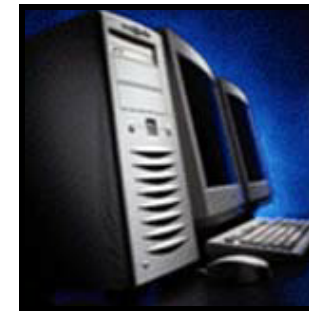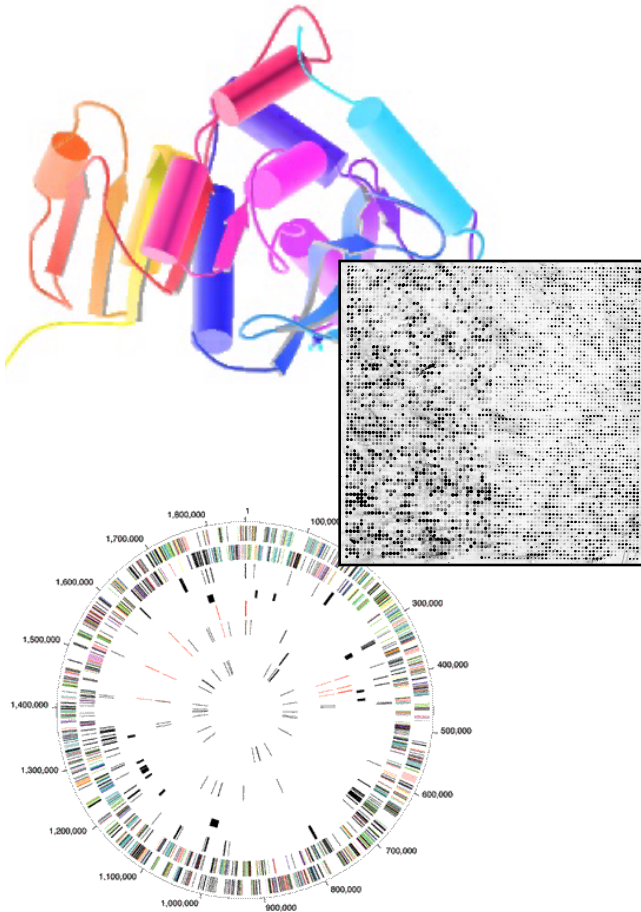# Variant Identification, Focusing on SVs
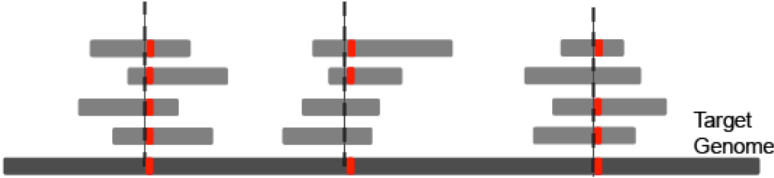


Mark Gerstein, Yale University

gersteinlab.org/courses/452

(last edit in spring '20, pack #6)

