# GRCh38 Centromere Reference Models
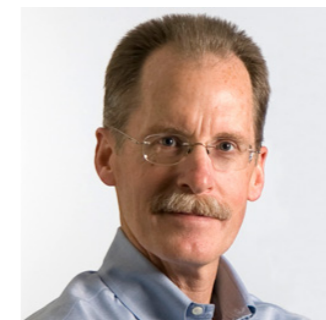
Karen H. Miga

University of California, Santa Cruz
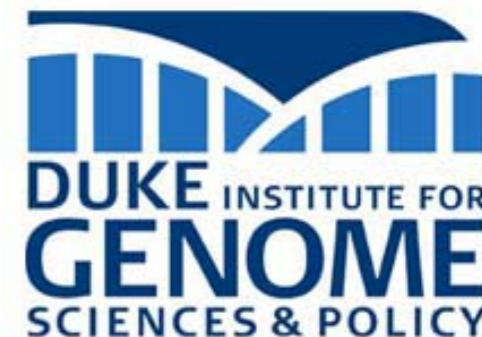
**Baskin Engineering UC SANTA CRUZ**
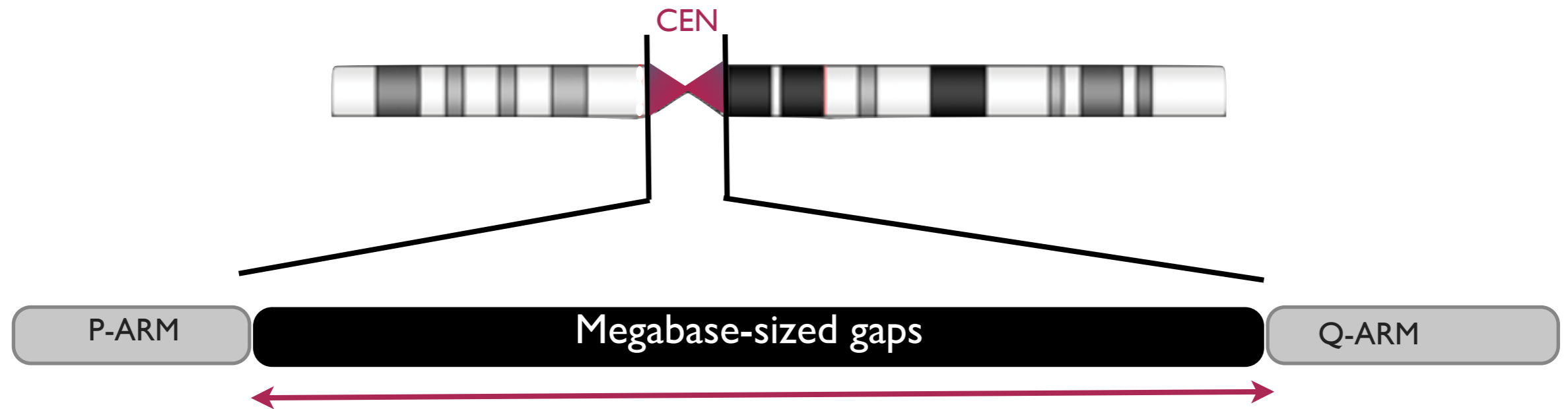
Jim Kent          Huntington F. Willard

**DUKE INSTITUTE FOR GENOME SCIENCES & POLICY**
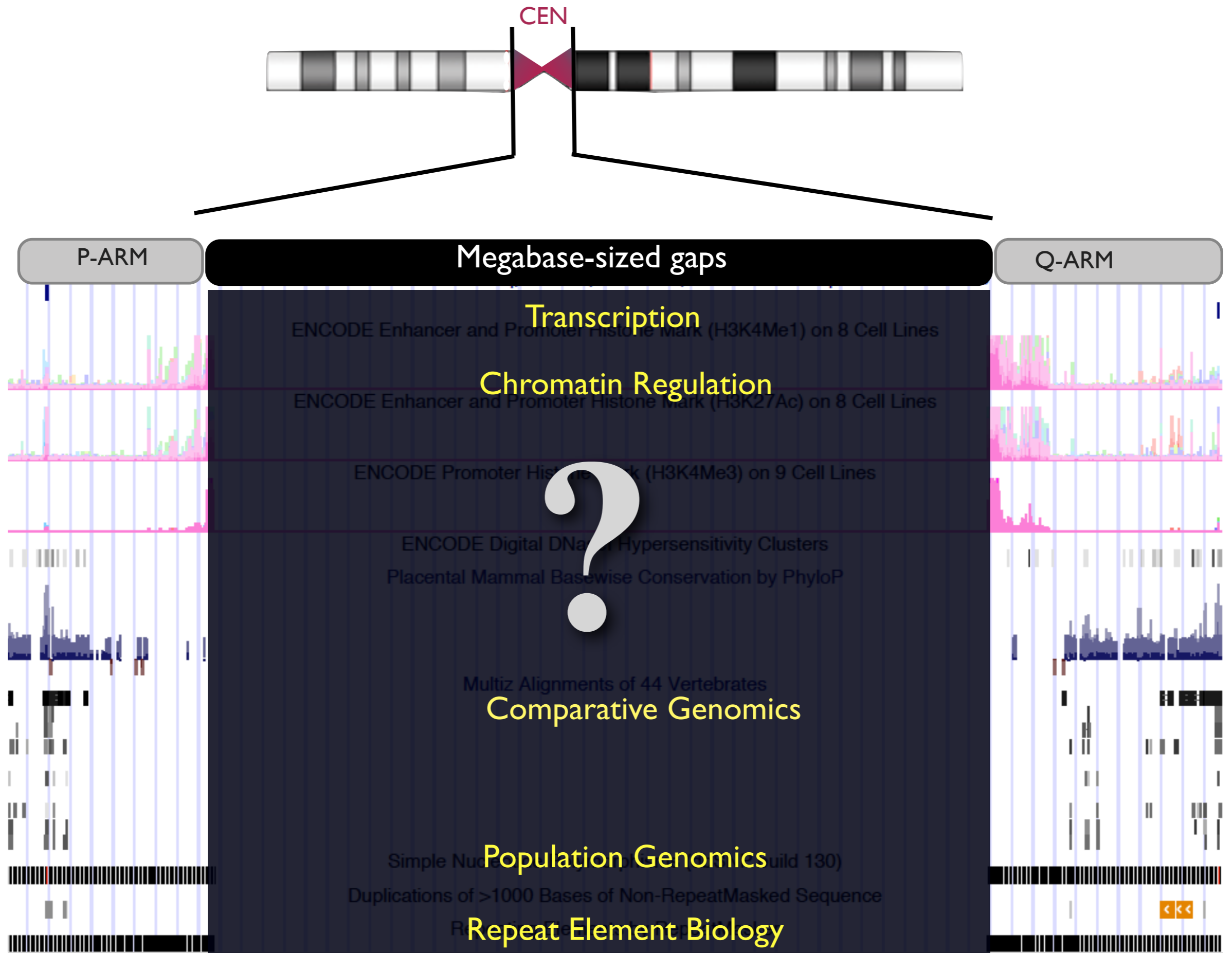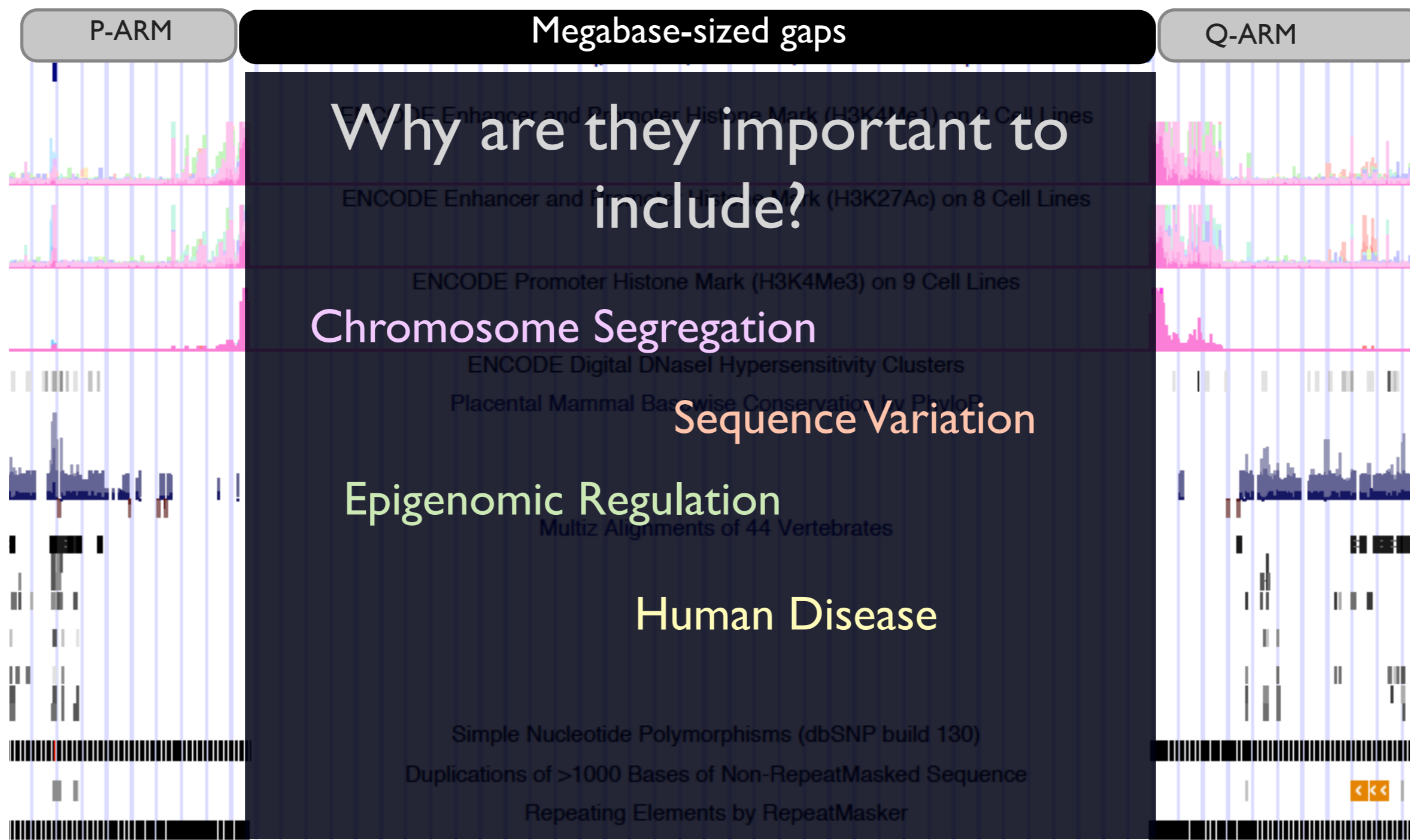
# Human centromeres are currently defined by gaps in the reference assembly

# Human centromeres are currently defined by gaps in the reference assembly



CEN

| P-ARM | Megabase-sized gaps | Q-ARM |

**Transcription**

ENCODE Enhancer and Promoter Histone Mark (H3K4Me1) on 8 Cell Lines

**Chromatin Regulation**

ENCODE Enhancer and Promoter Histone Mark (H3K27Ac) on 8 Cell Lines

ENCODE Promoter Histone Mark (H3K4Me3) on 9 Cell Lines

?

