

Analysis of Next Generation Sequence Data

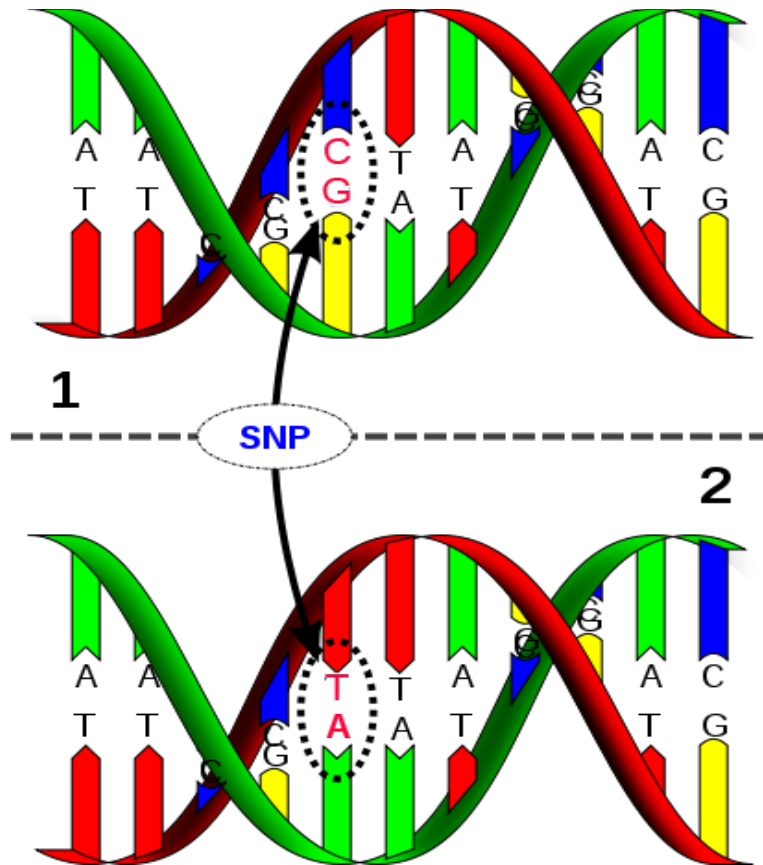
BIOST 2055

03/16/2012

Outline

- Background
- From sequence data to genotype
 - Alignment
 - SNP detection and genotype calling
 - Genotype refinement
- A walk-through example
 - Get familiar with various file formats
 - Get familiar with popular programs

Human Genome and Single Nucleotide Polymorphisms (SNPs)



- 23 chromosome pairs
- 3 billion bases
- A single nucleotide change between pairs of chromosomes
- E.g.

Haplotype1: AAGG**G**ATCCAC

Haplotype2: AAGG**A**ATCCAC

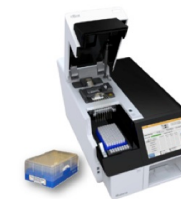
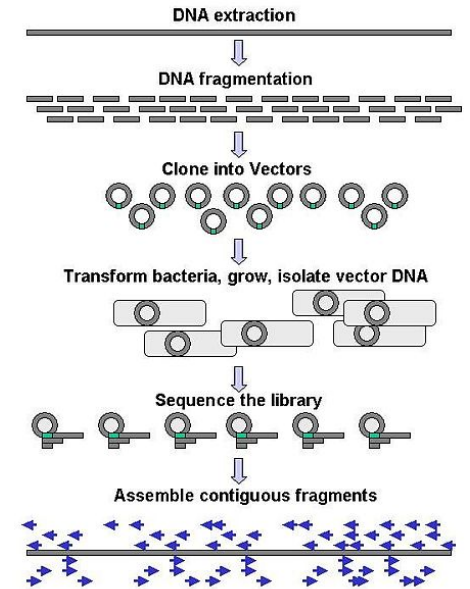
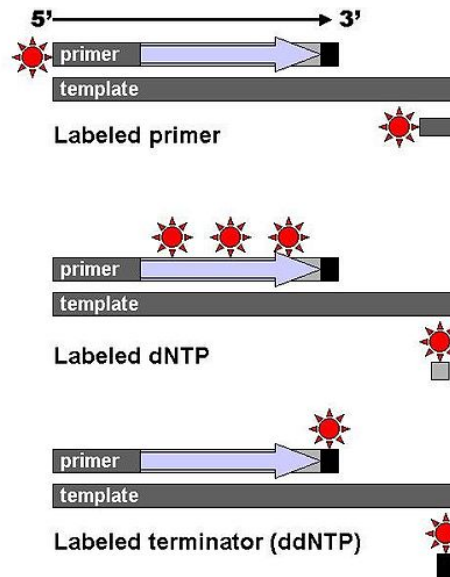
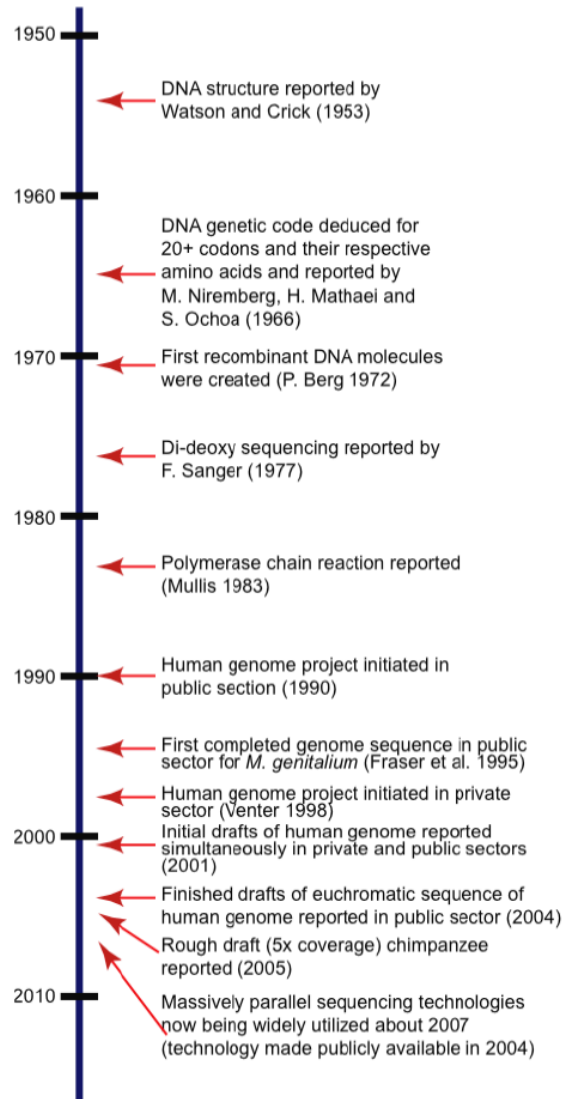
SNPs in Population



