

Methods for Exploring Large-Scale Genetic Datasets

Ryan M. Layer
ryan.layer@gmail.com
[@ryanlayer](https://twitter.com/ryanlayer)
University of Utah
quinlanlab.org

Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTCAGGGATCACTG



CGTGTCAATA

GTGCATTCGG

AGCGTGCATT

ATTCGGCGCG

GGCGCGGTCC

CGTCGAAATA

TCAATAGAC

CAACCTTCT ACACAGTTTC

CCGTCCAAGG

ATAGACAACC

ATTCTCGTC

AGCG

CGAAATACAC

CTGGGAGTGA

GAGTGATGTG

GCTACGGAAT

TACGGAATCA

GTTCCTGGGA

ATGTGGCTAC

CTG

AATCACTG

Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTCAGGGATCACTG



AGCGTGCATT GGCGCGGTCC TCAATAGAC ATTCTCGTC ACACAGTTT GAGTGATGTG TACGGAATCA
GTGCATT CGG CCGTCCAAGG ATAGACAACC CGTCGAAATA GTTCTGGGA ATGTGGCTAC AATCACTG
AGCG ATTGGCGCG AGTCAATA CAACCTTCT CGAAATACAC CTGGGAGTGA GCTACGGAAT CTG

Reference:



 G C G C A C C
G G C A C C
AGCGTGCATT GGCGCGGTCC TCAATAGAC ATTTCTCGTC ACACAGTTC GAGTGATGTG TACGGAAATCA
GTGCATTCCG CCGTCCAAGG ATAGACAACC CGTCGAAATA GTTCTGGGA ATGTGGCTAC AATCACTG
AGCG ATTGGCGCGC CGTGTCAATA CAACCTTCT CGAAATACAC CTGGGAGTGA GCTACGGAAT CTG

Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



G
G

C
G

G
C

C
A

C
C



Reference:

AGCGTGCATTGCCGGTCCAACGTGTCAATAGACAACATTCTCGTCAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



G
G

C
G

G
C

C
A

C
C



Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



G G C G

C A

C C

CGTGCATTG

CATTCTCGTCG

TGGGAGTGAT

ChIP-seq

ATTCGGCGCG

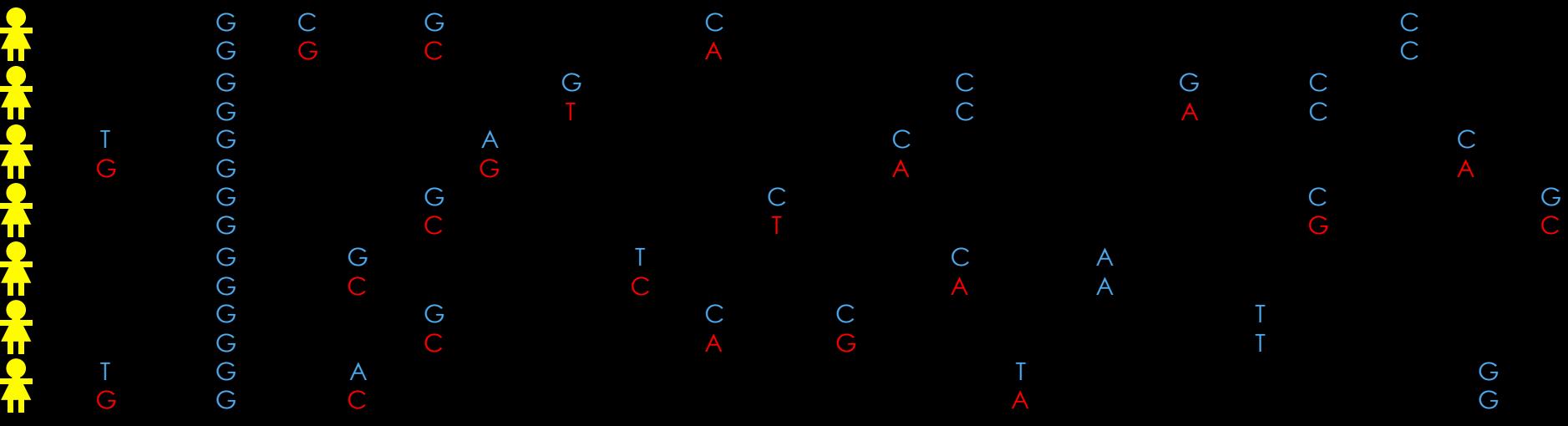
ACAACATTCTC

TTTCTGGGA

RNA-seq

ATTCGGCGCGACAACATTCTCTTCTGGGA
UAAGCCGCGCUGUUGUAAGAGAGUUUGACCU

Reference:



CGTGCGATTG

CATTCTCGTCG

TGGGAGTGAT

ChIP-seq

ATTCGGCGCG

ACAAACATTCTC

TTTCTGGGA

RNA-seq

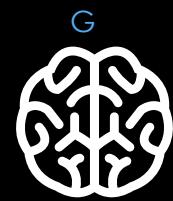
Reference:



CGTGCATTG
ATTCGGCGCG

CATTCTCGTCG
ACAAACATTCTC

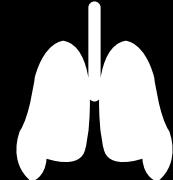
TGGGAGTGA
TTTCTGGGA



GCATTGCCCG
ATTGCCGGCG

ACAAACATTCTC

TTCTGGGA GATGTGGCT



GCATTGCGCCG
ATTGCCGGCGG

ATTCTCGTCGAAATA

TTTCTGGGA



CAATAGACAAACATT TCGAAATACACA



Reference:

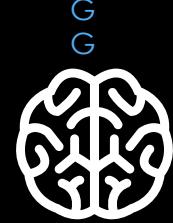


T
G

CGTGCATTG
ATTGGCGCG

CATTCTCGTCG
ACAAACATTTCTC

TGGGAGTGAT
TTTCTGGGA

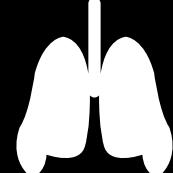


GC ATTGGCCG
ATTGGCCGCG

ACAAACATTCTC

TTTCTGGGA

GATGTGGC'



GOATTGCCG
ATTGCCGGCG

ATTCTCGTCGAAATA

TTTCTGGGA



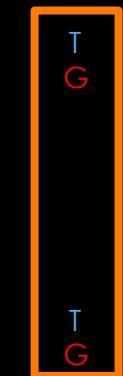
CAATAGACAAACATT

TCGAAATACACA



Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAACTACTG



CGTGCATTCG
ATTGGCGCG

GOATTGCCCG
ATTGCCCGCG

GOATTGCCCG
ATTGCCCGCGG

CAATAGACAACATT

ATTCTCGTCGAAATA

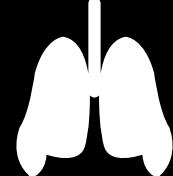
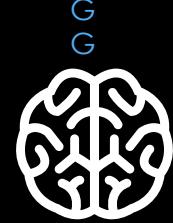
ACAACATTCTC