ENCODE Digital DNase Hypersensitivity Clusters

Placental Mammal Basewise Conservation by PhyloP

Multiz Alignments of 44 Vertebrates

**Comparative Genomics**

Simple Nucleotide Polymorphisms (dbSNP build 130)

**Population Genomics**

Duplications of >1000 Bases of Non-RepeatMasked Sequence

**Repeat Element Biology**

# Human centromeres are currently defined by gaps in the reference assembly

CEN

| P-ARM | Megabase-sized gaps | Q-ARM |

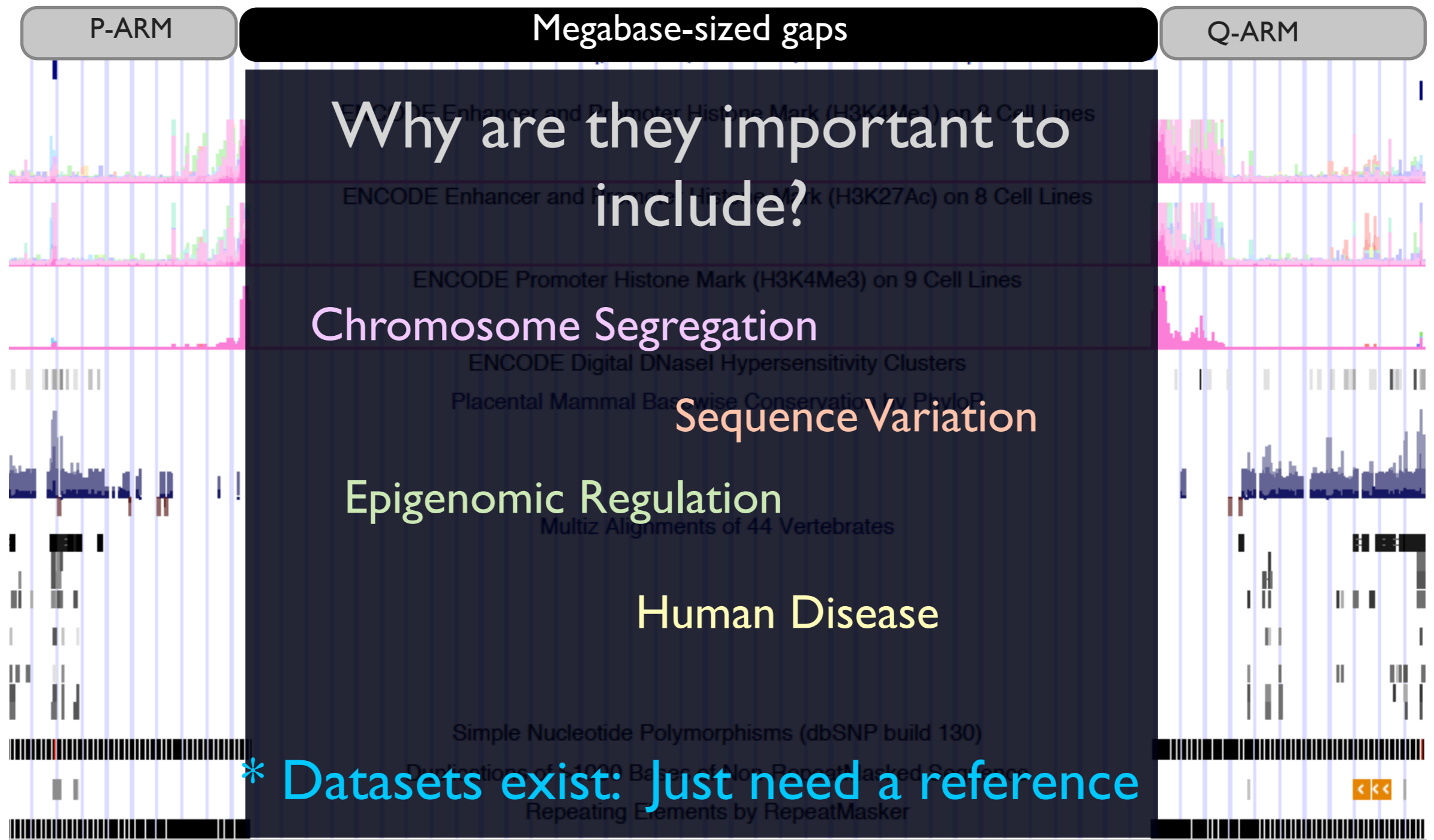## Why are they important to include?

Chromosome Segregation
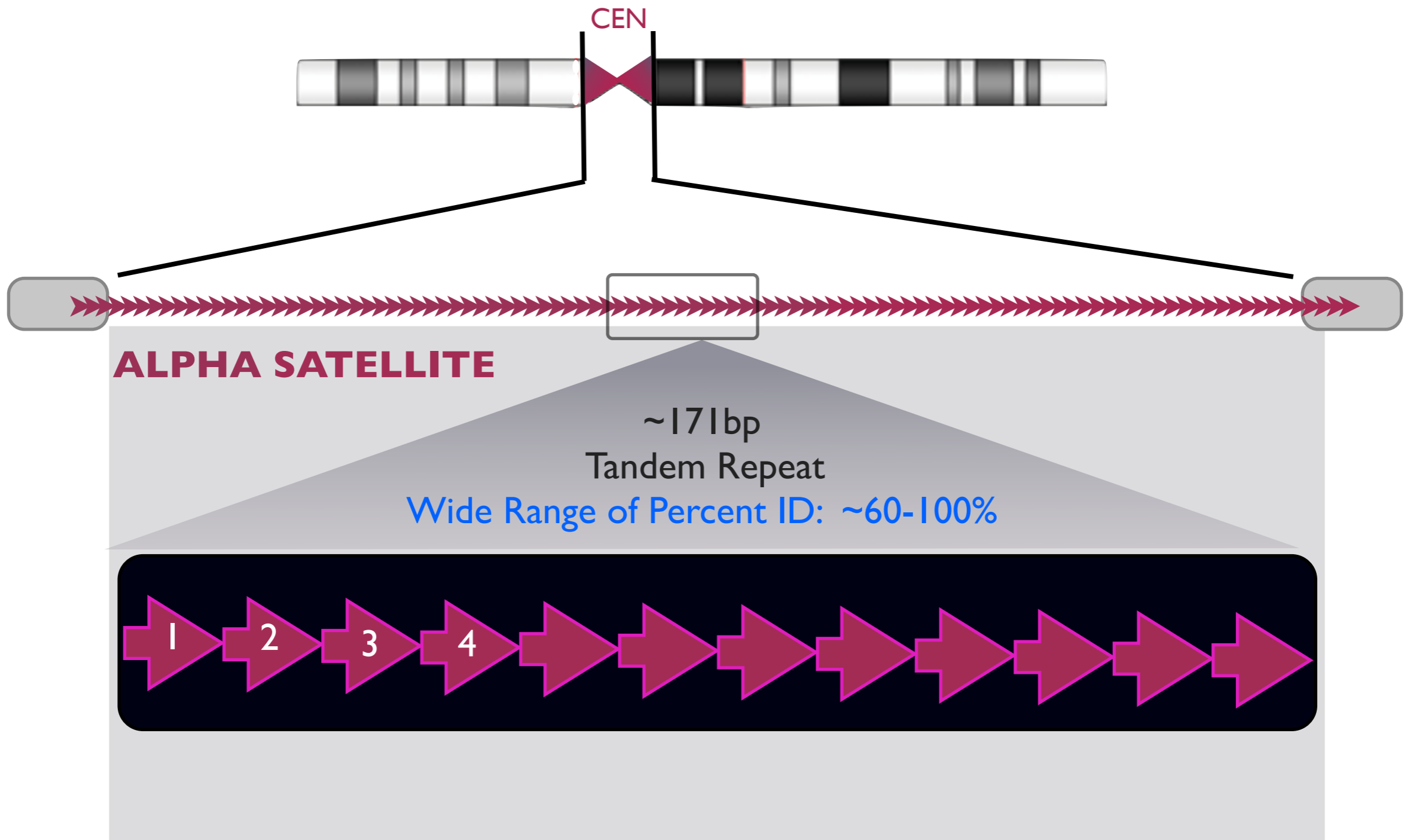
Sequence Variation

Epigenomic Regulation

Human Disease

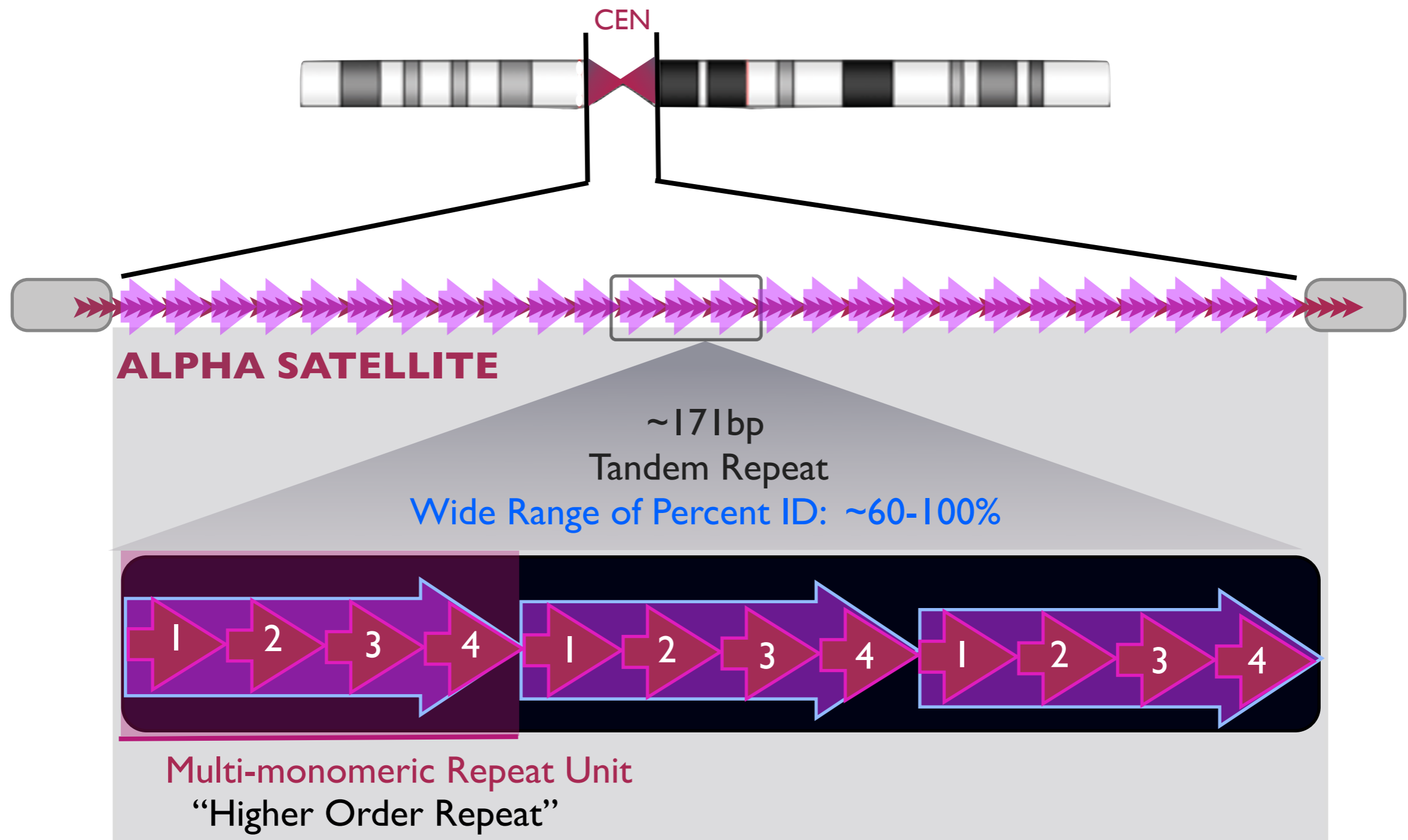# Human centromeres are currently defined by gaps in the reference assembly