Association Study in Case Control Samples

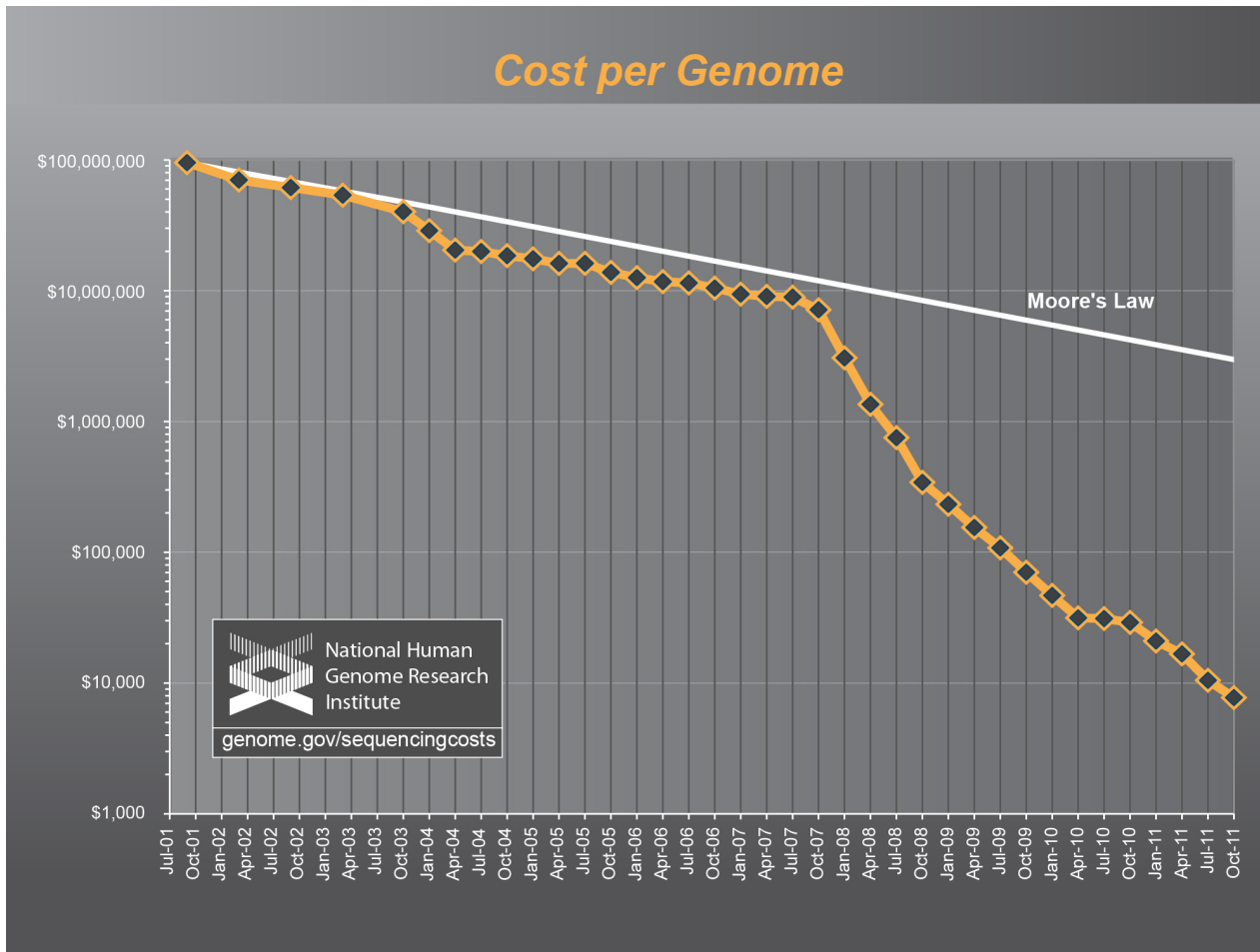


History of DNA Sequencing



- 2 flow cells (can be run independently)
- Up to 320Gb mapped sequence per FC
- 64Gb sequence per day (2 flow cells)

Sequencing Cost



A Road to Discover Human Genome



1990-2003



2002 -



2008 -

Different Approaches

- Deep whole genome sequencing
 - Expensive, only can be applied to limited samples currently
 - Most complete ascertainment of all variations
- Low coverage whole genome sequencing
 - Modest cost, typically 100-1000 samples
 - Complete ascertainment of common variations
 - Less complete ascertainment of rare variants
- Exome capture and targeted region sequencing
 - Modest cost, high coverage
 - Most interesting part of the genome

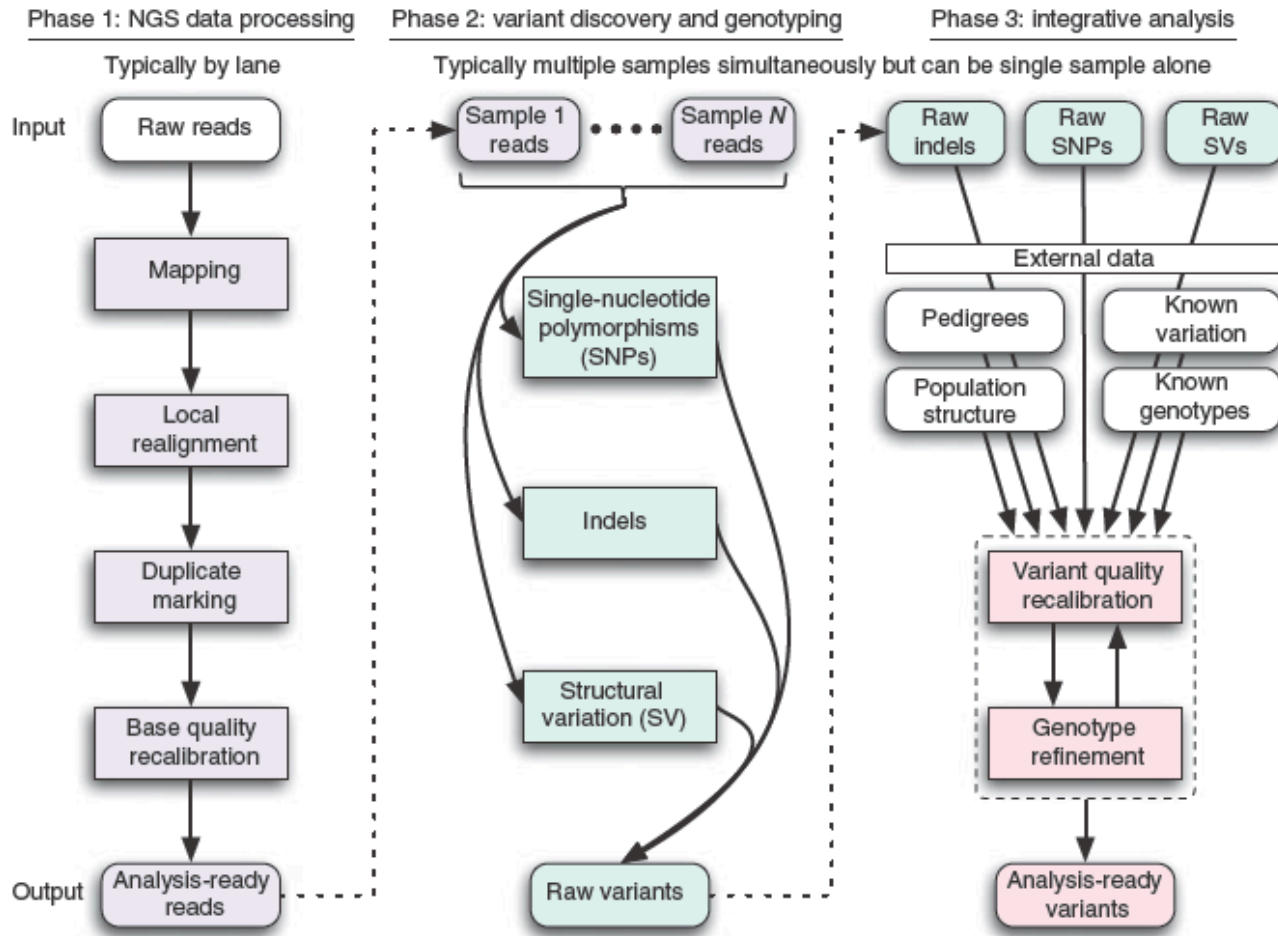
With Complete Sequence Data

- What is the contribution of each identified locus to a trait?
 - Multiple variants, common and rare
 - Effect size
- What is the mechanism? What happens when we knockout a gene?
 - Most often, the causal variant will not have been examined directly
 - Rare coding variants will provide important insights into mechanisms
- What is the contribution of structural variation to disease?
 - These are hard to interrogate using current genotyping arrays
- Are there additional susceptibility loci to be found?
 - Only subset of functional elements include common variants
 - Rare variants are more numerous and thus will point to additional loci