# Main Steps in Genome Resequencing

**[Snyder et al. Genes & Dev. ('10)]**
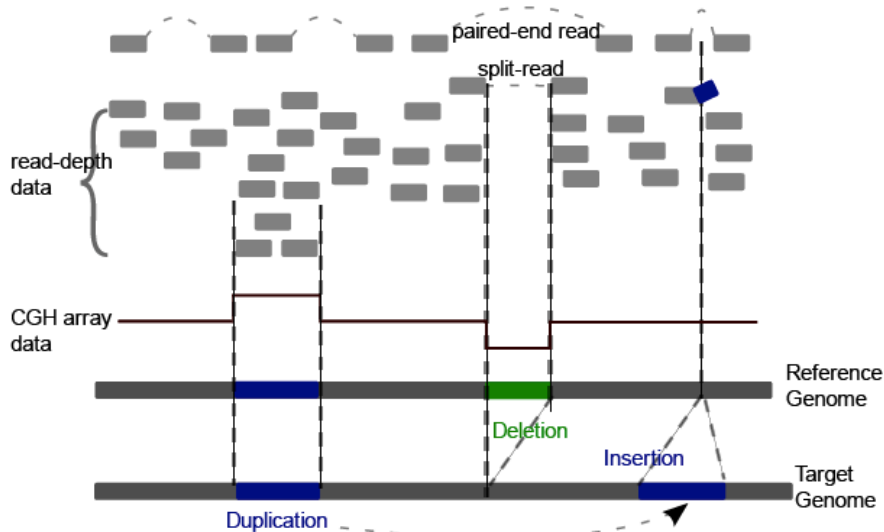
**Step 0: Generate Reads**

**Step 1: Call SNPs**

using uniquely and correctly mapped reads

Target Genome

**Step 2: Find SVs**

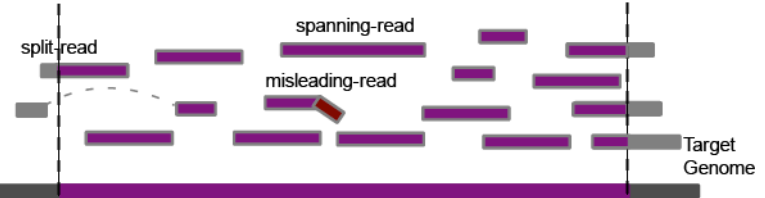with aberrant paired-end reads, split-reads, read-depth analysis and CGH array data

paired-end read

split-read

read-depth data

CGH array data

Reference Genome

Deletion

Insertion

Target Genome

Duplication

**Step 3: Assemble New Sequences**

with split-, spanning- and misleading-reads

split-read

spanning-read

misleading-read

Target Genome

**Step 4: Phasing**

mostly with paired-end reads

SNP / Indel

paired-end read

Insertion (heterozygous)

Inversion (heterozygous)

Target Diploid Genome

Duplication

# Main Steps in Genome Resequencing

**[Snyder et al. Genes & Dev. ('10)]**



**Step 0: Generate Reads**

**Step 1: Call SNPs**
using uniquely and correctly mapped reads

Target Genome

**Step 2: Find SVs**
with aberrant paired-end reads, split-reads, read-depth analysis and CGH array data
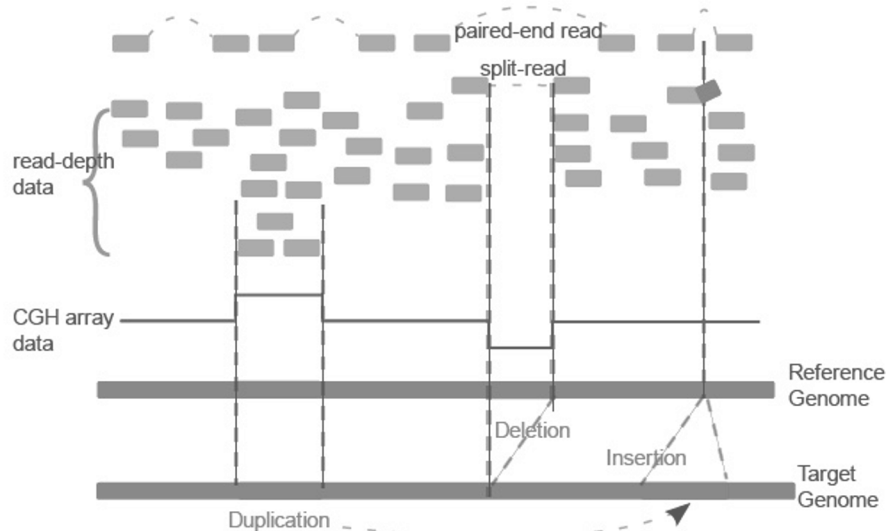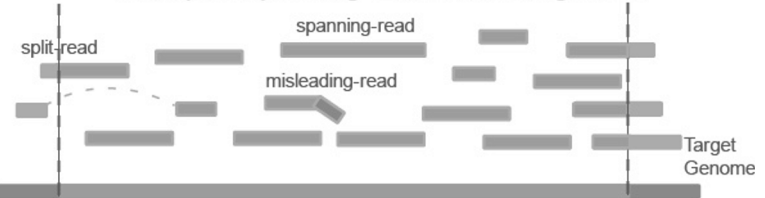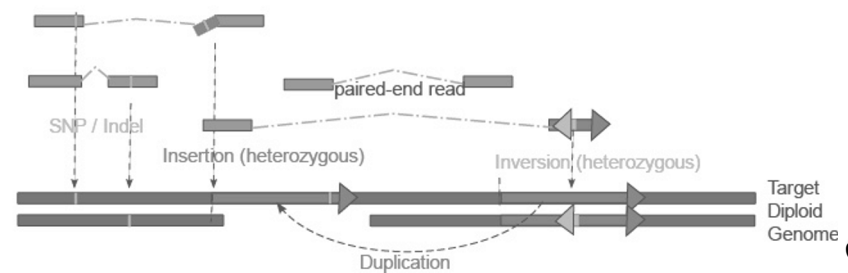
paired-end read
split-read
read-depth data
CGH array data
Reference Genome
Deletion
Insertion
Target Genome
Duplication

