

A Novel Genetic Algorithm for Detecting FLT3 Internal Tandem Duplications in Acute Myeloid Leukemia Patients

Jonathan King

Undergraduate, U.C. Berkeley

BS Bioengineering BA Computer Science

Intern, Plexxikon Inc.



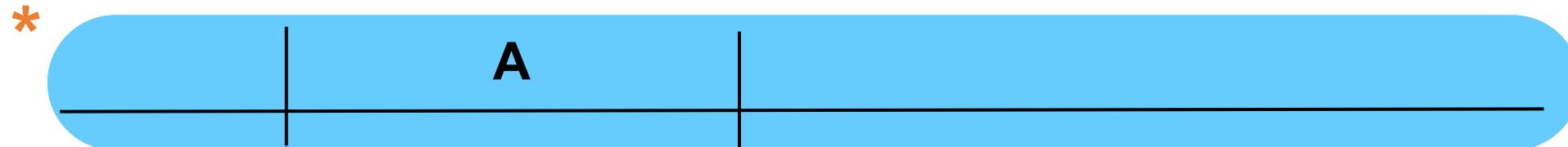
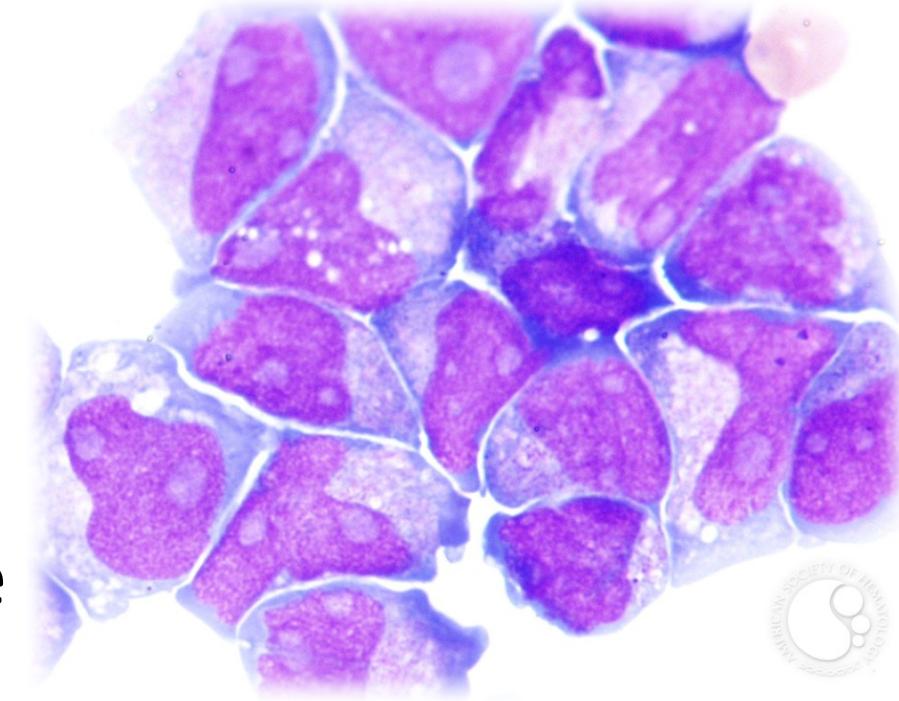
Berkeley
UNIVERSITY OF CALIFORNIA