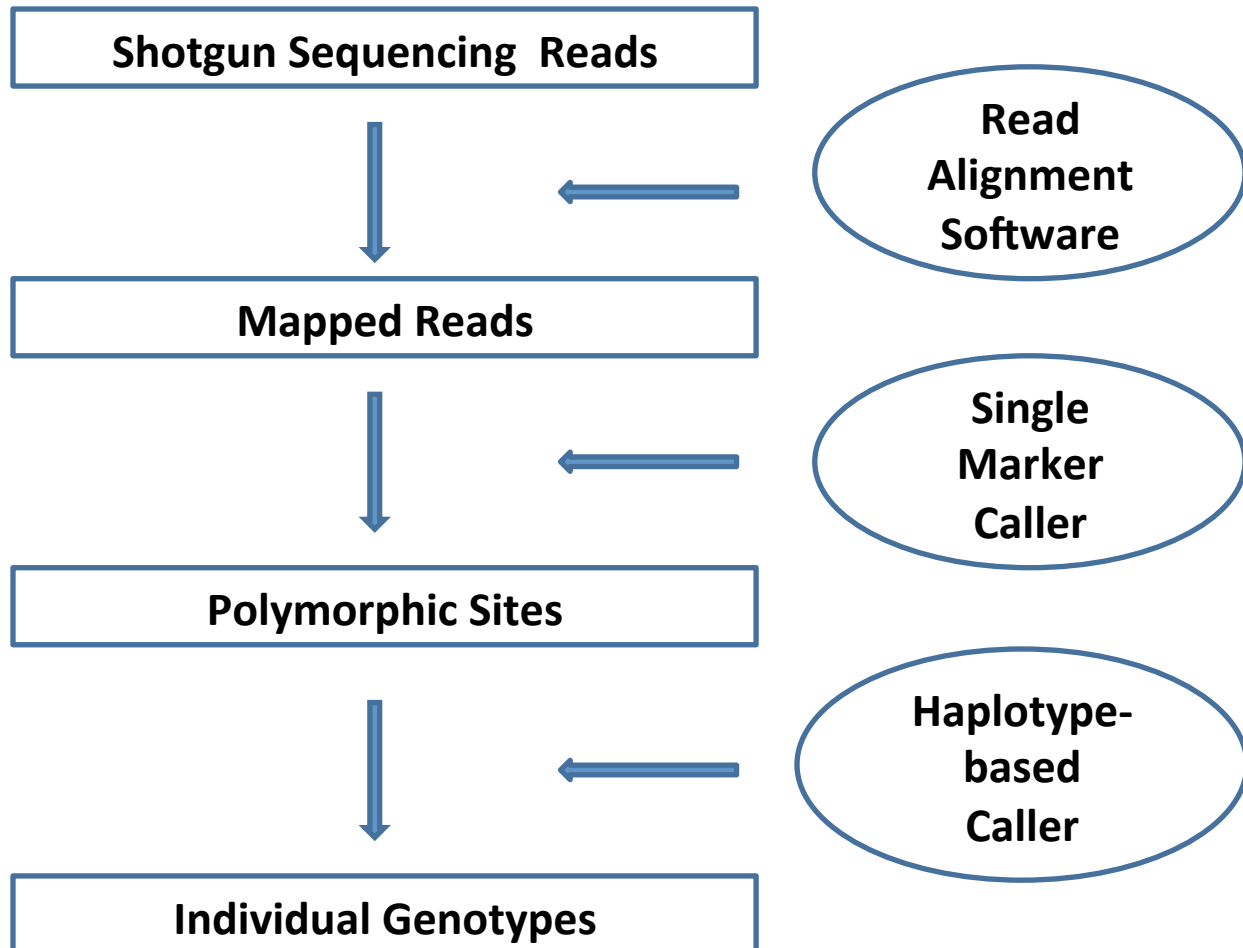
Next Generation Sequencing

- Commercial platforms produce gigabases of sequence rapidly and inexpensively
 - ABI SOLiD, Illumina Solexa, Roche 454, Complete Genomics, and others...
- Sequence data consist of thousands or millions of short sequence reads with moderate accuracy 0.5 – 1.0% error rates per base may be typical
- High-throughput but hard to assemble

A Typical Pipeline



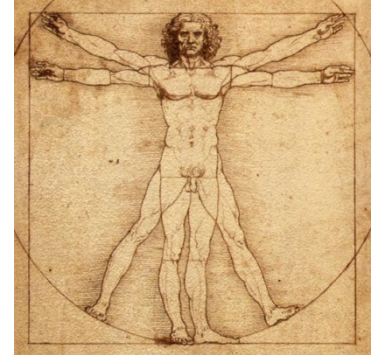
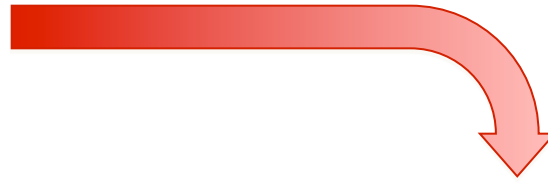
A Typical Pipeline