TCGAAATACACA

CATTCTCGTCG
ACAACATTCTC

TGGGAGTGAT
TTTCTGGGA

TTTCTGGGA
GATGTGGCT



Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



CGTGCATTG
ATTGGCGCG

GOATTGCCG
ATTGCCCGCG

GOATTGCCG
ATTGCCCGCGG

GENOTYPE QUERY TOOLS



C C

G A

C C

C A

G C

CATTCTCGTCG
ACAACATTCTC

ACAACATTCTC

ATTCTCGTCGAAATA

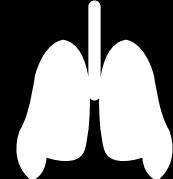
CAATAGACAACATT

TGGGAGTGAT
TTTCTGGGA

TTTCTGGGA

GATGTGGCT

TTTCTGGGA



Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



CGTGCATTG
ATTGGCGCG

GOATTGCCCG
ATTGCCCGCG

GOATTGCCCG
ATTGCCCGCGG

GENOTYPE QUERY TOOLS



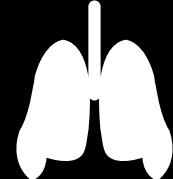
QIGGLE

CAATAGACAACATT

TCGAAATACACA

TGGGAGTGAT
TTTCTGGGA

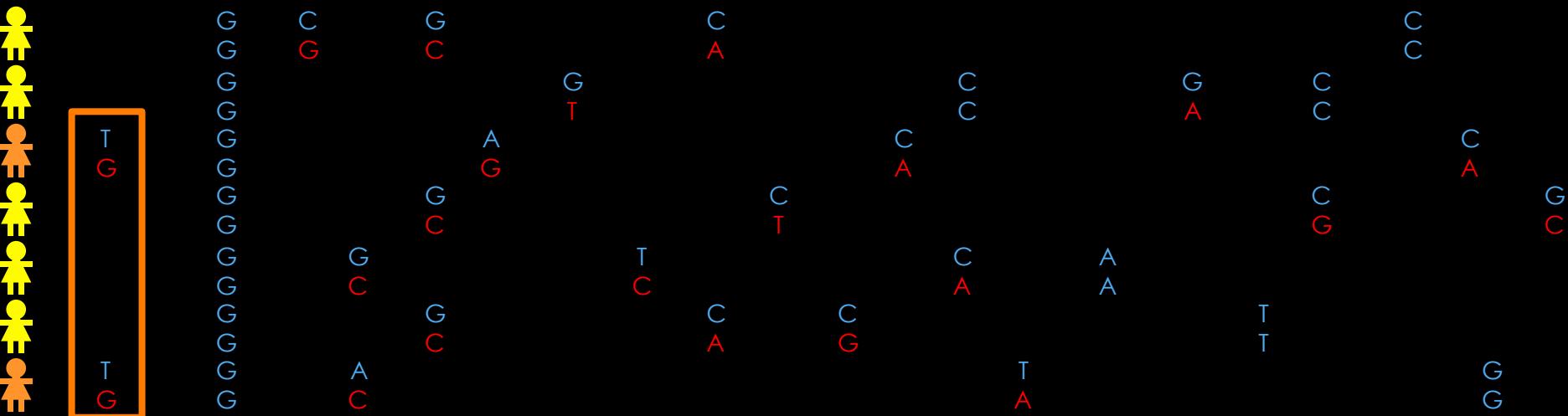
GATGTGGCT
TTTCTGGGA





GENOTYPE QUERY TOOLS

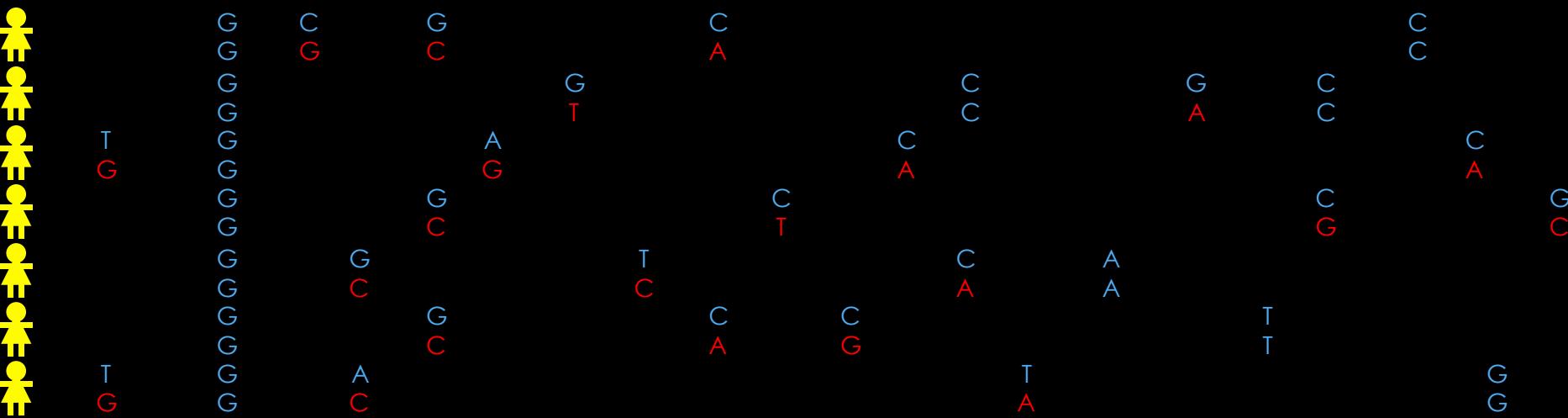
github.com/ryanlayer/gqt



Efficient genotype compression and analysis of large genetic-variation data sets
RM Layer, N Kindlon, Exome Aggregation Consortium, KJ Karczewski, AR Quinlan.
Nature methods, 2016

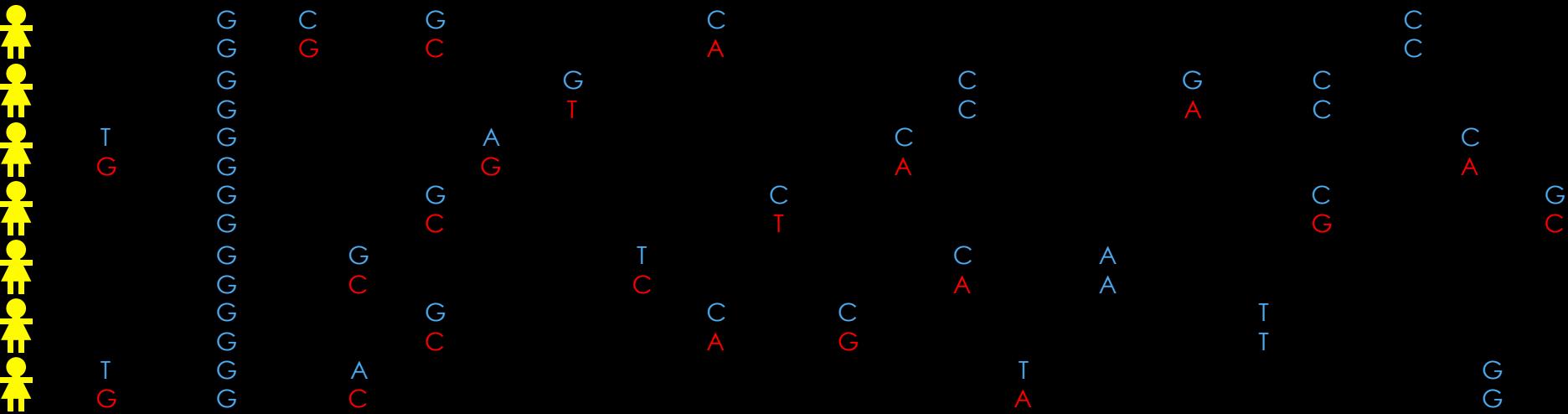
Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAACTACTG



Reference:

AGCGTGCATCGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



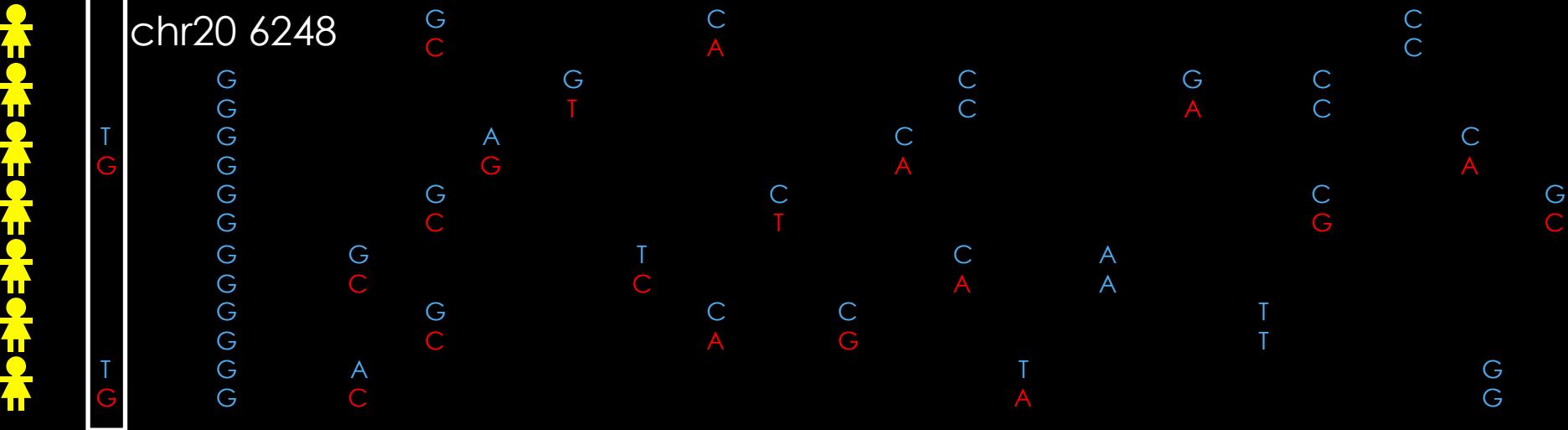
Variant Call Format (VCF)

Variant Data

Individual Genotypes

Reference:

AGCGTGCATCGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



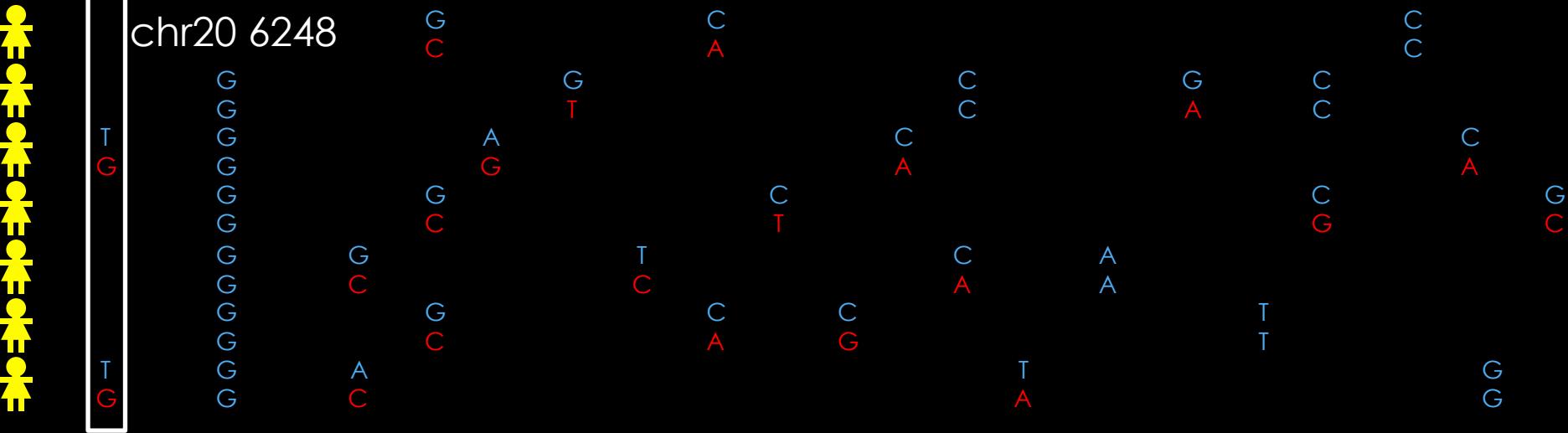
Variant Call Format (VCF)

Variant Data

Individual Genotypes

Reference:

AGCGTGCATCGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



Variant Call Format (VCF)

