal{L}[M_\delta^m] \cdot Pr[m, \delta]}{\mathcal{L}[M_0] \cdot (1 - Pr[m, f])}\right) \geq \log_{10} \delta$$

arbitrary base ... anything > 1 is ok, I guess.

It's easiest to think of $\delta$ as 0.
So if the numerator > denominator
the data supports the alt. model
more than the null model.

In practice $\delta = 6.3$ is used.
Implying the alt. model must be
$10^{6.3} : 1$ favoured over the null model.
More conservative that $\delta = 0 \Rightarrow$ fewer
sites are called somatic mutations.
(Don't worry about where 6.3 comes from)

For our site $i$, let $r \in \{A, C, G, T\}$ represent the reference allele for the woman.

Here we ~~are~~ are going to assume that she is homozygous for either $A, C, G,$ or $T$ at site $i$. (Analysis that allows heterozygosity is a bit more cumbersome but essentially the same.)

Suppose that we have $d$ reads that cover site $i$.

| | $b_i$ | $e_i$ |
|---|---|---|
| 1 | A | 0.1 |
| 2 | T | 0.01 |
| 3 | A | 0.3 |
| 4 | T | $\square$ |
| 5 | A | 0.9 —bad |
| 6 | T | 0.0001 |
| I | i | 0.000001 —good |

— A —
— T —
— A —
— T —
— A —
— T —
— T —

$e_4$ is a probability that the site is a sequencing error $0 \le e_i \le 1$. It is derived by a program called PHRED.

~~That is basically a program to equal~~

Now observe that

$$M_0 = M_f^m \quad \text{where} \quad f = 0$$

(so a variant $m$ exists at the site $i$ but it has frequency $0$ ... just a mathematical reformulation to save us a bit of time).

$$\mathcal{L}\left(M_f^m\right) = Pr\left[\{b_i\} \mid \{e_i\}, r, m, f\right]$$

over all the different reads.

(the prob. ~~of all~~ considering all read calls $b_i$ given what PHRED thinks ($\{e_i\}$), the reference allele of the site, and assuming the site has somatic mutation $m$ with frequency $f$).

$$= \prod_{i=1}^{d} Pr\left[b_i \mid e_i, r, m, f\right]$$

Let's assume all substitutions (mutations) away from the reference $r$ occur with equal probability $e_i/3$.

($e_i/3$ because there are 3 alternatives vs. the observed reference)

$$Pr\left[b_i \middle| e_i, r, m, f\right] = \begin{cases} f \cdot e_i/3 + (1-f)(1-e_i) & \text{if } b_i = r \\ f(1-e_i) + (1-f)(e_i/3) & \text{if } b_i = \\ e_i/3 & \text{otherwise} \end{cases}$$

$e_i/3 \rightarrow$ prob. of sequencing error

$1 - e_i/3 \rightarrow$ prob. of no sequencing error

$f \rightarrow$ frequency of the alt. allele $m$

$1-f \rightarrow$ frequency the site has the ref $r$.

$1-f$

$f$

—A—
—A—

—A—

ref — A —
Norm · — A —
        ⋮
       — A —          } r = A

ref = A

———— A ————        ——————————
——— T ———          —————————
  — A              —————————
——— T ———          ————————
———— T ——          ———
——— A ——           ————————
——— T ——           ——————————
—— G               ——————————
    ↑

$b_8 = G,$      $e_8 = 0.5$

Consider

$Pr[b_8 = G \mid r = A, m = T, f = 50\%]$     (case 3)
        $e_8 = 0.5$

$Pr[b_8 = G \mid r = A, m = C, f = 1\%]$     (case 3)

$Pr[b_8 = G \mid r = A, m = G, f = 20\%]$     (case 3)

Which do you think has the highest probability?

Consider

$$\left[ \begin{array}{l} Pr\left[ b_2 = T \;\middle|\; e_2 = 0.0001, \; \overset{r=A}{m=T}, \; f = 50\% \right] \\[6pt] vs. \\[6pt] Pr\left[ b_2 = T \;\middle|\; e_2 = 0.0001, \; \overset{r=A}{m=T}, \; f = 90\% \right] \\[6pt] vs \\[6pt] Pr\left[ y_2 = T \;\middle|\; e_2 = 0.9, \; \overset{r=A}{m=T}, \; f = 90\% \right] \end{array} \right.$$

Consider

$$\left[ \begin{array}{l} Pr\left[ b_3 = A \;\middle|\; e_2 = 0.5, \; \overset{r=A}{m=T}, \; f = 90\% \right] \\[6pt] vs. \\[6pt] Pr\left[ b_3 = A \;\middle|\; e_2 = 0.5, \; \overset{r=A}{m=T}, \; f = 50\% \right] \end{array} \right.$$

$$\left[ Pr\left[ b_3 = A \;\middle|\; e_2 = 0.000001, \; r = A, \; m = T, \; f = 10^5\% \right] \right.$$

Are we sure the sequencer is correct?