Short read alignment



Sequencer



Human source

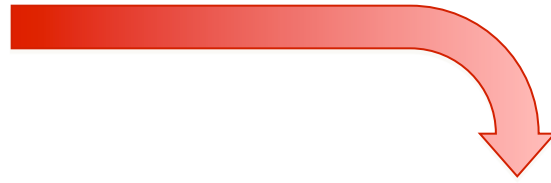


Reads from new sequencing machines are short: 30-400 bp

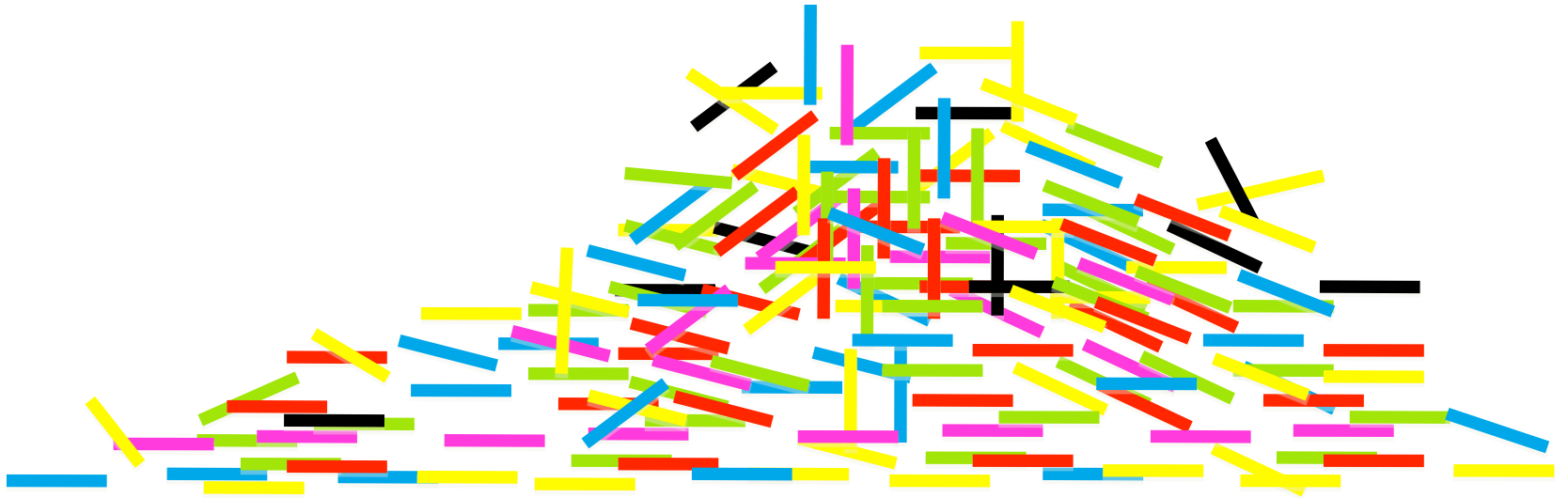
Short read alignment



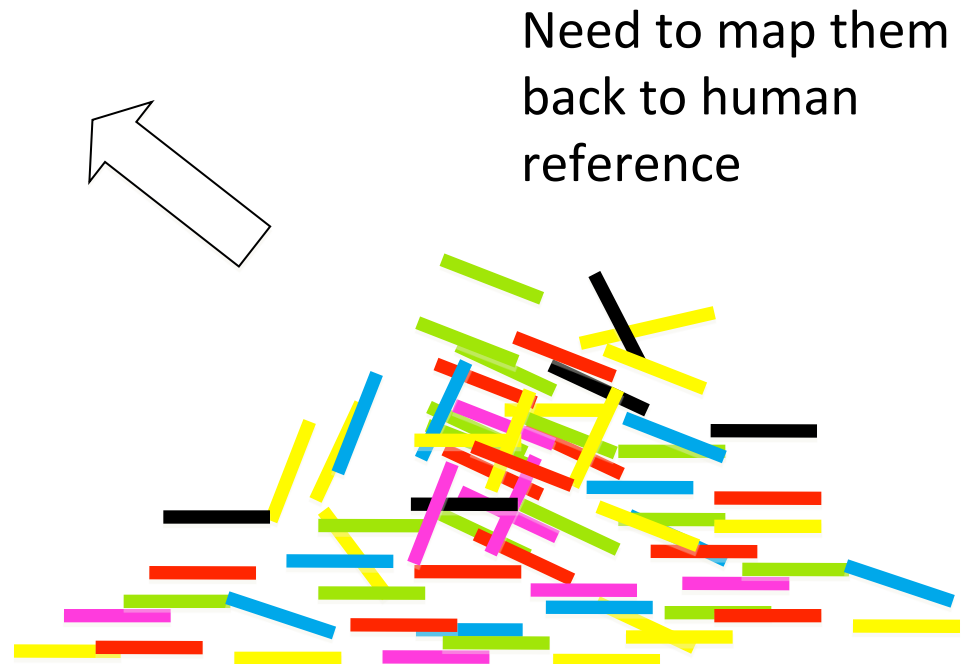
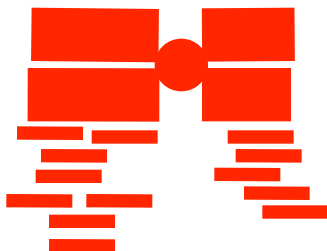
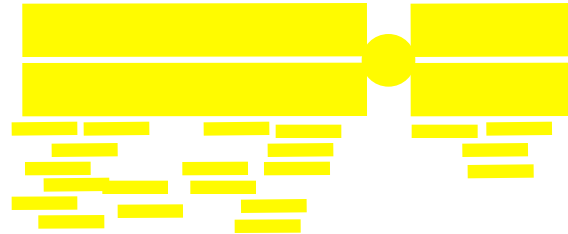
Sequencing machine



And you get
MILLIONS of them



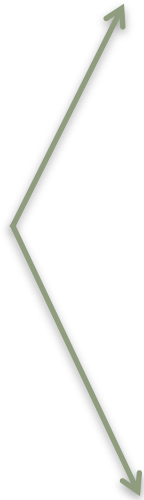
Short read alignment



Alignment

Reference sequence:

actgtagattag**ccgagtagctag**ctagtcgat



Find best match for each
read in a reference
sequence

ccgagaagctag

Which is better? Long words or Short words?

Existing Alignment by Category

- Hashing reference genome
 - SOAP1, MOSAIK, PASS, BFAST, ...
- Hashing short reads
 - Eland, MAQ, SHRiMP, ...
- Merge-sorting reference together with reads
 - Slider
- Based on Burrows-Wheeler Transform
 - BWA, SOAP2, Bowtie, ...

A list

Sean Penn	Michael Douglas	Paul Scofield	Ray Milland
Daniel Day Lewis	Paul Newman	Lee Marvin	Bing Crosby
Forest Whitaker	William Hurt	Rex Harrison	Paul Lukas
Philip S. Hoffman	F. Murray Abraham	Sidney Poitier	James Cagney
Jamie Foxx	Robert Duvall	Gregory Peck	Gary Cooper
Sean Penn	Ben Kingsley	Maximilian Schell	James Stewart
Adrien Brody	Henry Fonda	Burt Lancaster	Robert Donat
Denzel Washington	Robert De Niro	Charlton Heston	Spencer Tracy
Russell Crowe	Dustin Hoffman	David Niven	Spencer Tracy
Kevin Spacey	Jon Voight	Alec Guinness	Paul Muni
Roberto Benigni	Richard Dreyfuss	Yul Brynner	Victor McLaglen
Jack Nicholson	Peter Finch	Ernest Borgnine	Clark Gable
Geoffrey Rush	Jack Nicholson	Marlon Brando	Charles Laughton
Nicolas Cage	Art Carney	William Holden	Wallace Beery
Tom Hanks	Jack Lemmon	Gary Cooper	Fredric March
Tom Hanks	Marlon Brando	Humphrey Bogart	Lionel Barrymore
Al Pacino	Gene Hackman	José Ferrer	George Arliss
Anthony Hopkins	George C. Scott	Broderick Crawford	Warner Baxter
Jeremy Irons	John Wayne	Laurence Olivier	Emil Jannings
Daniel Day Lewis	Cliff Robertson	Ronald Colman	Emil Jannings
Dustin Hoffman	Rod Steiger	Fredric March	