CEN



| P-ARM | Megabase-sized gaps | Q-ARM |

## Why are they important to include?

**Chromosome Segregation**

**Sequence Variation**

**Epigenomic Regulation**

**Human Disease**

\* Datasets exist:  Just need a reference

Alpha Satellite define all normal human centromeres

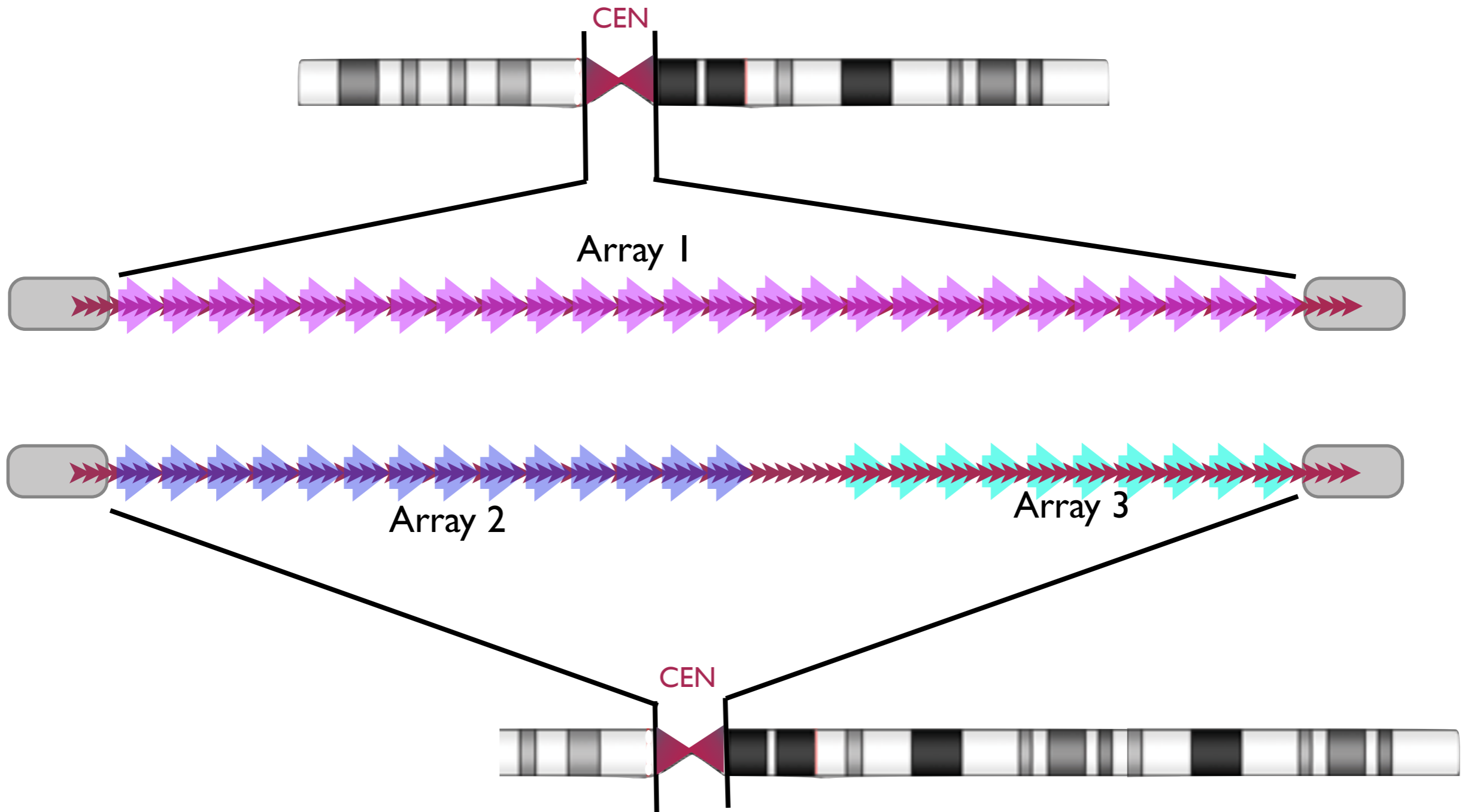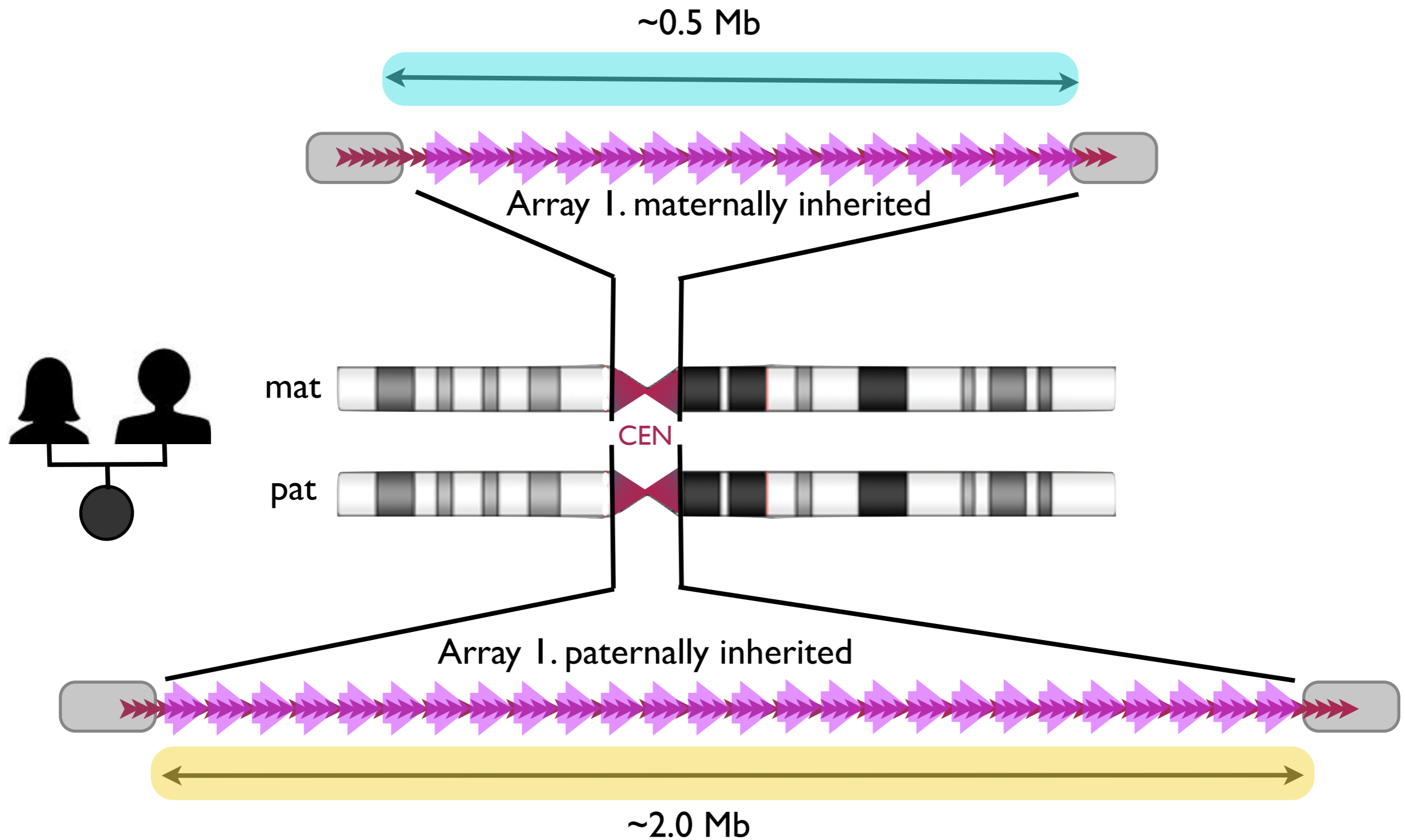Alpha Satellite repeats (or monomers) are commonly found in long arrays of near-identical higher order repeats

Alpha Satellite repeats (or monomers) are commonly found in long arrays of near-identical higher order repeats

Satellite DNA are the primary sequence in each gap



ALPHA SATELLITE

~171bp
Tandem Repeat
Wide Range of Percent ID: ~60-100%

Narrow Range of Percent ID: **94% - 100%**

Alpha Satellite repeats (or monomers) are commonly found in long arrays of near-identical higher order repeats

# Each chromosome has a different centromeric sequences

# Higher–order arrays vary between individuals



~0.5 Mb

Array 1. Individual A

CEN

Array 1. Individual B

~2.0 Mb

# Higher-order arrays can vary between homologous chromosomes in the same individual



~0.5 Mb

Array 1. maternally inherited

mat

CEN

pat

Array 1. paternally inherited

~2.0 Mb

# Model of Centromere Sequence Organization



CEN

Array 1

Array 2

Repeat Elements

Inter-Array (Non-satellite) Sequence

# Model of Centromere Sequence Organization



Goal: To generate a reference that models alpha satellite (and adjacent non−satellite) sequences within each centromeric gap

# 2. Reformat sequences observed in each read library into linear reference model



① Constructing Read Libraries for each HOR array

② LinearSat Software to Convert Reads to Linear Reference Models

③ Scaffold Reference Models and HuRef assembled contigs using mate pairs

# Constructing Read Libraries for each HOR array

**Whole Genome Sequencing Data Single Individual**

**HuRef Genome**

Centromeric database construction from reads containing alpha satellite repeats.
(2.6% of the human genome)

Determine chromosome-specific organization of alpha variants into higher order repeats.