Plexxikon

Background

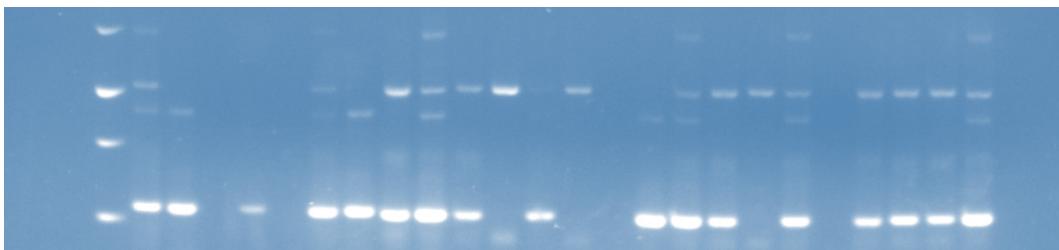
- Treating Acute Myeloid Leukemia
 - 20,000 expected cases in 2016
- Target: “Unfavorable” FLT3-ITD⁺ subtype
 - ITD = Internal Tandem Duplication *
- **Goal:** Identify ITD⁺ individuals for future clinical trials



Identifying ITD⁺ Drug Candidates

PCR

- Lacks potentially useful info:
 - Sequence
 - Precise location, length
- Inaccessible



Next Generation Sequencing

- Best at finding small variants:
 - SNPs
 - Short insertions & deletions
- Limited software for ITDs

```
ACTGTACCAGGACGCGCGCAGGGGGAGGGTACCTCGGTCGACA  
ACACCAGGTGGAGCAGGGCAATCTTACGCTGCCGGCAACTCCCC  
TGTGTATAGTAGTATAGCGGCTAATGATCTACGGTTACGTTATAA  
AGGGAGCAGCGTTCGCTTATTGGAAGTAGTATAGGTTGGCGATAT  
CTCCCGATAGGAGCAGTATCGAGGGTCACTGAACGTTCAGGGGC
```

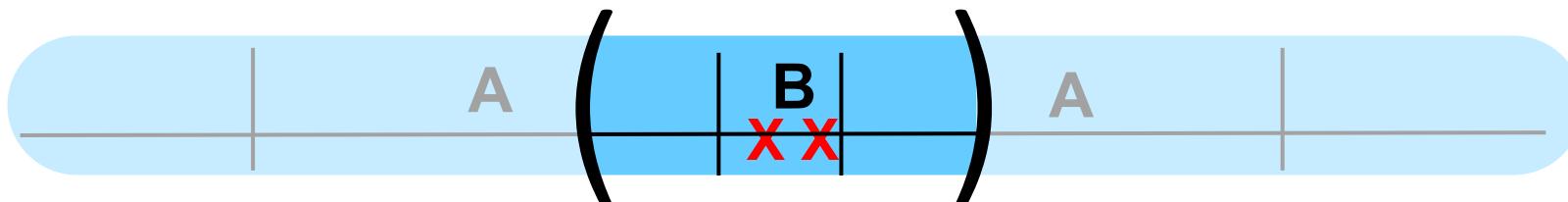
Challenges to ITD Discovery

- Vary in location, length, and frequency!

Read

ITD

Claim: The ~ 20bp “Junction of Duplication” uniquely identifies FLT3-ITDs.



Probability of same sequence occurring by random chance in the region of interest is:

$$\frac{1}{4^{20}} \times 10^3 \approx \frac{1}{10^9}$$

An Algorithm for Detecting ITDs in FLT3

Normal DNA Sequence (Reference)