An Ordered List

Adrien Brody	F. Murray Abraham	José Ferrer	Robert Donat
Al Pacino	Forest Whitaker	Kevin Spacey	Robert Duvall
Alec Guinness	Fredric March (2)	Laurence Olivier	Roberto Benigni
Anthony Hopkins	Gary Cooper (2)	Lee Marvin	Rod Steiger
Art Carney	Gene Hackman	Lionel Barrymore	Ronald Colman
Ben Kingsley	Geoffrey Rush	Marlon Brando (2)	Russell Crowe
Bing Crosby	George Arliss	Maximilian Schell	Sean Penn (2)
Broderick Crawford	George C. Scott	Michael Douglas	Sidney Poitier
Burt Lancaster	Gregory Peck	Nicolas Cage	Spencer Tracy (2)
Charles Laughton	Henry Fonda	Paul Lukas	Tom Hanks (2)
Charlton Heston	Humphrey Bogart	Paul Muni	Victor McLaglen
Clark Gable	Jack Lemmon	Paul Newman	Wallace Beery
Cliff Robertson	Jack Nicholson (2)	Paul Scofield	Warner Baxter
Daniel Day Lewis (2)	James Cagney	Peter Finch	William Holden
David Niven	James Stewart	Philip S. Hoffman	William Hurt
Denzel Washington	Jamie Foxx	Ray Milland	Yul Brynner
Dustin Hoffman (2)	Jeremy Irons	Rex Harrison	
Emil Jannings (2)	John Wayne	Richard Dreyfuss	
Ernest Borgnine	Jon Voight	Robert De Niro	

Hash Table

- $\text{Function}(\text{key}) = \text{value}$
- $F(\text{name}) =$
 $\text{index}(\text{last letter of the first name}) * 26 +$
 $\text{index}(\text{first letter of the last name})$
- $F(\text{John Wayne}) = F(n,w) = 13 * 26 + 22 = 360$

A Hash Table

58: Alec Guinness	271: Jack Lemmon	345: Dustin Hoffman (2)	496: Art Carney
64: Fredric March (2)	273: Jack Nicholson (2)	345: Charlton Heston	497: Robert Duvall
80: Ronald Colman	287: Lionel Barrymore	348: Ben Kingsley	497: Robert De Niro
81: Richard Dreyfuss	287: Yul Brynner	353: Sean Penn (2)	497: Robert Donat
91: David Niven	288: Russell Crowe	356: Maximilian Schell	505: Burt Lancaster
96: Rod Steiger	289: Michael Douglas	356: Kevin Spacey	516: Forest Whitaker
104: George Arliss	295: Emil Jannings (2)	359: Jon Voight	605: Rex Harrison
105: Wallace Beery	297: Daniel Day Lewis (2)	360: John Wayne	625: Humphrey Bogart
109: José Ferrer	297: Paul Lukas	365: Roberto Benigni	626: Gary Cooper (2)
109: Jamie Foxx	298: Paul Muni	397: Philip S. Hoffman	629: Henry Fonda
111: Gene Hackman	299: Paul Newman	443: Warner Baxter	631: Anthony Hopkins
116: Lee Marvin	301: Al Pacino	447: Peter Finch	632: Jeremy Irons
118: Laurence Olivier	304: Paul Scofield	454: Victor McLaglen	636: Ray Milland
122: George C. Scott	308: Denzel Washington	461: Spencer Tracy (2)	639: Gregory Peck
130: F. Murray Abraham	319: William Holden	470: Nicolas Cage	639: Sidney Poitier
147: Cliff Robertson	319: Tom Hanks (2)	470: James Cagney	641: Geoffrey Rush
158: Bing Crosby	319: William Hurt	479: Charles Laughton	
262: Broderick Crawford	339: Marlon Brando (2)	486: James Stewart	
266: Clark Gable	339: Adrien Brody	495: Ernest Borgnine	

Search In A Hash Table

58: Alec Guinness	271: Jack Lemmon	345: Dustin Hoffman (2)	496: Art Carney
64: Fredric March (2)	273: Jack Nicholson (2)	345: Charlton Heston	497: Robert Duvall
80: Ronald Colman	287: Lionel Barrymore	348: Ben Kingsley	497: Robert De Niro
81: Richard Dreyfuss	287: Yul Brynner	353: Sean Penn (2)	497: Robert Donat
91: David Niven	288: Russell Crowe	356: Maximilian Schell	505: Burt Lancaster
96: Rod Steiger	289: Michael Douglas	356: Kevin Spacey	516: Forest Whitaker
104: George Arliss	295: Emil Jannings (2)	359: Jon Voight	605: Rex Harrison
105: Wallace Beery	297: Daniel Day Lewis (2)	360: John Wayne	625: Humphrey Bogart
109: José Ferrer	297: Paul Lukas	365: Roberto Benigni	626: Gary Cooper (2)
109: Jamie Foxx	298: Paul Muni	397: Philip S. Hoffman	629: Henry Fonda
111: Gene Hackman	299: Paul Newman	443: Warner Baxter	631: Anthony Hopkins
116: Lee Marvin	301: Al Pacino	447: Peter Finch	632: Jeremy Irons
118: Laurence Olivier	304: Paul Scofield	454: Victor McLaglen	636: Ray Milland
122: George C. Scott	308: Denzel Washington	461: Spencer Tracy (2)	639: Gregory Peck
130: F. Murray Abraham	319: William Holden	470: Nicolas Cage	639: Sidney Poitier
147: Cliff Robertson	319: Tom Hanks (2)	470: James Cagney	641: Geoffrey Rush
158: Bing Crosby	319: William Hurt	479: Charles Laughton	
262: Broderick Crawford	339: Marlon Brando (2)	486: James Stewart	
266: Clark Gable	339: Adrien Brody	495: Ernest Borgnine	

John Wayne = $(n,w) = 13*26+22 = 360 = \text{FOUND}$

Search In A Hash Table

58: Alec Guinness	271: Jack Lemmon	345: Dustin Hoffman (2)	496: Art Carney
64: Fredric March (2)	273: Jack Nicholson (2)	345: Charlton Heston	497: Robert Duvall
80: Ronald Colman	287: Lionel Barrymore	348: Ben Kingsley	497: Robert De Niro
81: Richard Dreyfuss	287: Yul Brynner	353: Sean Penn (2)	497: Robert Donat
91: David Niven	288: Russell Crowe	356: Maximilian Schell	505: Burt Lancaster
96: Rod Steiger	289: Michael Douglas	356: Kevin Spacey	516: Forest Whitaker
104: George Arliss	295: Emil Jannings (2)	359: Jon Voight	605: Rex Harrison
105: Wallace Beery	297: Daniel Day Lewis (2)	360: John Wayne	625: Humphrey Bogart
109: José Ferrer	297: Paul Lukas	365: Roberto Benigni	626: Gary Cooper (2)
109: Jamie Foxx	298: Paul Muni	397: Philip S. Hoffman	629: Henry Fonda
111: Gene Hackman	299: Paul Newman	443: Warner Baxter	631: Anthony Hopkins
116: Lee Marvin	301: Al Pacino	447: Peter Finch	632: Jeremy Irons
118: Laurence Olivier	304: Paul Scofield	454: Victor McLaglen	636: Ray Milland
122: George C. Scott	308: Denzel Washington	461: Spencer Tracy (2)	639: Gregory Peck
130: F. Murray Abraham	319: William Holden	470: Nicolas Cage	639: Sidney Poitier
147: Cliff Robertson	319: Tom Hanks (2)	470: James Cagney	641: Geoffrey Rush
158: Bing Crosby	319: <u>William Hurt</u>	479: Charles Laughton	
262: Broderick Crawford	339: Marlon Brando (2)	486: James Stewart	
266: Clark Gable	339: Adrien Brody	495: Ernest Borgnine	