Build statistical models to generate faux centromere sequence that will serve as a target for mapping centromeric reads.

# Higher Order Repeat Prediction



**HuRef Genome (8x Coverage)**

**DATA COMPRESSION**

**IDENTICAL MONOMERS**

**1 base change**

Similarity Clustering:
Defining Epsilon Neighborhood

**SEQUENCE RELATIONSHIP**

# Higher Order Repeat Prediction

**PHYSICAL RELATIONSHIP**

**SEQUENCE RELATIONSHIP**

# Higher Order Repeat Prediction



**ADJACENCY MATRIX**
Alpha Satellite Monomers

2-away

1-away

Sanger Read (Ave: 3 monomers)

# Higher Order Repeat Prediction

# Higher Order Repeat Prediction



**Viterbi greedy-algorithm second order markov model**

# Higher Order Repeat Prediction
## Determine Chromosome Specificity:

# Higher Order Repeat Prediction
## Determine Chromosome Specificity:

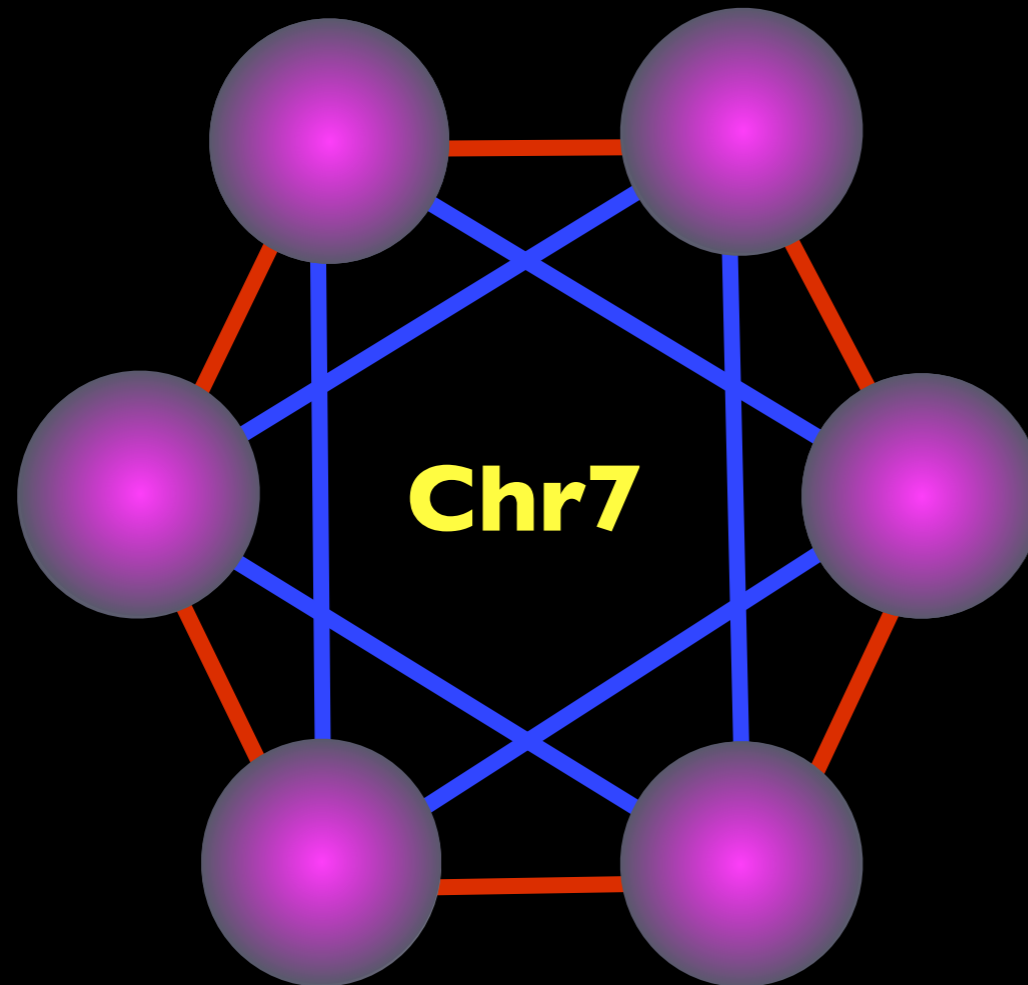Flow Sorted Chromosome Alignment/Enrichment
344 Mb of Alpha Satellite from 15 Chromosomes
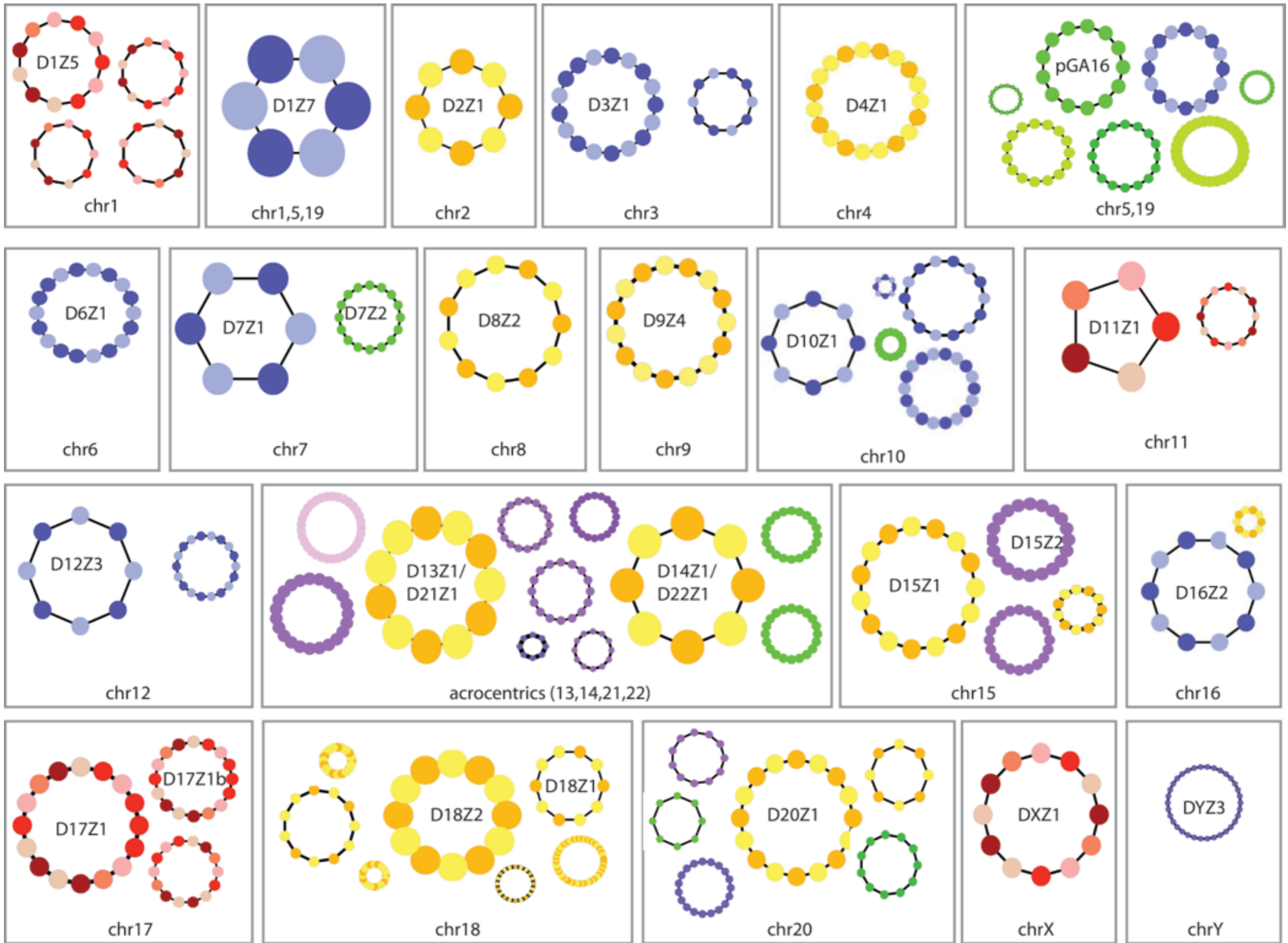
Chr7