**Step 3: Assemble New Sequences**
with split-, spanning- and misleading-reads

split-read
spanning-read
misleading-read
Target Genome

**Step 4: Phasing**
mostly with paired-end reads

SNP / Indel
paired-end read
Insertion (heterozygous)
Inversion (heterozygous)
Target Diploid Genome
Duplication

# Characterization of genomic variations: somatic vs germline



**Sequencing tumor and normal samples from cancer patients provide insight into somatic and germline variation profile.**

4

*Samuel et al, Clinical chemistry, 2013*

# Bayes' Theorem to detect genomic variant

A AGCTTGAC TCCA TGATGATT
B AGCTTGAC GCCA TGATGATT
C AGCTTGAC TCCC TGATGATT
D AGCTTGAC GCCC TGATGATT
E AGCTTGAC TCCA TGATGATT
F AGCTTGAC GCCA TGATGATT
G AGCTTGAC TCCC TGATGATT
H AGCTTGAC GCCC TGATGATT

$$P(G|D) = \frac{P(D|G)P(G)}{P(D)}$$

$$= \frac{P(D|G)\,P(G)}{\sum_{i=1}^{n} P(D|G_i)\,P(G_i)}$$

In the above equation:

- $D$ refers to the observed data
- $G$ is the genotype whose probability is being calculated
- $G_i$ refers to the $i$th possible genotype, out of n possibilities

Calculating the conditional distribution $P(D|G)$:

Assuming an error free model, for each heterozygous SNP site of the diploid genome, covered by K reads, the number of reads $i$ representing one of the two alleles follows binomial distribution.

$$P_{err\_free}(D|G) = f(i|k, 0.5) = \binom{k}{i} 0.5^k$$

With errors, the calculation is more complicated.

In general:

$$P(D|G) = P_{err\_free}(D|G) + P_{err}(D|G)$$
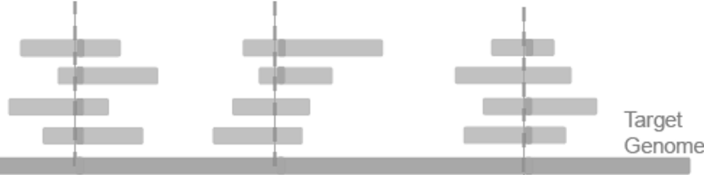
**Main Steps in Genome Resequencing**

[Snyder et al. Genes & Dev. ('10)]

**Step 0: Generate Reads**

**Step 1: Call SNPs**
using uniquely and correctly mapped reads

Target Genome

**Step 2: Find SVs**
with aberrant paired-end reads, split-reads, read-depth analysis and CGH array data

paired-end read
split-read
read-depth data
CGH array data
Reference Genome
Deletion
Insertion
Target Genome
Duplication

**Step 3: Assemble New Sequences**
with split-, spanning- and misleading-reads

split-read
spanning-read
misleading-read
Target Genome

**Step 4: Phasing**
mostly with paired-end reads

SNP / Indel
paired-end read
Insertion (heterozygous)
Inversion (heterozygous)
Target Diploid Genome
Duplication

# Methods to **Find SVs**

## 1. Paired ends

Deletion

Reference

Genome

Sequenced    paired-ends

**Mapping** →

Reference

## 2. Split read

Deletion

Reference

Genome

Read

**Mapping**

Reference

## 3. Read depth (or aCGH)

Deletion

Reference

Genome

Reads

**Mapping**

Read count

Zero level

## 4. Local Reassembly

# Read Depth

**Array Signal**

**Read depth**

Patient 98-135

**Individual genome**

**Reads**

Mapping

**Reference genome**

Counting mapped reads

**Read depth signal**

**Zero level**

[Urban et al. ('06) PNAS; Wang et al. Gen. Res. ('09); Abyzov et al. Gen. Res. ('11)]