Adam Sandler = $(m,s) = 12 * 26 + 18 = 320 = \text{NOT FOUND}$

Back to Alignment

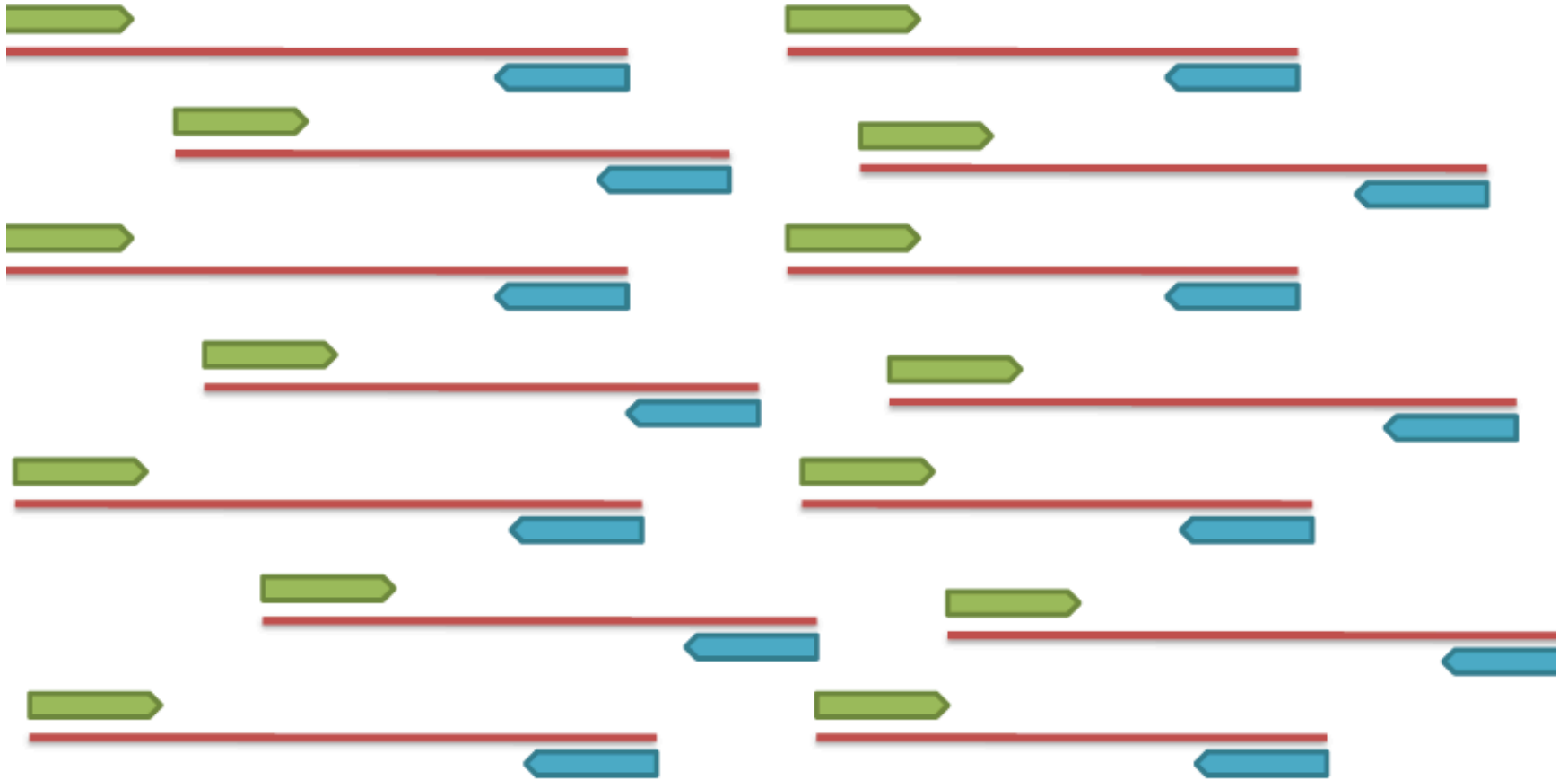
Reference sequence:

actgtagattag**ccgagtagctag**ctagtcgat



- Hashing is time and memory consuming for millions of reads and billion-base long reference
- Each read may be mapped to multiple positions

Paired End Reads



Paired End Reads

Paired Reads



Initial alignment to the reference genome



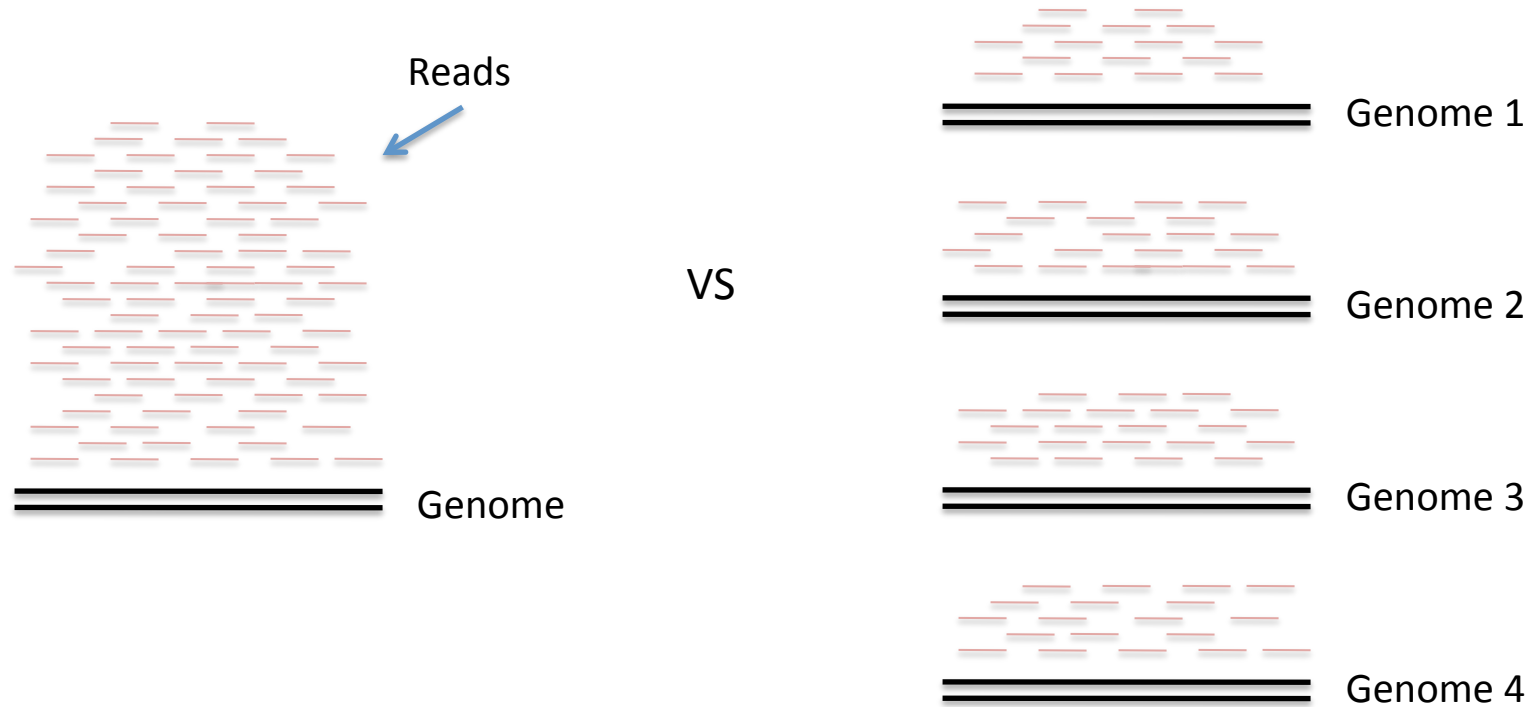
Paired end resolution



After Alignment

- Each read is mapped to reference genome with tolerated number of mismatches
 - Mismatches allow us to discover the individual variation
- Each site of reference genome is covered by multiple un-evenly distributed read
 - Some sites might not be covered

Coverage (High vs Low)



- Which one has more power to detect variations?