1 2 3 4 5 6 1 2 3 4 5 6

# Higher Order Repeat Prediction
## Determine Chromosome Specificity:



## Flow Sorted Chromosome Alignment/Enrichment
344 Mb of Alpha Satellite from 15 Chromosomes

## Experimental Evidence
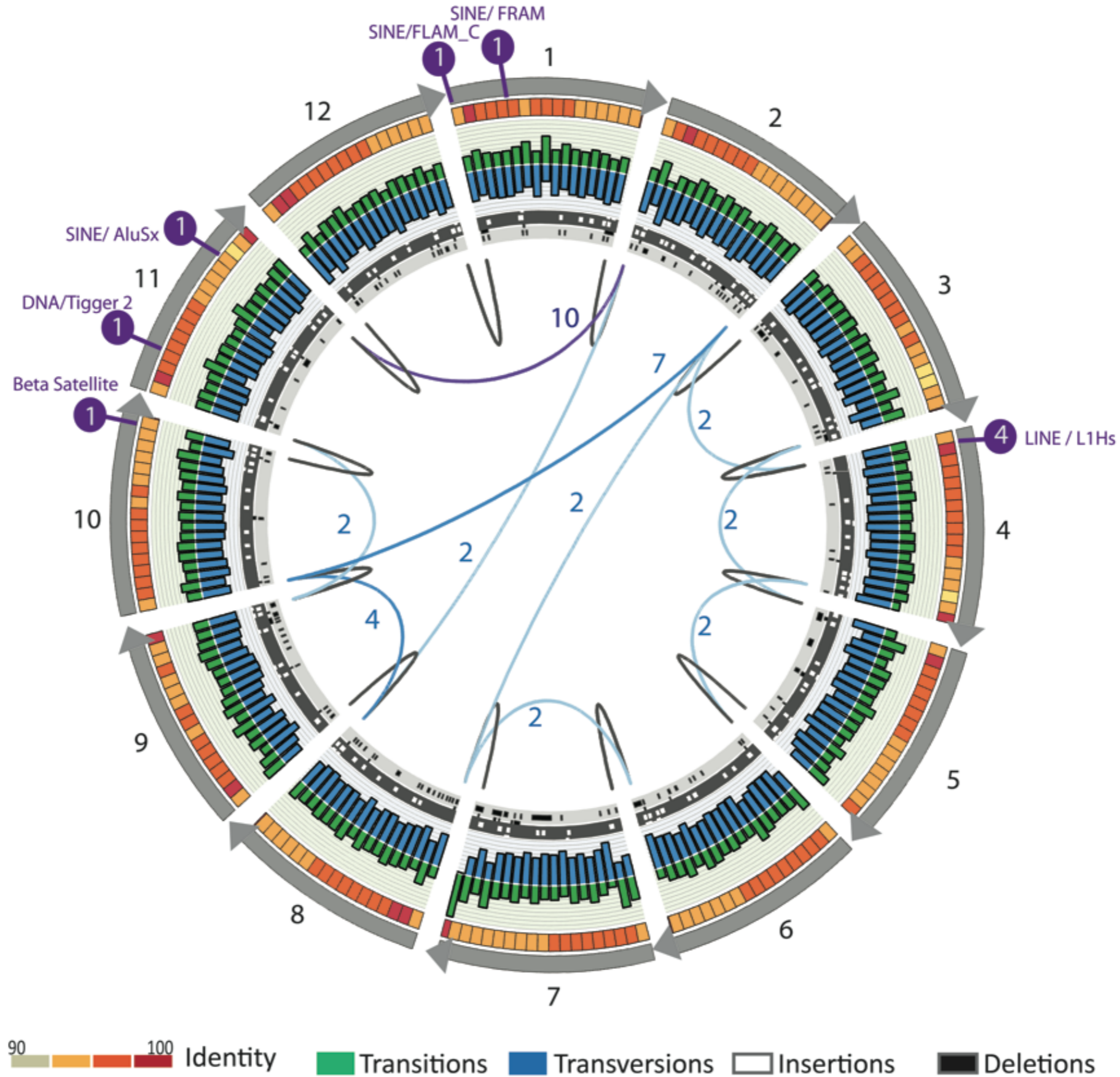FISH Hybridization and Screening Somatic Cell Hybrid Panel

# Higher Order Repeat Prediction
## Determine Chromosome Specificity:

**Chr7**

### Flow Sorted Chromosome Alignment/Enrichment
344 Mb of Alpha Satellite from 15 Chromosomes

### Experimental Evidence
FISH Hybridization and Screening Somatic Cell Hybrid Panel

### Paired Reads
"Anchor" to adjacent mapped HuRef contigs

1 2 3 4 5 6 1 2 3 4 5 6

chr1

D1Z5

D1Z7

chr1,5,19

D2Z1

chr2

D3Z1

chr3

D4Z1

chr4

pGA16

chr5,19

D6Z1

chr6

D7Z1

D7Z2

chr7

D8Z2

chr8

D9Z4

chr9

D10Z1

chr10

D11Z1

chr11

D12Z3

chr12

D13Z1/D21Z1

D14Z1/D22Z1

acrocentrics (13,14,21,22)

D15Z1

D15Z2

chr15

D16Z2

chr16

D17Z1

D17Z1b

chr17

D18Z2

D18Z1

chr18

D20Z1

chr20

DXZ1

chrX

DYZ3

chrY

# Alpha Satellite Array (DXZ1) on Chromosome X
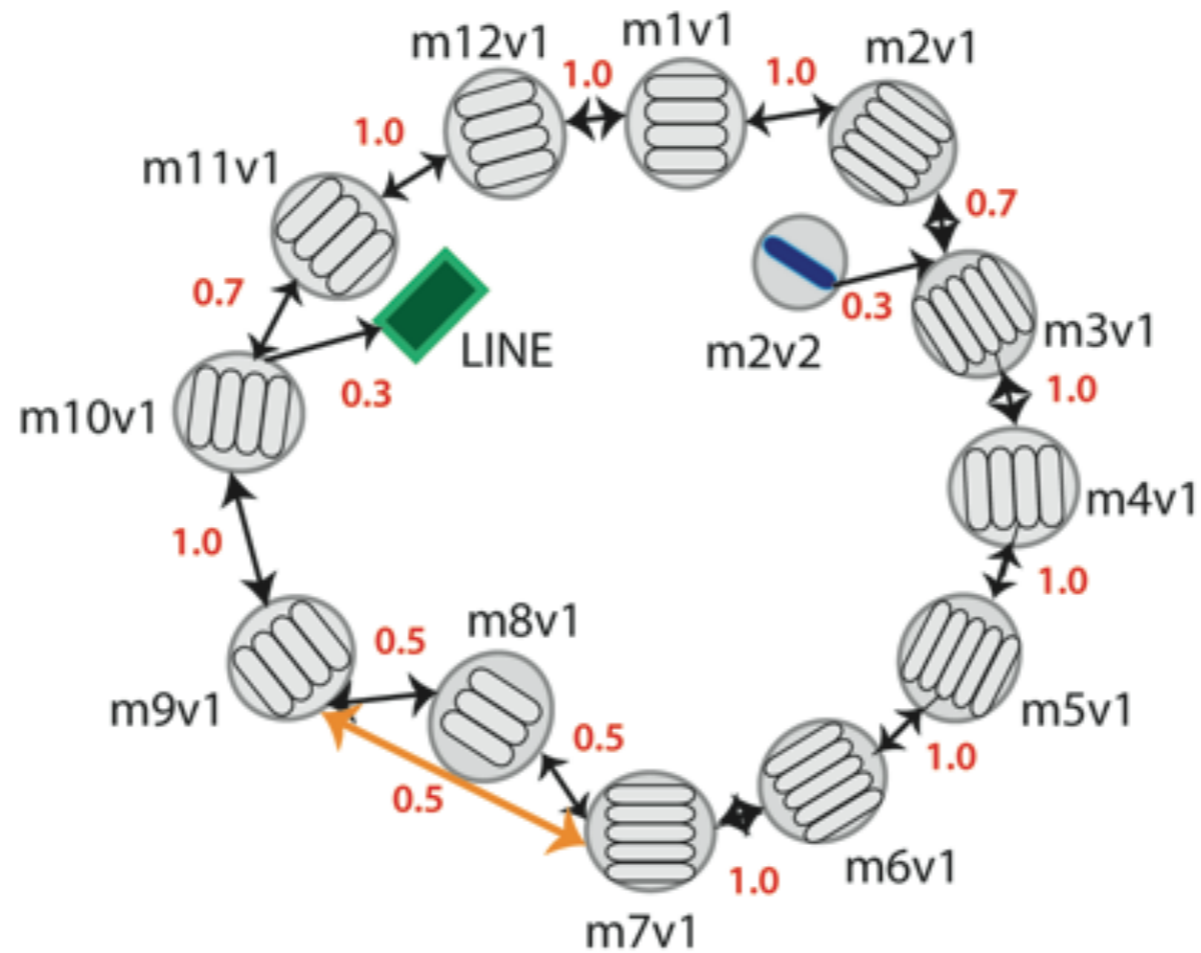
# LinearSat Software to Convert Reads to Linear Reference Models