ITD Positive Sample



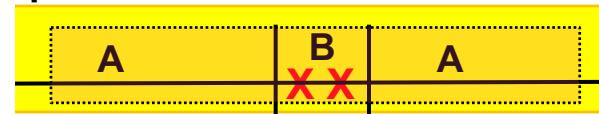
Read 1



Read 2

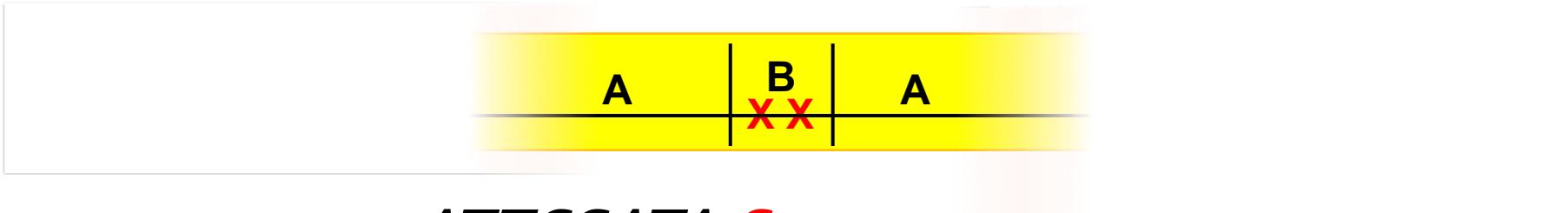


Search Sequence



~ 20 bp

Flexible Searching with Regular Expressions



ATTCGATA (C [ACTG]{0,m} A) CGAAGTTC

Left Probe [0:n - 1] *LP [n]* *between 0 and m bases* *RP [0]* *Right Probe [1:n]*

Example parameters:

$$n = 13$$

$$m = 30$$



FLT3 ITD Identification Results

Plexxikon

07/19/2016 11:45:39 AM

9996_S7 is **ITD positive**, with **168** supporting reads and **2** possible duplications.

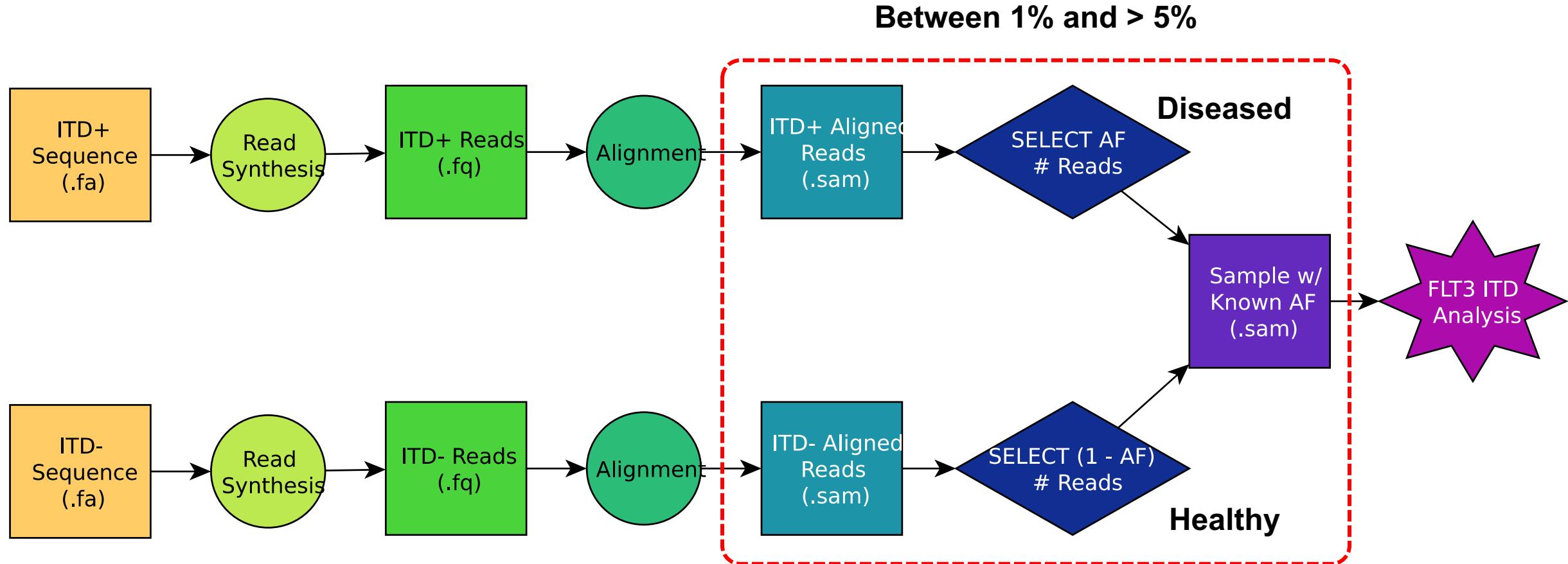
Duplication 1

Location	Length	Reads	Non-ITD Reads	Allele Frequency	Reads with entire ITD
28608262	51	118	769	0.13303	80