Genotype Calling from Sequence Data



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT
ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC
ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC
AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG
GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

A/C

Predicted Genotype

A Simple Model

At one site, N_A reads carry A, N_B reads carry B

$$N_A \sim \begin{cases} \text{Binomial}(n_A + n_B, 1 - \delta) & G = A / A \\ \text{Binomial}(n_A + n_B, 0.5) & G = A / B \\ \text{Binomial}(n_A + n_B, \delta) & G = B / B \end{cases}$$

Inference with no reads

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$$P(\text{reads}|\text{A/A})= 1.0$$

$$P(\text{reads}|\text{A/C})= 1.0$$

$$P(\text{reads}|\text{C/C})= 1.0$$

Possible Genotypes

Inference with short read data



GCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads}|\text{A/A}) = P(\text{C observed, read maps } |A/A)$

$P(\text{reads}|\text{A/C}) = P(\text{C observed, read maps } |A/C)$

$P(\text{reads}|\text{C/C}) = P(\text{C observed, read maps } |C/C)$

Possible Genotypes

Inference assuming error of 1%



GCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$$P(\text{reads}|\text{A/A})= 0.01$$

$$P(\text{reads}|\text{A/C})= 0.50$$

$$P(\text{reads}|\text{C/C})= 0.99$$

Possible Genotypes

As data accumulate ...



AGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTG
GCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$$P(\text{reads}|\text{A/A})= 0.0001$$

$$P(\text{reads}|\text{A/C})= 0.25$$

$$P(\text{reads}|\text{C/C})= 0.98$$

Possible Genotypes

As data accumulate ...



ATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCC
AGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTG
GCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$$P(\text{reads}|\text{A/A})= 0.000001$$

$$P(\text{reads}|\text{A/C})= 0.125$$

$$P(\text{reads}|\text{C/C})= 0.97$$

Possible Genotypes

As data accumulate ...



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT

ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC

AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG

GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads}|\text{A/A}) = 0.00000099$

$P(\text{reads}|\text{A/C}) = 0.0625$

$P(\text{reads}|\text{C/C}) = 0.0097$

Possible Genotypes

In the “end”



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT
ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC
ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC
AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG
GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGT**C**GATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCT**C**GACG-3'

Reference Genome

$$P(\text{reads}|\text{A/A})= 0.00000098$$

$$P(\text{reads}|\text{A/C})= 0.03125$$

$$P(\text{reads}|\text{C/C})= 0.000097$$

Not the “end” yet



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT
ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC
ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC
AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG
GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$$P(\text{reads}|\text{A/A}) = 0.00000098$$

$$P(\text{reads}|\text{A/C}) = 0.03125$$

$$P(\text{reads}|\text{C/C}) = 0.000097$$

Making a genotype call requires
combining sequence data with prior information

$$P(\text{Genotype}|\text{reads}) = \frac{P(\text{reads}|\text{Genotype})\text{Prior}(\text{Genotype})}{\sum_G P(\text{reads}|G)\text{Prior}(G)}$$

Not the “end” yet



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT
 ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC
 ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC
 AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG
 GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5' -ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads}|A/A) = 0.00000098$

$P(\text{reads}|A/C) = 0.03125$

$P(\text{reads}|C/C) = 0.000097$

$\text{Prior}(A/A) = 0.00034$

$\text{Prior}(A/C) = 0.00066$

$\text{Prior}(C/C) = 0.99900$

$P(A/A|\text{reads}) < 0.01$

$P(A/C|\text{reads}) = 0.175$

$P(C/C|\text{reads}) = 0.825$

Base Prior: every site has 1/1000 probability of varying

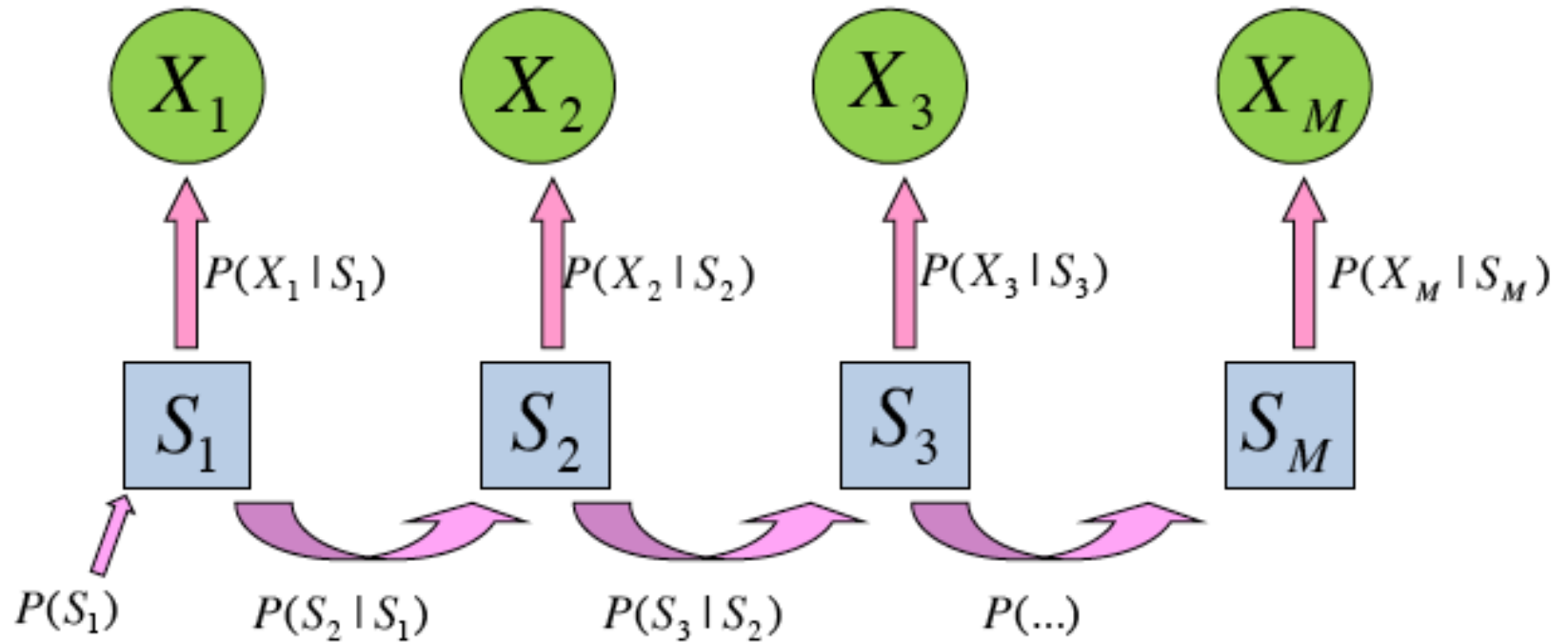
Prior Information

- Individual based prior
 - Equal probability of showing polymorphism
 - 1/1000 bases different from reference
 - Error Free and Poisson distribution
 - Single sample, single site
- Population based prior
 - Estimate frequency from many individuals
 - Multiple sample, single site
- Haplotype/Imputation based prior
 - Jointly model flanking SNPs, use haplotype information
 - Important for low coverage sequence data
 - Multiple samples, multiple sites