## ② LinearSat Software to Convert Reads to Linear Reference Models

# LinearSat

- 2nd Order Markov Chain
- Length determined by Normalized Read Coverage
- Sensitive to low coverage
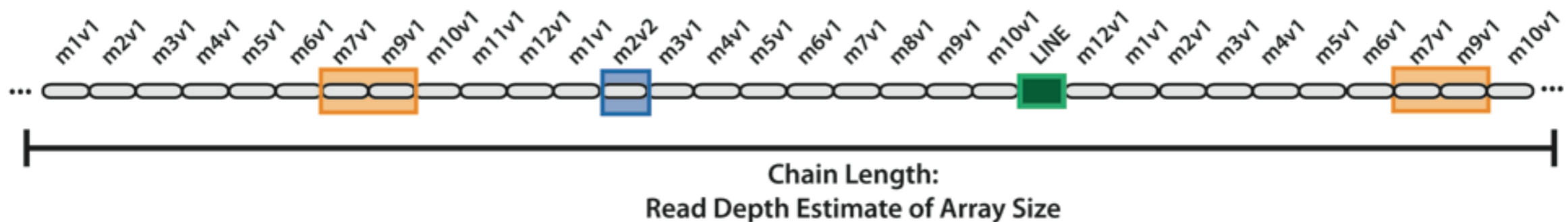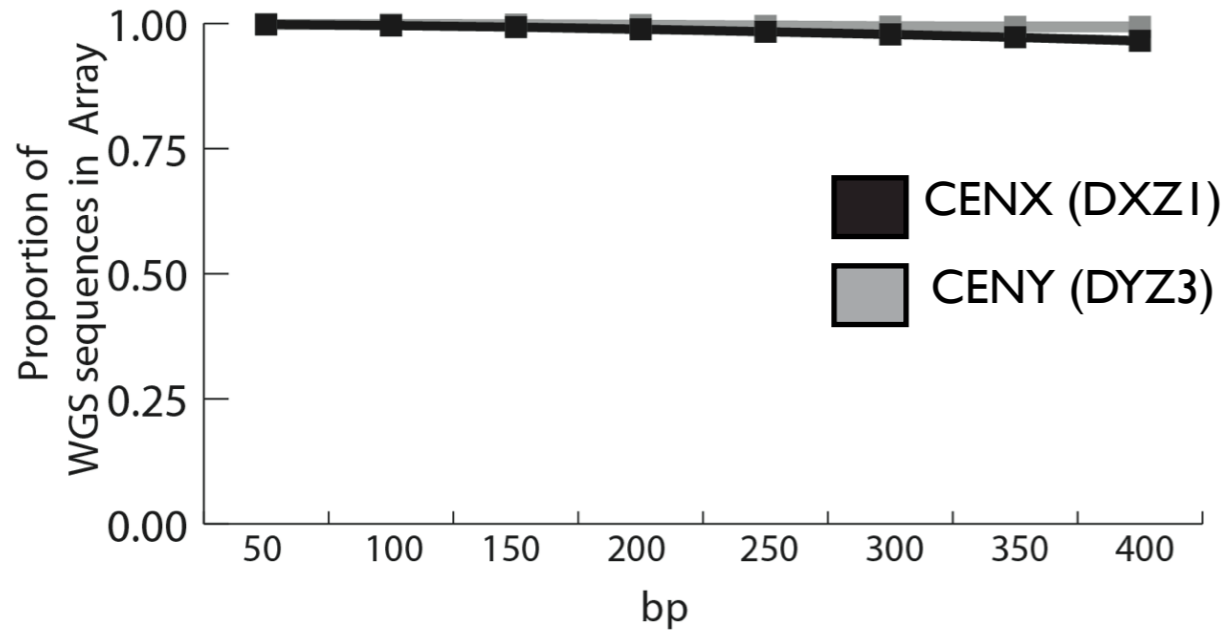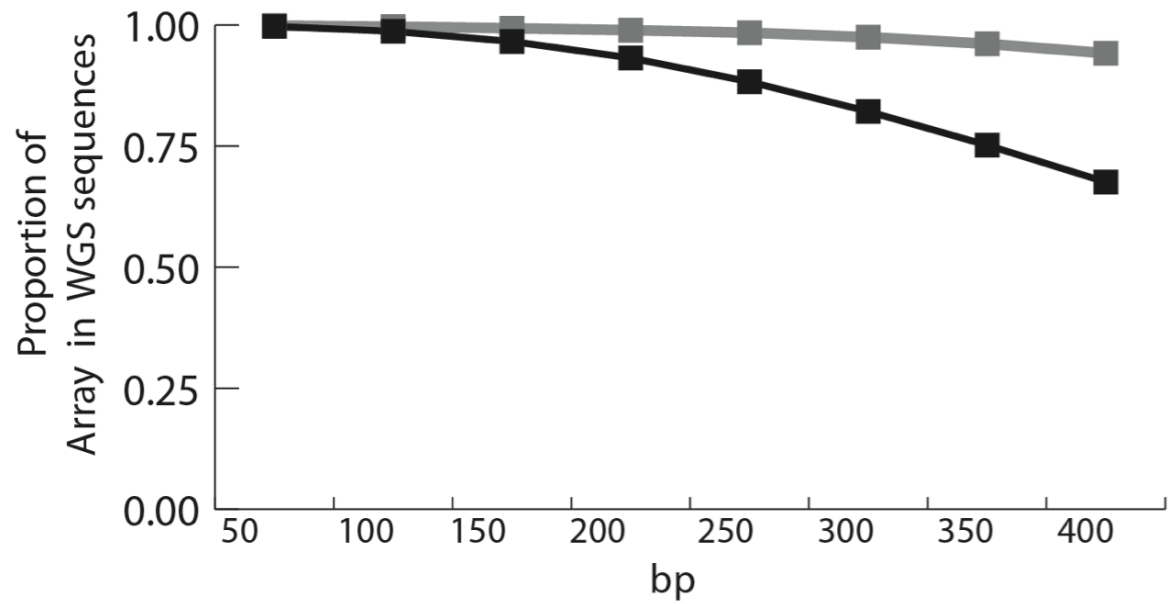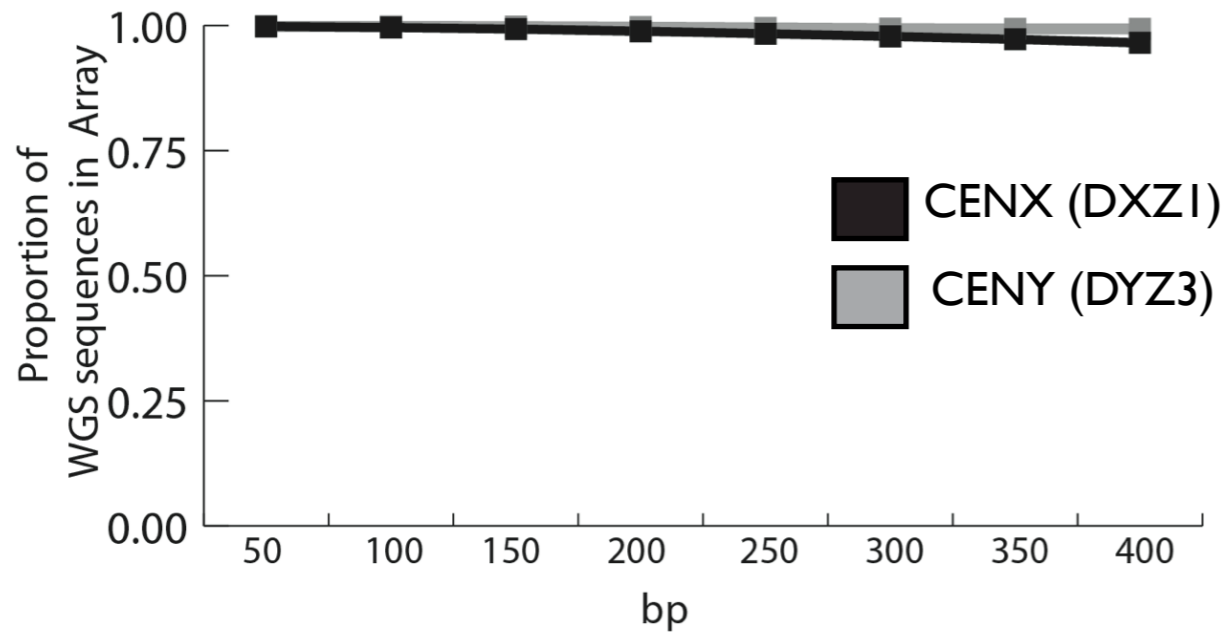-  Implemented with information on HOR repeat structure

**Chain Length:**
**Read Depth Estimate of Array Size**

## LinearSat

- 2nd Order Markov Chain
- Length determined by Normalized Read Coverage
- Sensitive to low coverage
- Implemented with information on HOR repeat structure

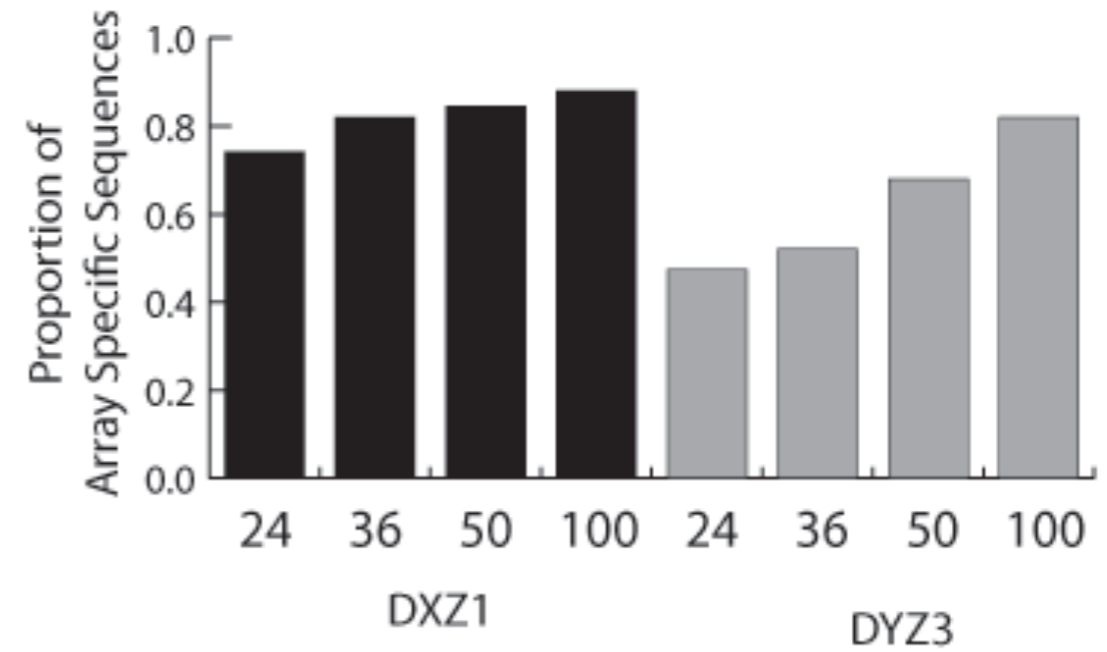Not the "true" long-range organization, yet adequately represents the alpha satellite array sequence