Sequence:

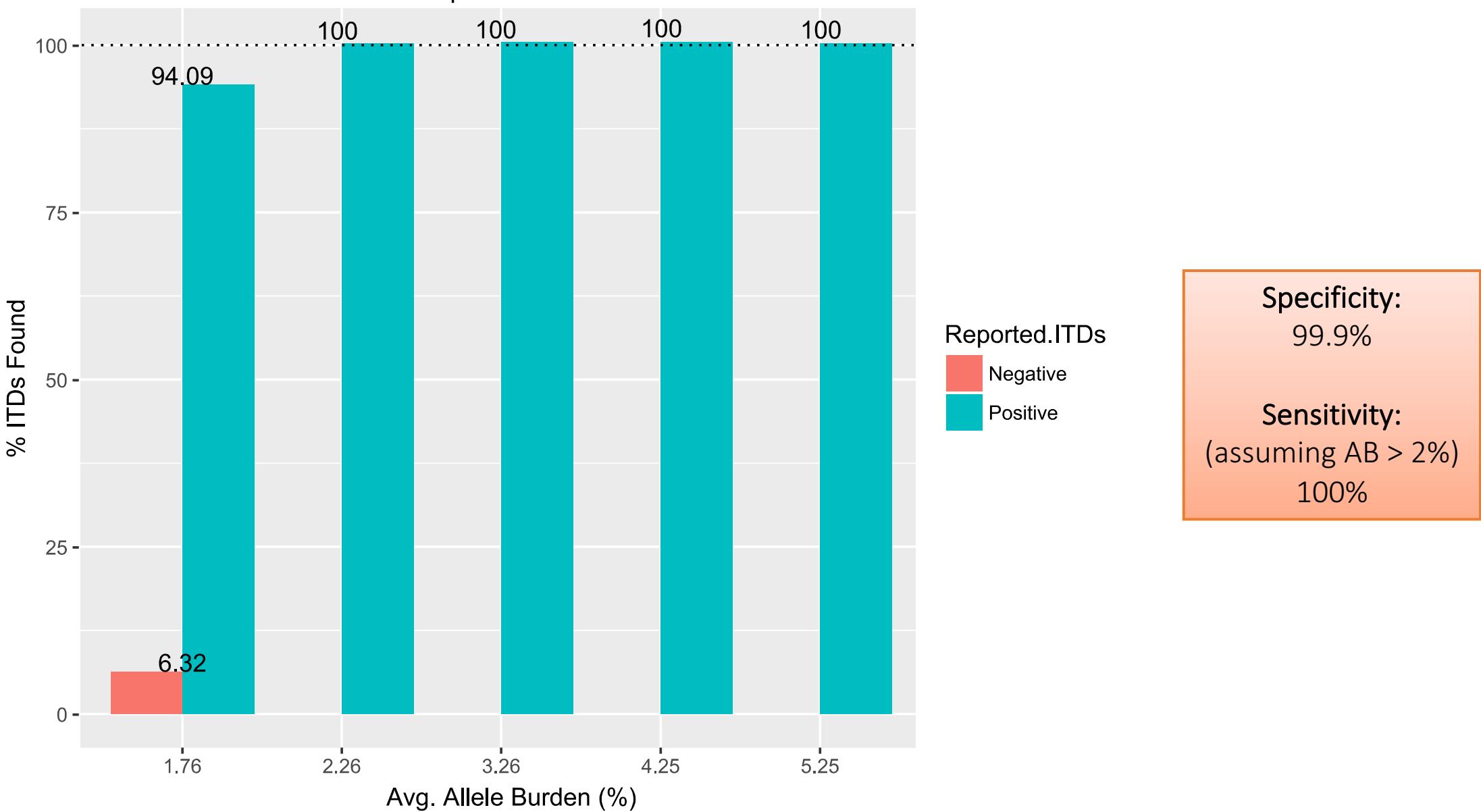
ATATGATCTCAAATGGGAGTTCCAAGAGAAAATTAGAGTTGGattccc