Haplotype-based Inference

- Start with some plausible configuration for each individual
- Use Markov model to update one individual condition on all others
- Repeat previous step many times
- Generate a consensus set of genotypes and haplotypes for each individual

Hidden Markov Model



Cartoon View of Shotgun Data

c	G	a	G	A	t	c	T	c	C	t	T	c	T	t	c	t	g	T	G	c
C	g	A	g	a	t	C	T	C	C	C	g	a	c	C	t	c	a	t	g	g
C	C	A	a	G	c	t	C	T	t	t	t	c	t	t	c	t	g	T	G	c
c	g	a	a	g	c	t	C	T	T	T	t	C	t	t	c	t	g	t	g	c
c	g	a	g	a	c	T	c	t	C	c	g	A	C	C	t	t	A	T	G	c
t	g	g	g	a	t	C	t	C	C	c	G	A	C	C	t	C	A	t	G	G
C	G	A	g	A	t	c	t	c	c	c	G	a	C	c	t	T	g	T	g	c
c	g	a	g	a	c	t	C	t	T	t	T	c	t	t	t	t	g	t	A	c
C	G	a	g	A	c	t	C	T	c	c	g	a	c	C	T	c	G	t	g	c
C	G	A	A	g	c	T	c	t	T	t	T	c	T	t	C	T	g	t	G	C
c	G	A	g	A	T	C	t	c	C	t	T	c	T	T	c	t	g	t	G	c
c	g	A	g	a	t	c	t	c	C	C	g	A	C	c	T	C	A	T	G	g
c	c	A	a	G	c	t	C	t	T	T	t	c	t	T	c	T	G	t	G	C
C	G	A	a	g	c	T	c	t	T	t	t	c	T	T	c	T	g	t	G	C
c	g	a	G	A	C	t	C	t	c	c	g	a	c	c	t	t	a	T	G	c
T	g	g	g	a	T	c	t	C	c	c	g	a	C	C	t	c	a	t	g	g
c	g	a	G	A	T	C	t	C	C	c	G	a	c	C	T	T	g	t	G	C
c	g	a	G	A	c	T	c	T	T	t	T	c	T	T	t	T	g	t	a	c
c	G	A	G	a	c	T	c	T	c	c	G	A	c	c	T	C	G	t	g	C
c	g	A	A	g	c	T	c	t	t	t	t	c	t	t	c	t	g	t	G	c

Cartoon View of the Method

Current update Sample

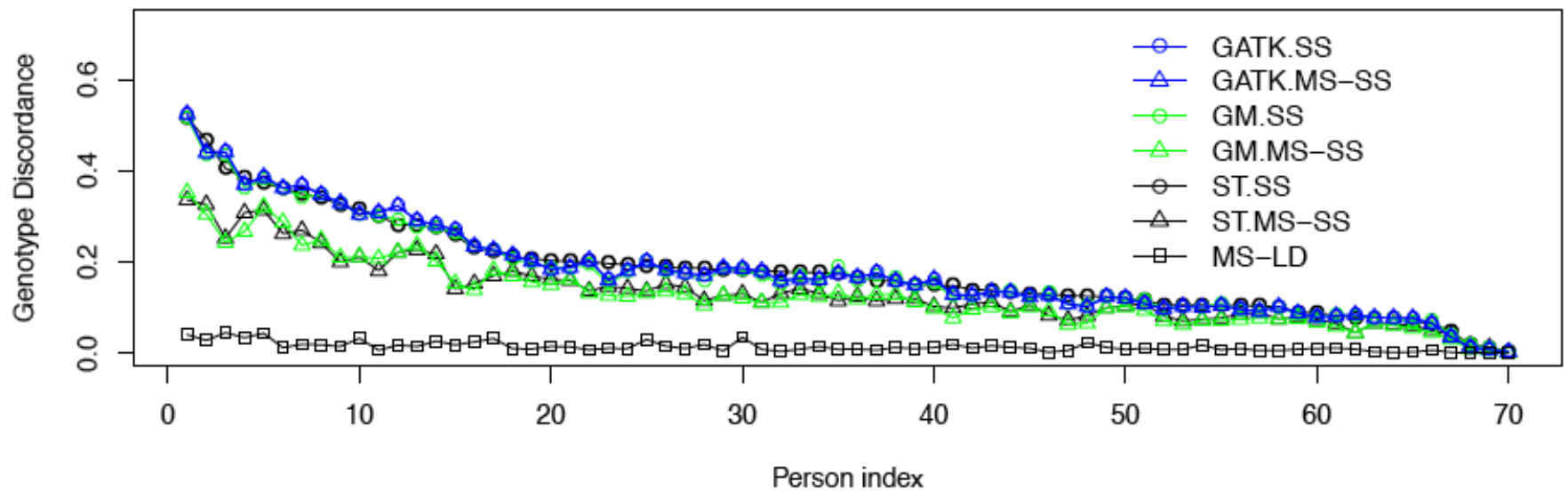
C	T	C	C	C	G	A	C	C	T	T	G	T	G	C
T	C	T	T	T	T	A	C	C	T	C	A	T	G	G

Known Haplotype Panel

C	T	C	C	C	G	A	C	C	T	C	A	T	G	G
C	T	C	C	C	G	A	C	C	T	T	G	T	G	C
T	C	T	T	T	T	C	T	T	T	T	G	T	A	C
T	C	T	C	C	G	A	C	C	T	C	G	T	G	C
T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
C	T	C	C	T	T	C	T	T	C	T	G	T	G	C
C	T	C	C	C	G	A	C	C	T	C	A	T	G	G
T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
T	C	T	T	T	T	C	T	T	C	T	G	T	G	C
T	C	T	C	C	G	A	C	C	T	T	A	T	G	C

Comparisons of Different Genotype Calling Methods

Low-Coverage (Overlapping Sites)



Simulation Results: Common Sites

- Detection and genotyping of sites
 - MAF > 5%, 2116 simulated SNP sites/Mb
- Detected polymorphic sites:

Sample	Detected Sites	Accuracy (all)	Accuracy(Het)
100	2102	98.5%	90.6%
200	2115	99.6%	99.4%
400	2116	99.8%	99.7%

Simulation Results: Rarer Sites

- Detection and genotyping of sites
 - MAF > 1-2%, 425 simulated SNP sites/Mb
- Detected polymorphic sites:

Sample	Detected Sites	Accuracy (all)	Accuracy(Het)
100	139	98.6%	92.9%
200	213	99.4%	95.0%
400	343	99.6%	95.9%