*10* – Lectures.GersteinLab.org

# Example of Application to RD data



**NA12878, Solexa 36 bp paired reads, ~30x coverage**

# HMM

- To get highest resolution on breakpoints need to smooth & segment the signal
- BreakPtr: prediction of breakpoints, dosage and cross-hybridization using a system based on Hidden Markov Models



Korbel*, Urban* *et al.,* PNAS (2007)

# Statistically integrates array signal and DNA sequence signatures
## (using a discrete-valued bivariate HMM)



Korbel*, Urban* *et al.,* PNAS (2007)

# Mean-shift-based (MSB) segmentation: no explicit model

- For each bin attraction (mean-shift) vector points in the direction of bins with most similar RD signal

- No prior assumptions about number, sizes, haplotype, frequency and density of CNV regions

- Not Model-based (e.g. like HMM) with global optimization, distr. assumption & parms. (e.g. num. of segments).

- Achieves discontinuity-preserving smoothing

- Derived from image-processing applications



**CNVnator**

RD signal

Bins

[Abyzov et al. Gen. Res. ('11)]

# Intuitive Description of MSB

**Observed depth of coverage counts as samples from PDF**

**Kernel-based approach to estimate local gradient of PDF**

**Iteratively follow grad to determine local modes**

**Region of interest**

**Center of mass**

**Mean Shift vector**

<u>Objective</u> : **Find the densest region**

**Distribution of identical billiard balls**

# Paired-End

# Paired-End Mapping



- Both paired-ends map within repeats.
- Limited the distance between pairs; therefore, neither large nor very small rearrangements can be detected