in silico Testing Procedure

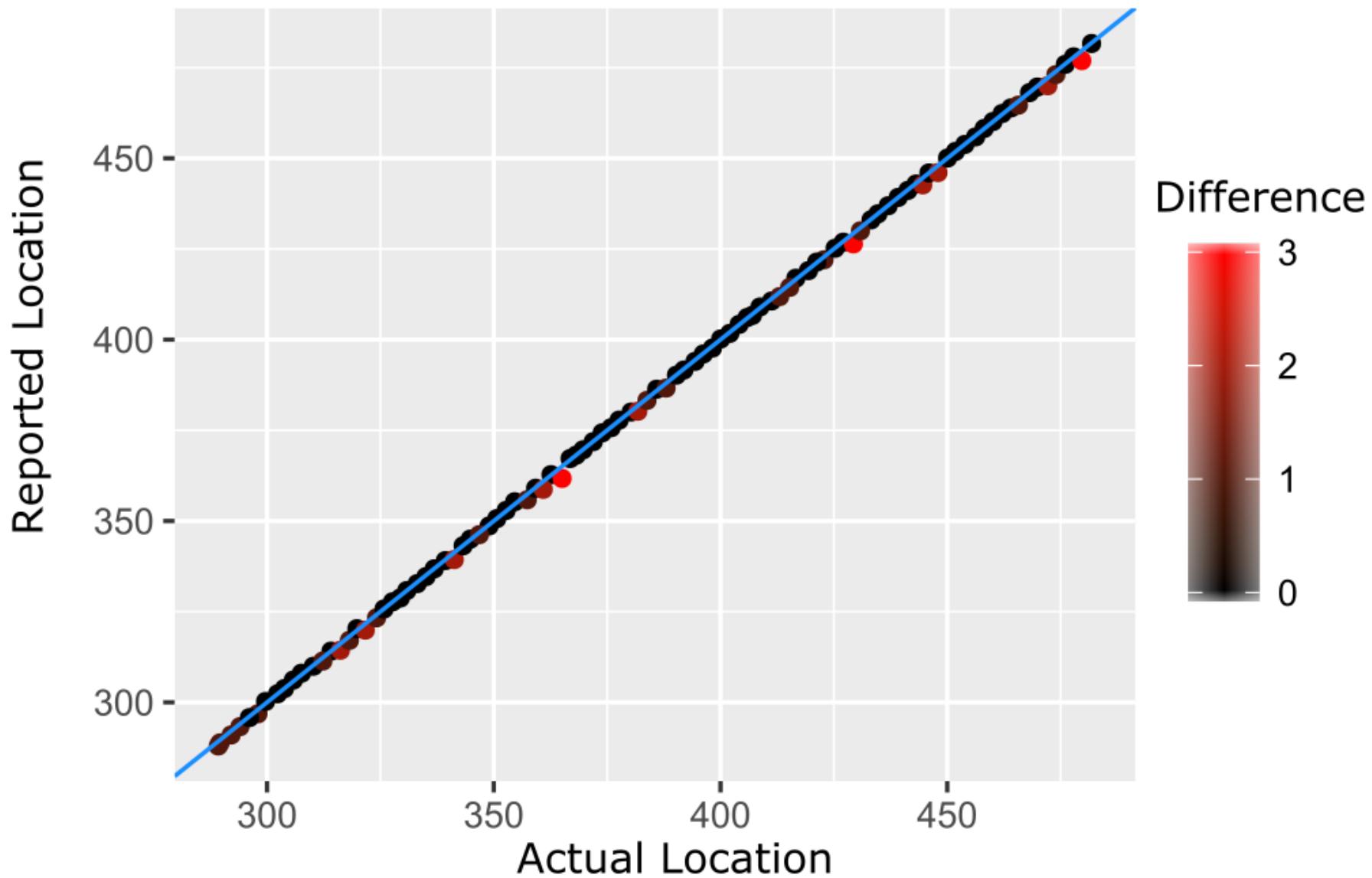


Detection of *in silico* ITDs

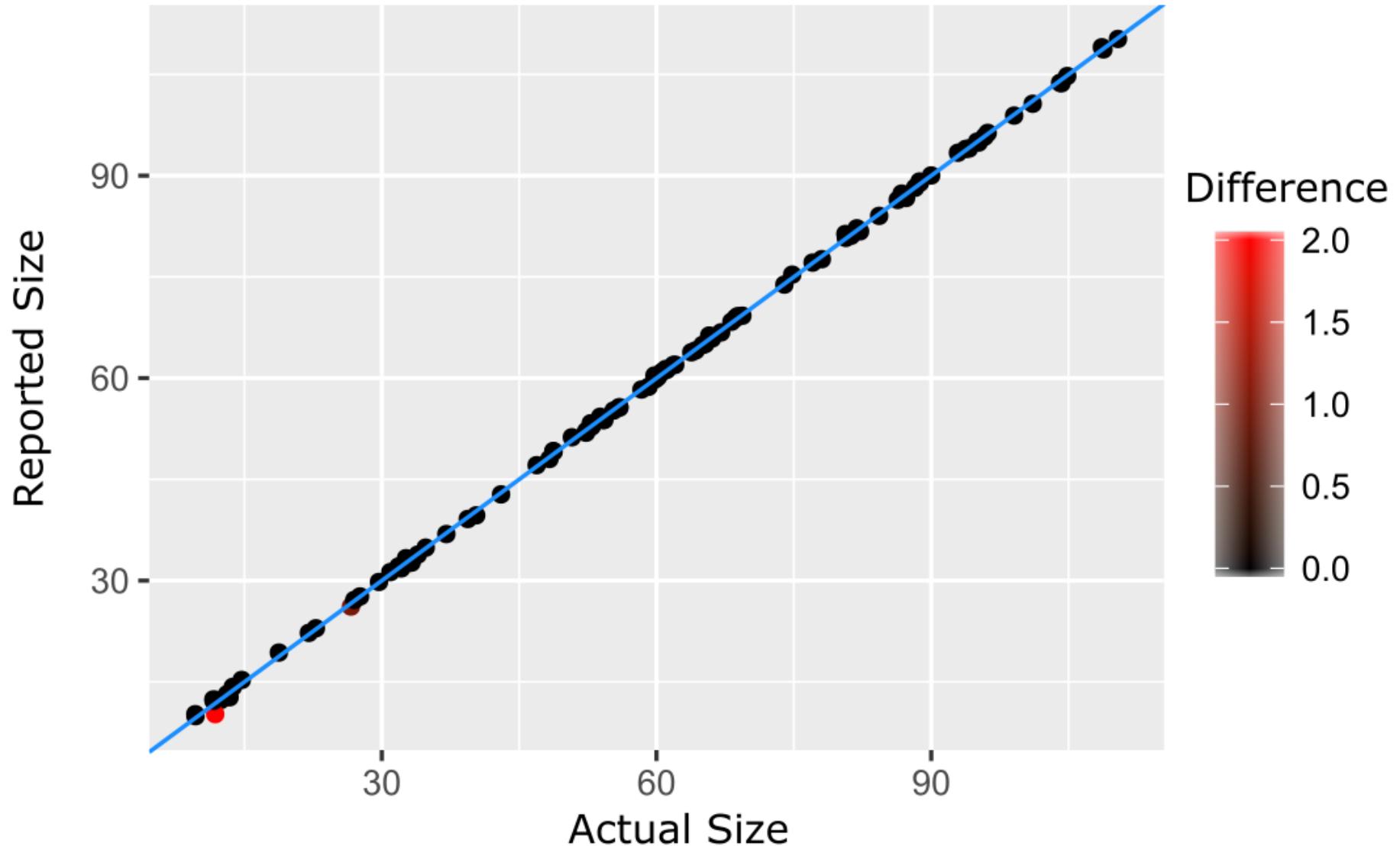
300 bp simulated reads



Location Reporting Accuracy



Size Reporting Accuracy



Using the Algorithm Offsite

Microbial Genomes RNA Methylation ChIP 16S rRNA DNA
Health Research miRNA Cancer Genetic Disease Discovery
Disease Discovery Metagenomics Microbial Genomes ChIP 16S rRNA DNA Development Cell Public Health Research