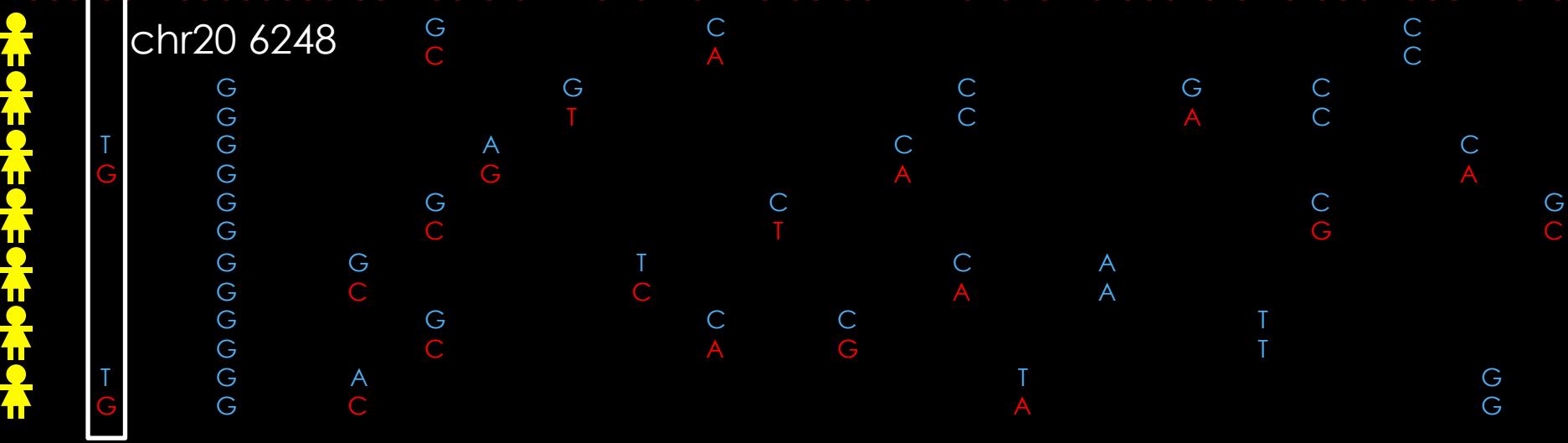
Variant Data

Individual Genotypes

chr20 6248

Reference:

AGCGG[GATTGCCGGTCCAACTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGA[TGGCTAGGGAAATCACTG



Variant Call Format (VCF)

Variant Data

Individual Genotypes

chr20 6248

G T

Reference:

Variant Call Format (VCF)

Variant Data

Individual Genotypes

chr20

6248

T

0 / 0

0

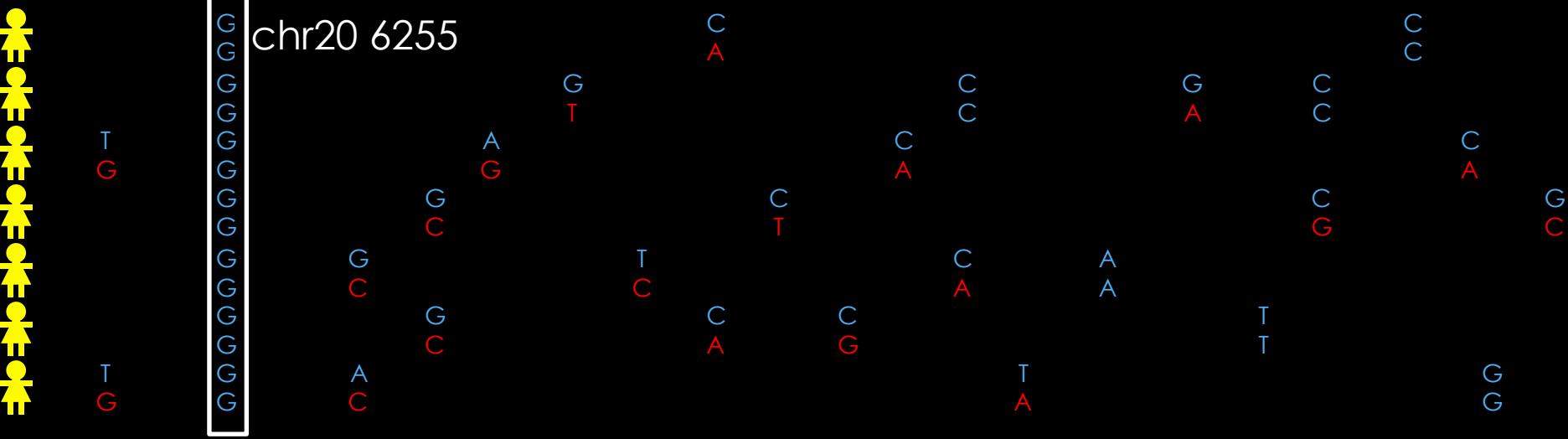
0/1

/0

1 /

Reference:

AGCGTGCATTCCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



Variant Call Format (VCF)

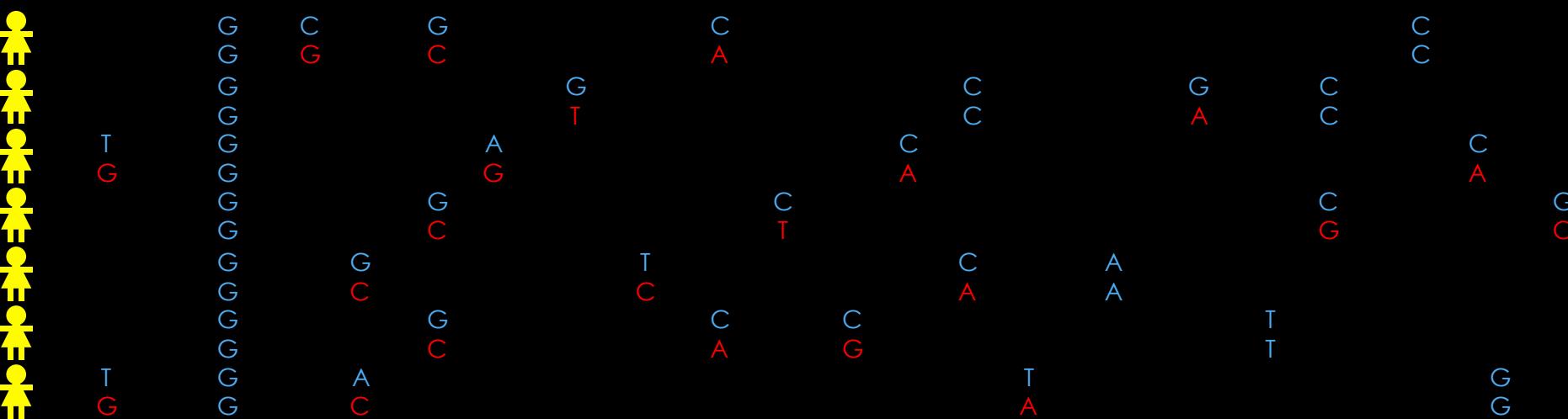
Variant Data

Individual Genotypes

chr20	6248	G	T	0/0	0/0	1/0	0/0	0/0	0/0	1/0
chr20	6255	C	G	1/1	1/1	1/1	1/1	1/1	1/1	1/1

Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



Variant Call Format (VCF)

Variant Data

chr20	6248	G	T	0/0	0/0	1/0	0/0	0/0	0/0	1/0
chr20	6255	C	G	1/1	1/1	1/1	1/1	1/1	1/1	1/1
chr20	6259	G	C	1/0	0/0	0/0	0/0	0/0	0/0	0/0
chr20	6262	C	G	0/0	0/0	0/0	0/0	1/0	0/0	1/0
chr20	6266	C	G	1/0	0/0	0/0	1/0	0/0	1/0	0/0
chr20	6269	G	A	0/0	0/0	1/0	0/0	0/0	0/0	0/0
chr20	6274	T	G	0/0	1/0	0/0	0/0	0/0	0/0	0/0
chr20	6278	C	T	0/0	0/0	0/0	0/0	1/0	0/0	0/0
chr20	6282	A	C	1/0	0/0	0/0	0/0	1/0	1/0	0/0
chr20	6287	T	C	0/0	0/0	0/0	1/0	0/0	0/0	0/0
chr20	6291	G	C	0/0	0/0	0/0	0/0	0/0	1/0	0/0
chr20	6294	A	C	0/0	0/0	0/0	0/0	1/1	0/0	0/0
chr20	6298	A	G	0/0	0/0	1/0	0/0	0/0	0/0	0/0
chr20	6301	T	A	0/0	1/1	0/0	0/0	1/0	0/0	0/0
chr20	6306	A	G	1/0	0/0	0/0	0/0	1/0	1/0	0/0
chr20	6311	G	C	0/0	1/1	0/0	1/0	0/0	0/0	0/0

Individual Genotypes

1,000,000 Individuals

100,000,000 Variants

1.0e14 genotypes

1 bit/genotype =
12.5 Terabytes

Text

char = 1 byte (8 bits)

0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/1 1/1

24 bit/genotype
300 TB

Binary(PLINK)