# Split Read

# Split-read Analysis

Deletions are the Easiest to Identify
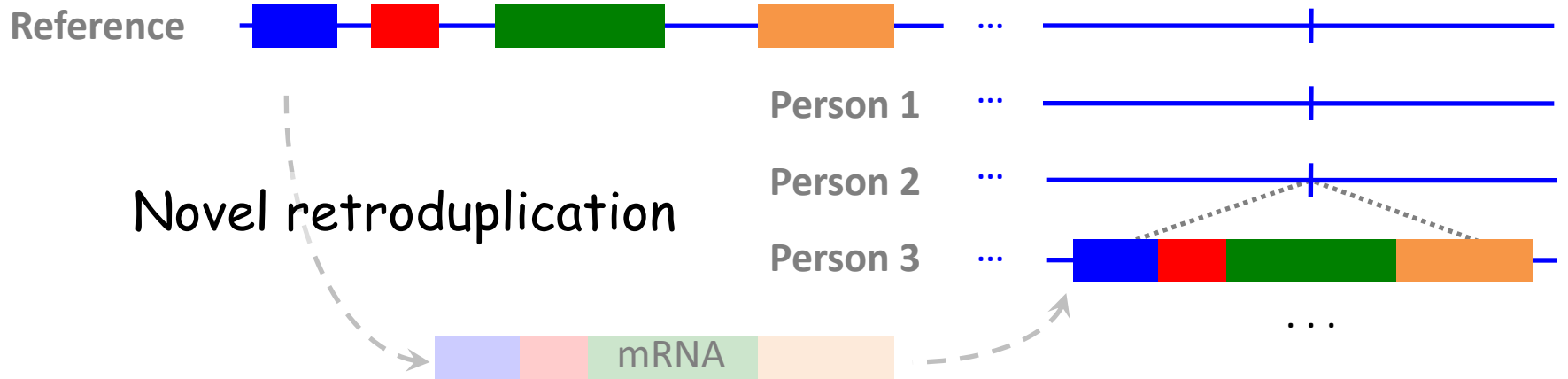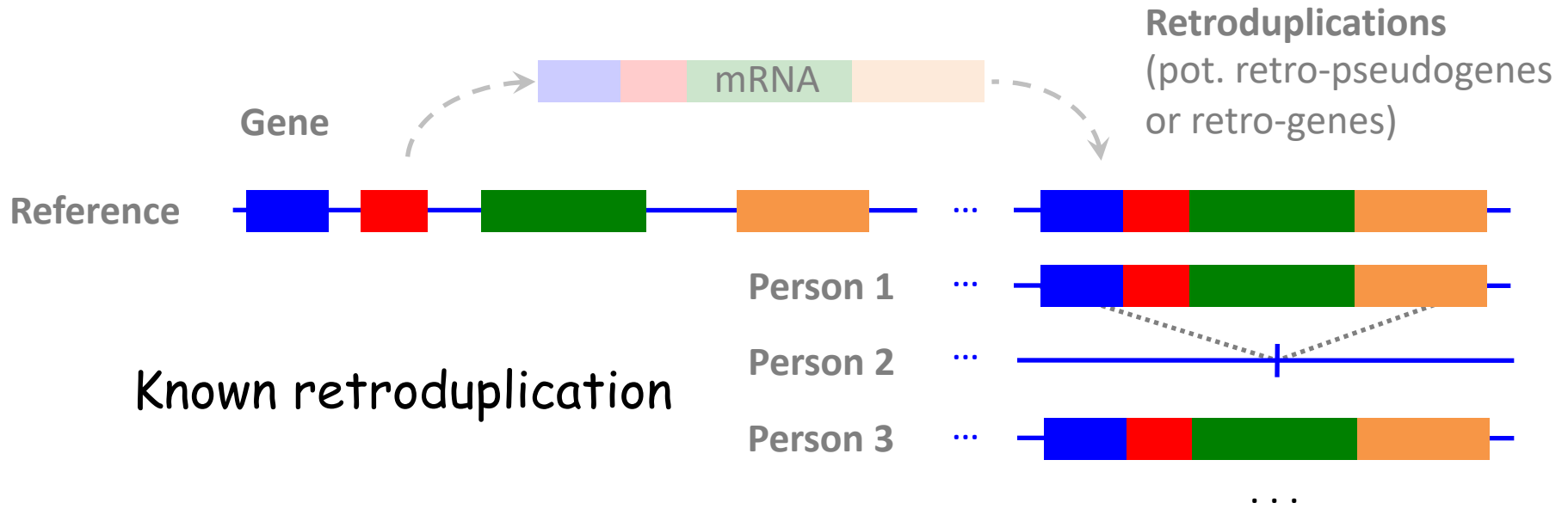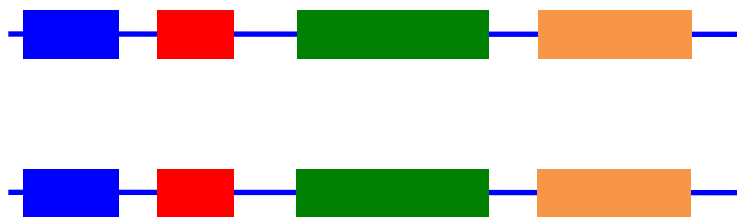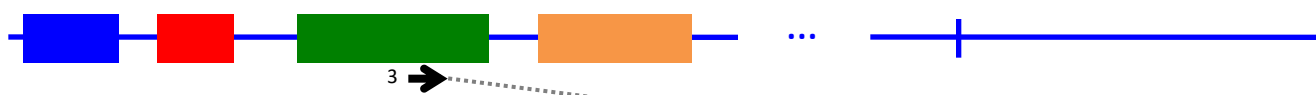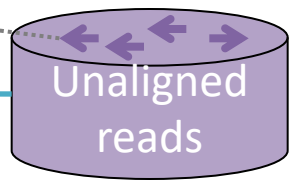
# RDV & Mobile Elements

# Retroduplication variation (RDV)

[Abyzov et al. Gen. Res. ('13) ]

**Gene**

**Novel retroduplication**

Read pairs

Reference

Alignment to the reference

**1**

Aligned reads

Evidence from alignment

Unaligned reads

Splice-junction library

**2**

Evidence from cluster

**3**

Evidence from read depth

Zero level

[Abyzov et al. Gen. Res. ('13) ]

**Pipeline to identify novel retro-dups. from 3 evidence sources**