BaseSpace®

Basespace SEQUENCE HUB DASHBOARD PREP RUNS PROJECTS APPS PUBLIC DATA

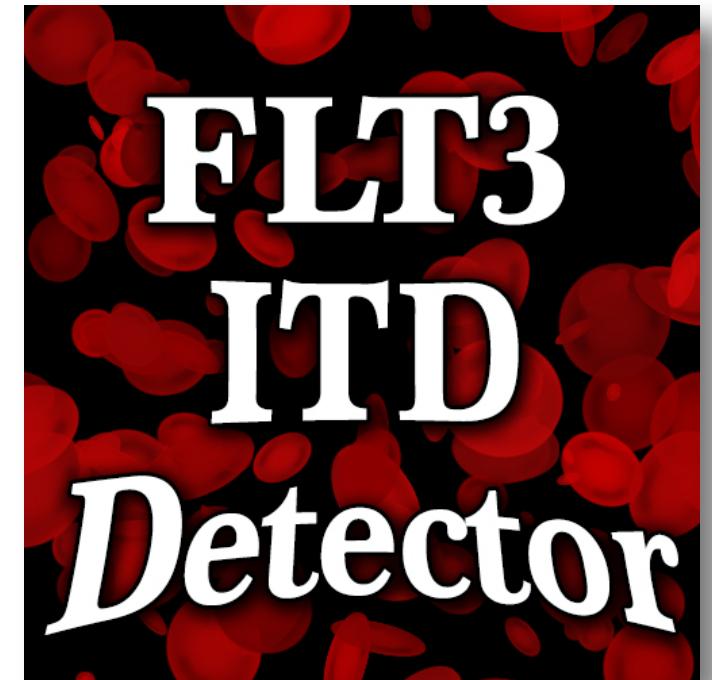
Applications

 BWA Enrichment Illumina, Inc.	 Cancer Variant Caller Samsung SDS
 DeepChek®-HBV ABL S.A.	 DeepChek®-HCV ABL S.A.
	

Search Apps

Categories

ChIP-Seq (5)	De Novo Assembly (7)
Differential Expression (14)	Gene Fusion Detection (7)
HIPAA (7)	Metagenomics (12)
Methyl-Seq (6)	Proteomics (10)
Quality (9)	Resequencing (18)



Thank you!



Plexxikon



California Institute for
Quantitative Biosciences

Paul Severson

Brian West

Gideon Bollag, Chao Zhang, Paul Lin

Marguerite Hutchinson, Robert Bernstein

Interns!

Prof. Susan Marqusee
Donna Hendrix

Friends and family!



Berkeley EECS
ELECTRICAL ENGINEERING & COMPUTER SCIENCES

BIOENGINEERING
UNIVERSITY of CALIFORNIA, BERKELEY