00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 01 10

21

2 bit/genotype
25 TB

Compressed

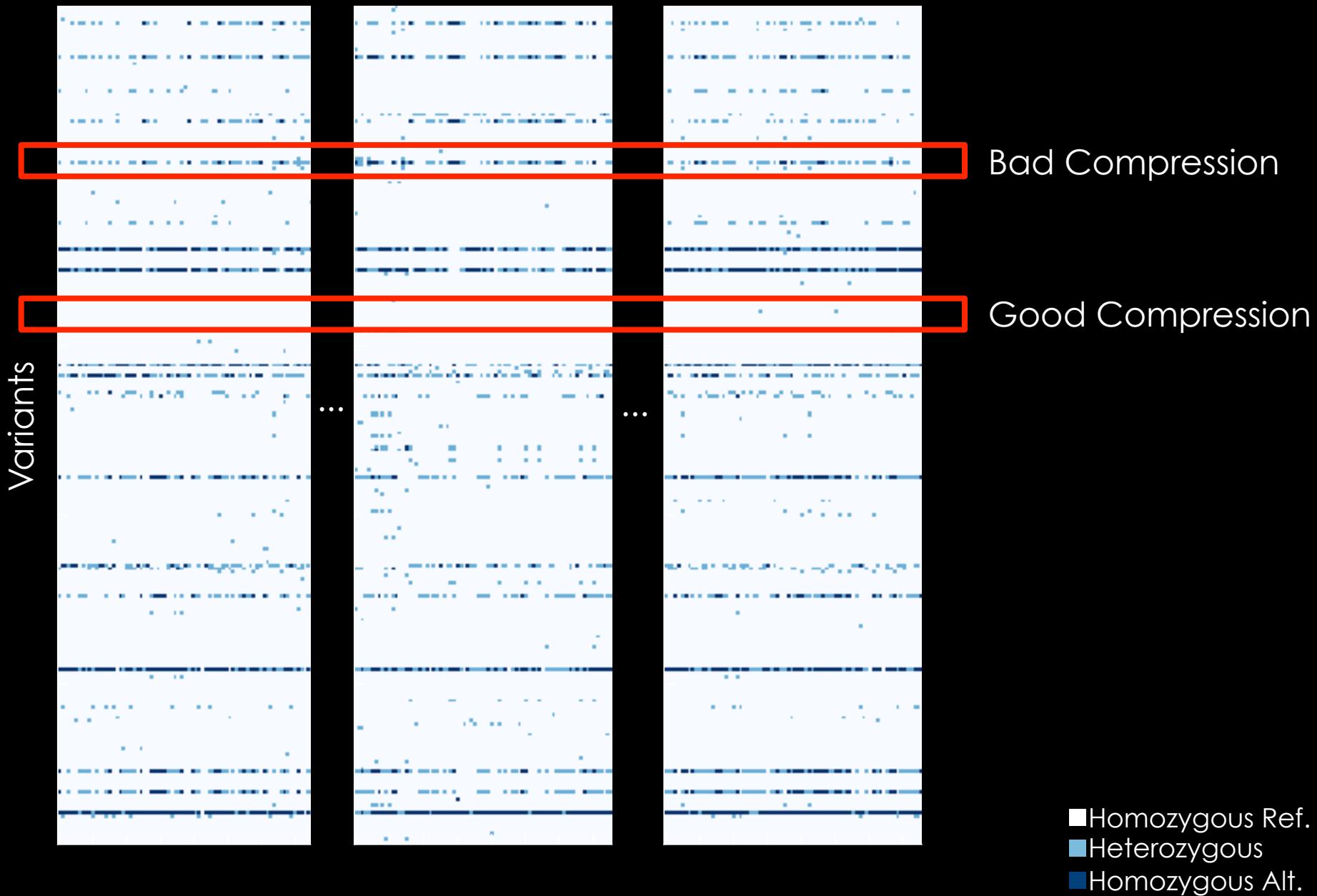
Run-length Encoding
value : length

00 00 00 00 00 00 00 00 00 00 00 00 0 → 0 : 21

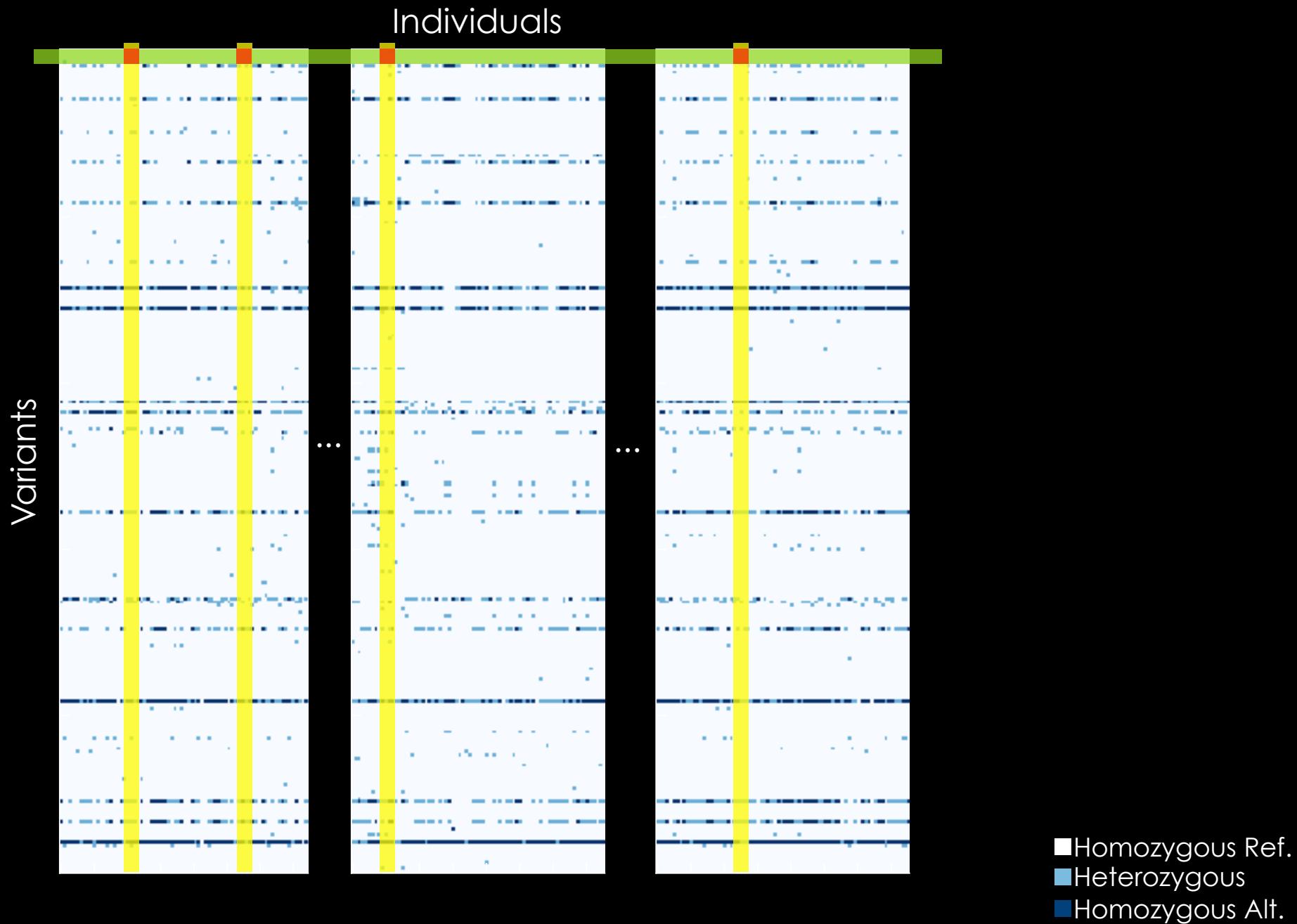
00010101

1000 Genomes chr20

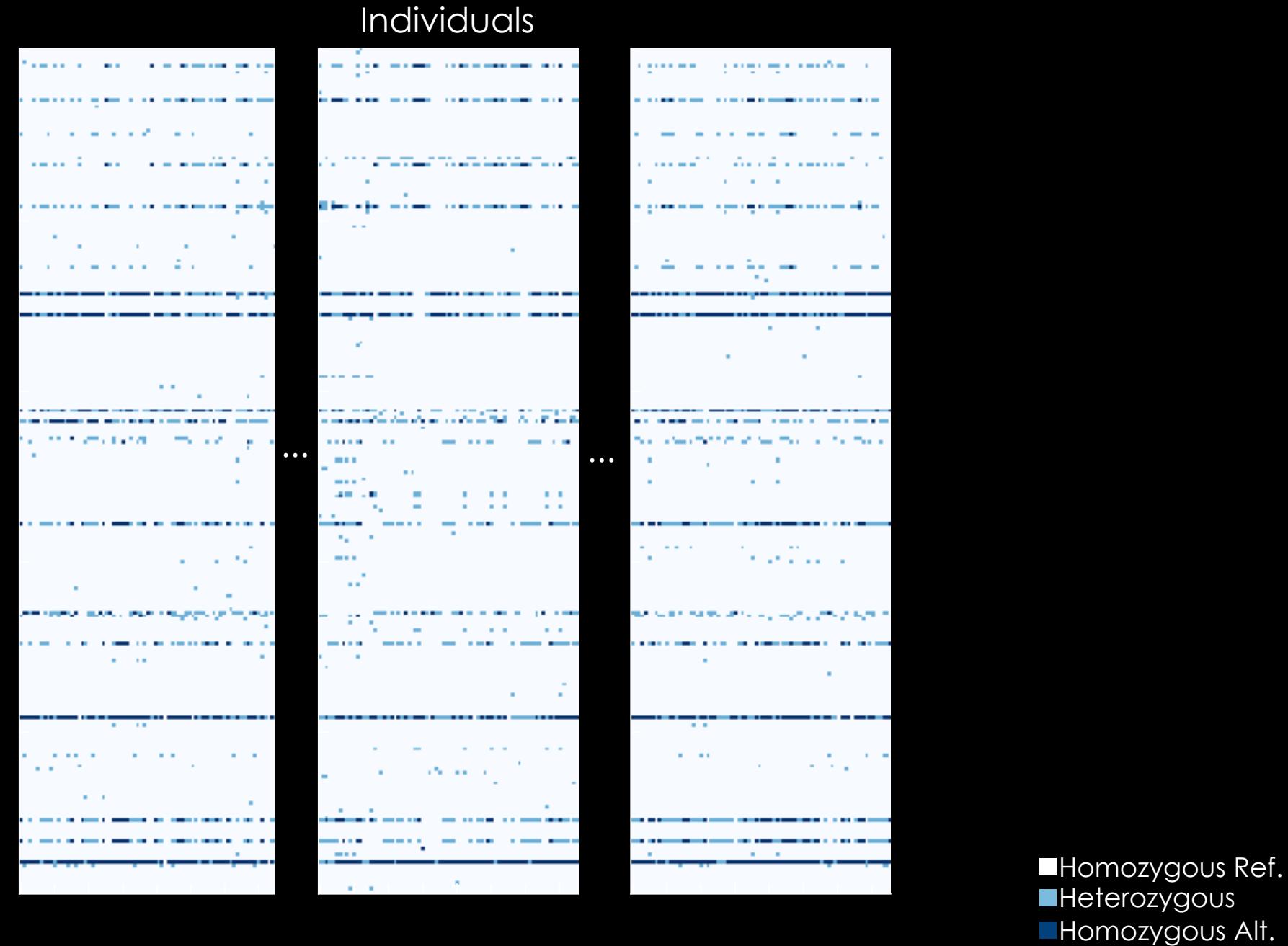