Test each satellite reference model to ensure that sequence variation is observed as expected within the initial read dataset
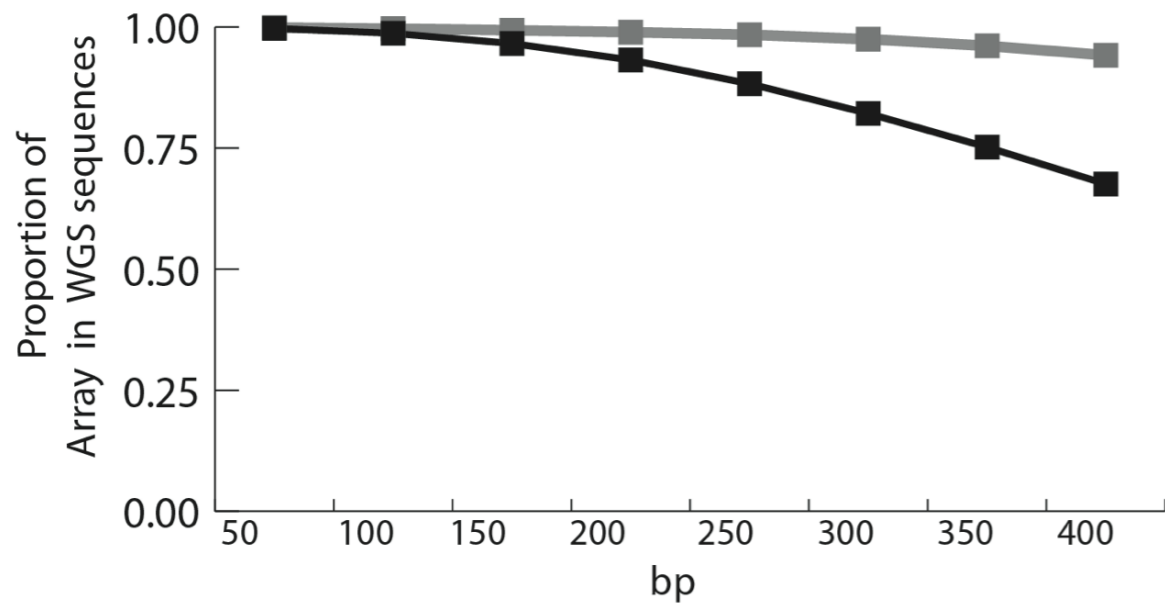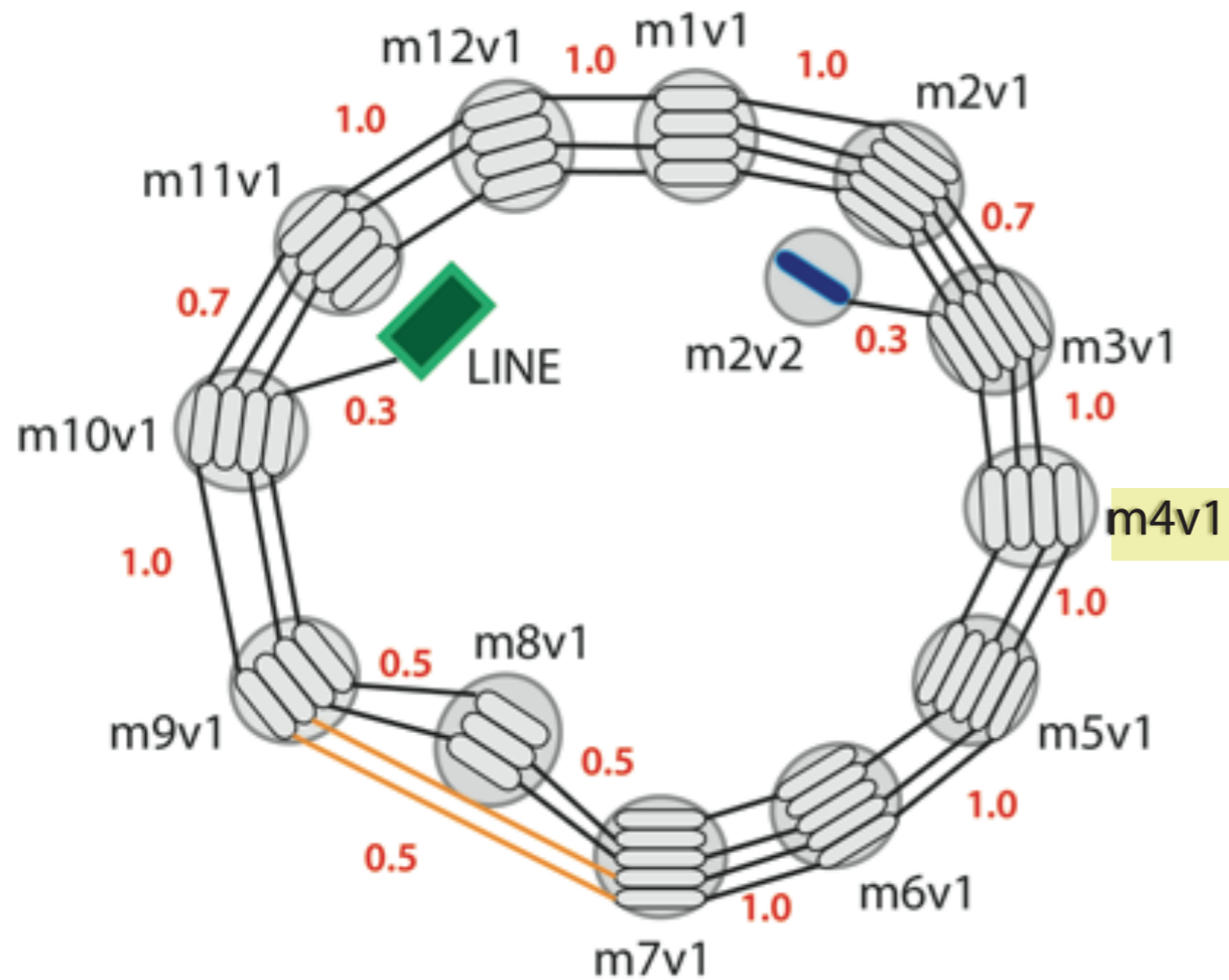
sensitivity

Proportion of WGS sequences in Array

CENX (DXZ1)
CENY (DYZ3)

Mappability

Proportion of Array Specific Sequences

24 36 50 100 — DXZ1
24 36 50 100 — DYZ3

specificity

Proportion of Array in WGS sequences

Evaluate as a Potential Read Mapping Target

Test each satellite reference model to ensure that sequence variation is observed as expected within the initial read dataset

# GRCh38 Data Structure
## Level 1: Repeat Components



Database all unique sequence in each array graph

>m4v1 4 identical monomers

CACTTGCAGATTCTACAAAAAGAGTGCTTCAAAAC
TGCTCTGTCAAAAGGAAGGTTCAACTCTGTTACTT
GAGTACACACATCACAAGGAAGTTTCTGAGAATGC
TTCTGTCTGGTTTTTAGGAGAAGATATTTCCTTTT
TCAACATAGGCCTCAAAGCGCTGCAAATGTCCACT
TCC

## Deposit (NCBI, TPA) individual component fasta sequence of each centromere reference model

# GRCh38 Data Structure
# Level 2: AGP describing the order of sequence components



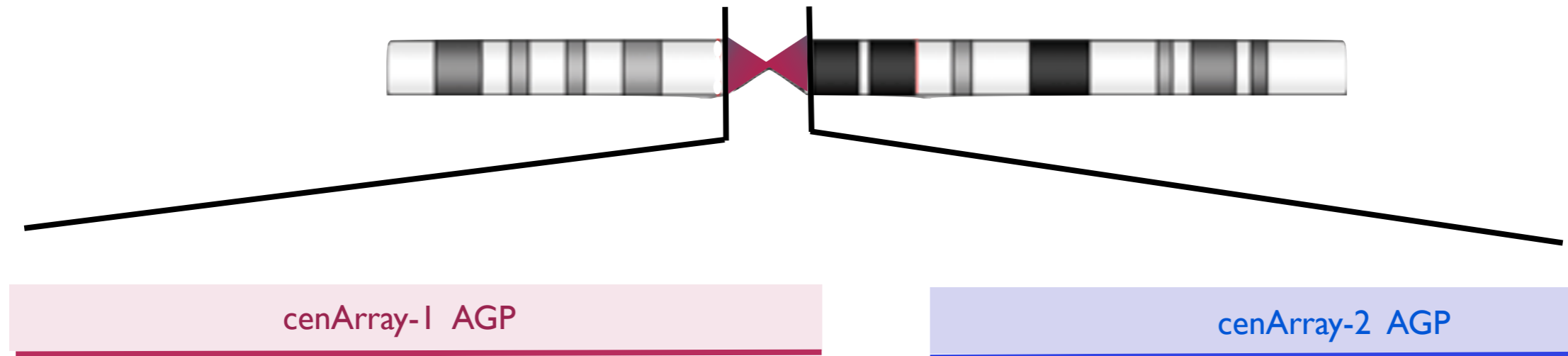Level 2: Centromere Reference Model "cenArray"  AGP

| | Array Name | Array Start | Array End | UID | UID | Level 1 Entry | L1 Start | L1 End | Level 1 | Ori |
|---|---|---|---|---|---|---|---|---|---|---|
| → | cenArray | 129970 | 130138 | 759 | O | m7v1 | 1 | 169 | | + |
| → | cenArray | 130139 | 130309 | 760 | O | m9v1 | 1 | 171 | | + |
| → | cenArray | 130310 | 130608 | 761 | O | m10v1 | 1 | 170 | | + |
| → | cenArray | 130609 | 130708 | 762 | N | m11v1 | 1 | 171 | | + |
| → | cenArray | 130709 | 130878 | 763 | O | m12v1 | 1 | 170 | | + |
| → | cenArray | 130879 | 131049 | 764 | O | m1v1 | 1 | 171 | | + |

Array Coordinates

Level 1 Sequence

# GRCh38 Data Structure
## Level 3: AGP describing the order of Array components



cenArray-1  AGP

cenArray-2  AGP

Single centromeric gap can contain more than one array

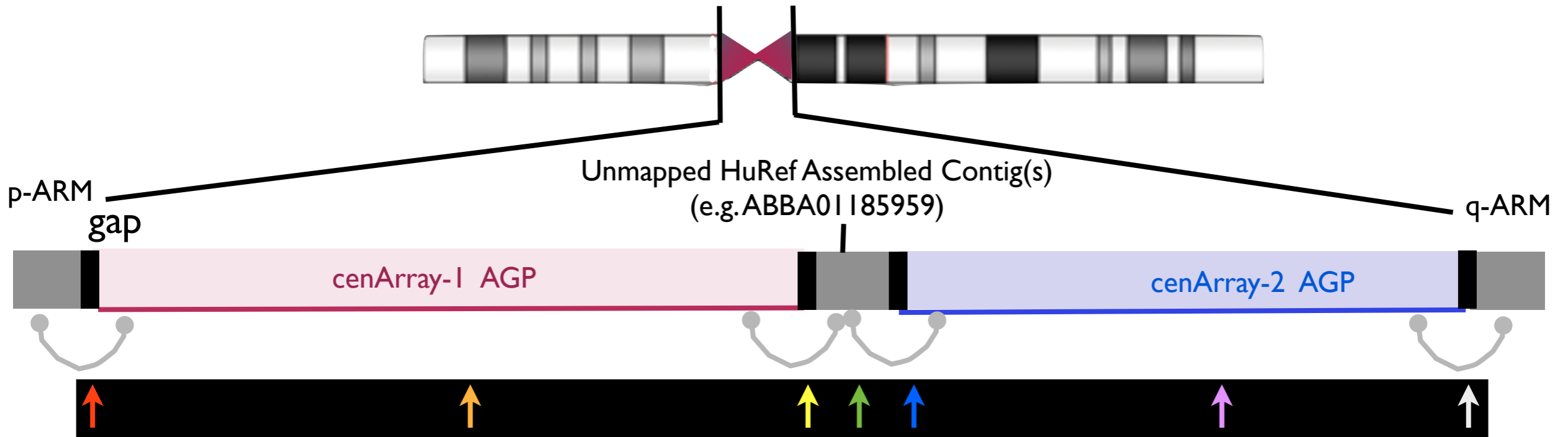**3** Scaffold Reference Models and HuRef assembled contigs using mate pairs

Single centromeric gap can contain more than one array

Scaffolding Order: Weighted by Mate Pairs

-- Bundled paired read information informs array component order
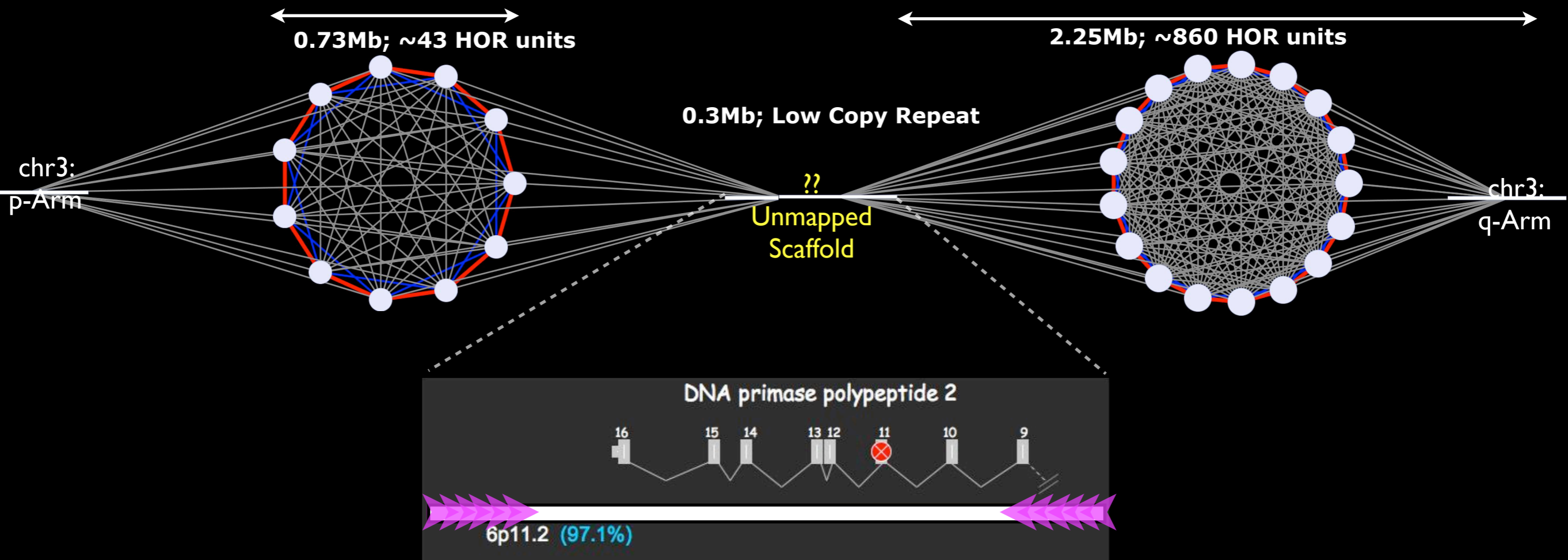
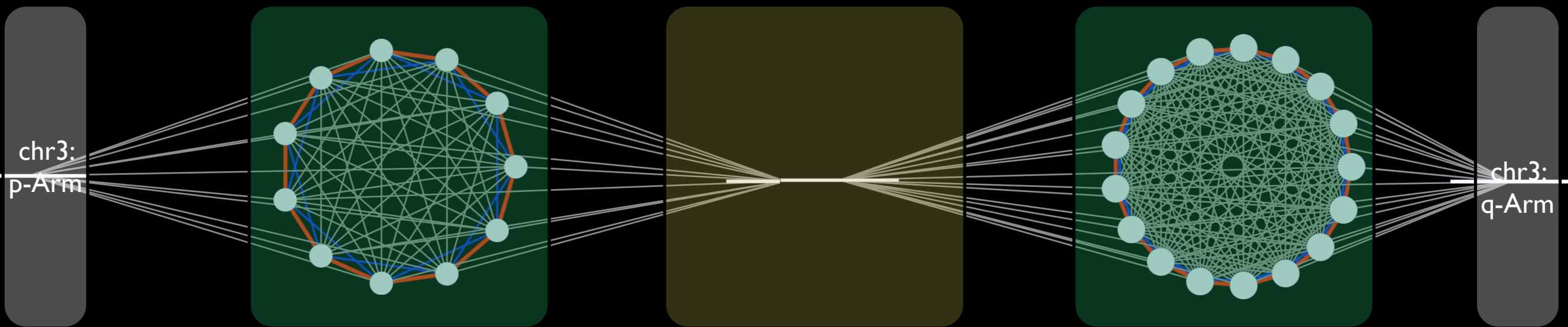# GRCh38 Data Structure
# Level 3: AGP describing the order of Array components

Unmapped HuRef Assembled Contig(s)
(e.g. ABBA01185959)

p-ARM
gap
q-ARM

cenArray-1 AGP

cenArray-2 AGP

| Array Name | Array Start | Array End | UID | UID | Level 1 Entry | L1 Start | L1 End | Level 1 Ori |
|---|---|---|---|---|---|---|---|---|
| cenArray | 0 | 100 | 1 | N | p-ARM gap | 1 | 100 | paired-read |
| cenArray | 101 | 2836899 | 2 | O | cenArray-1 | 1 | 2836798 | + |
| cenArray | 2836899 | 2837000 | 3 | N | gap | 1 | 100 | paired-read |
| cenArray | 2837000 | 2842055 | 4 | O | ABBA011859591 | 1 | 5055 | paired-read |
| cenArray | 2836899 | 2837000 | 5 | N | gap | 1 | 100 | paired-read |
| cenArray | 2837001 | 4369982 | 6 | O | cenArray-2 | 1 | 1532981 | + |
| cenArray | 4369983 | 4370083 | 7 | N | gap | 1 | 100 | paired-read |

CEN Coordinates

Level 2 Array Sequences

# Scaffolding Problem:
# Order Elements by Paired Reads

# Scaffolding Problem:
# Order Elements by Paired Reads

An Initial Draft of Human Centromere Sequence Composition

Alpha Satellite Reference Models: ~60 Mb (59571670 bp)

An Initial Draft of Human Centromere Sequence Composition

Redundant Arrays: Cannot assign to a specific chromosome that is normalized appropriately

# Adding Genome Annotation in Centromere Regions

CEN



Contribute more than just read mapping targets!

Genomic Reference Sequence

P-ARM

Q-ARM

Transcription

ENCODE Enhancer and Promoter Histone Mark (H3K4Me1) on 8 Cell Lines

Chromatin Regulation

ENCODE Enhancer and Promoter Histone Mark (H3K27Ac) on 8 Cell Lines

ENCODE Promoter Histone Mark (H3K4Me3) on 9 Cell Lines

?

ENCODE Digital DNase Hypersensitivity Clusters

Placental Mammal Basewise Conservation by PhyloP

Multiz Alignments of 44 Vertebrates

Comparative Genomics

Population Genomics

Simple Nucleotide Polymorphisms (dbSNP build 130)

Duplications of >1000 Bases of Non-RepeatMasked Sequence

Repeat Element Biology

# Query existing datasets that contain centromeric sequence



# K-mer frequency comparison and confident assignment of annotation back to the reference coordinates

## Y Chromosome DNA Haplotyping Suggests That Most European and Asian Men Are Descended from One of Two Males

REBECCA OAKEY[1] AND CHRIS TYLER-SMITH[2]

CRC Chromosome Molecular Biology Group, Department of Biochemistry, University of Oxford,
South Parks Road, Oxford OX1 3QU, United Kingdom

Received November 15, 1989; revised February 23, 1990

HuRef k-mers (24mers) useful in predicting array length across ~400 male individuals



European
Asian

Size of CEN Y array (Mb)

# of Individuals

Size of CEN Y array (Mb)

# Y Chromosome DNA Haplotyping Suggests That Most European and Asian Men Are Descended from One of Two Males

REBECCA OAKEY[1] AND CHRIS TYLER-SMITH[2]

CRC Chromosome Molecular Biology Group, Department of Biochemistry, University of Oxford, South Parks Road, Oxford OX1 3QU, United Kingdom

HuRef k-mer profiles are useful in predicting array classification across ~400 male individuals into two distinct groups

Group 1

Group 2

Clustergram: K-mer Identity Matrix between Male Individuals

Group 1
🔵

Group 2
🔴

Clustergram: K-mer Identity Matrix between Male Individuals

For each feature → Iteratively run SVM (LIBSVM R package)

24-mer Library

Iteratively run SVM (LIBSVM R package) → Leave-one-out validation for each sample → Compute accuracy based on how many times each sample is correctly predicted

Sort features based on accuracy → Select features in top X percentile of accuracy

Group 1

Group 2

GATGTTTGCATTCACCTGACAGAG

Group 1
Group 2

Catalogue a new source of human sequence variation

Survey those k-mers that are enriched in one array group

MON 1

1          171 bp

GAAGATATTTCCTTTCTCACCTTA

Base Position 100 - 123

ENCODE Tier 1: Human Embryonic Stem Cell (H1-hESC)

MON 1

1                                            171 bp

GAAGATATTTCCTTTCTCACCTTA

Base Position 100 - 123

$$\frac{\text{Norm Freq H3K9me3 (IP)}}{\text{Norm Freq H3K9me3 (M)}} = \text{Relative Enrichment}$$

ENCODE Tier 1: Human Embryonic Stem Cell (H1-hESC)

# H1hESC Histone Profile of DYZ3 Array



Active Chromatin

- H3K4me1
- H3K4me2
- H3K4me3

Histone Variants

- H2Az

Repressive Chromatin

- H3K27me3
- H3K9me3

# H1hESC Histone Profile of DYZ3 Array



Active Chromatin
- H3K4me1
- H3K4me2
- H3K4me3

Histone Variants
- H2Az

Repressive Chromatin
- H3K27me3
- H3K9me3

# H1hESC Transcription Factor Enrichment Profile



**Transcription Factor**
- EZH2
- HDAC6
- PLU1
- JARID1A
- SUZ12

**Active Chromatin**
- H3K4me1
- H3K4me2
- H3K4me3

**Histone Variants**
- H2Az

**Repressive Chromatin**
- H3K27me3
- H3K9me3

# Adding Custom Datasets or "Tracks"



Transcription Factor
- EZH2
- HDAC6
- PLU1
- JARID1A
- SUZ12

**Active Chromatin**
- H3K4me1
- H3K4me2
- H3K4me3

**Histone Variants**
- H2Az

**Repressive Chromatin**
- H3K27me3
- H3K9me3

H3K9me3
[+ Chaetocin]

Non-polyA Transcript Mapping

Transcription Factor
- EZH2
- HDAC6
- PLU1
- JARID1A
- SUZ12

Active Chromatin
- H3K4me1
- H3K4me2
- H3K4me3

Histone Variants
- H2Az

Repressive Chromatin
- H3K27me3
- H3K9me3

H3K9me3
[+ Chaetocin]

# UCSC: Centromere Annotation and Tool Development

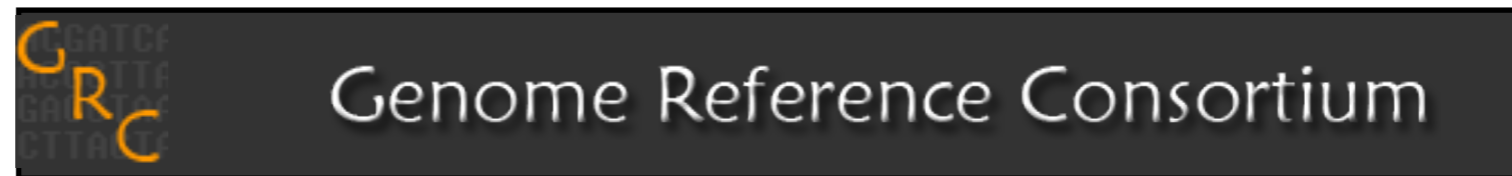# UCSC: Centromere Annotation and Tool Development

# Acknowledgements



**Jim Kent**
**David Haussler**
Max Hauesslar
Miten Jain
Yulia Newton
Dave Greenberg

**Hunt Willard**
Nick Altemose



Deanna Church
Valerie Schneider
Karen Clarke