Individuals



Query: Find the variants that are heterozygous in all affected individuals

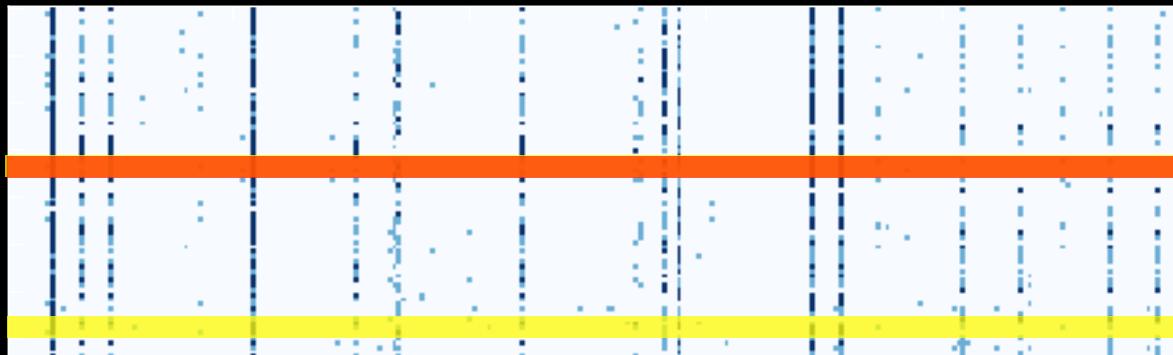


Query: Find the variants that are heterozygous in all affected individuals



Query: Find the variants that are heterozygous in all affected individuals

Variants



...



...

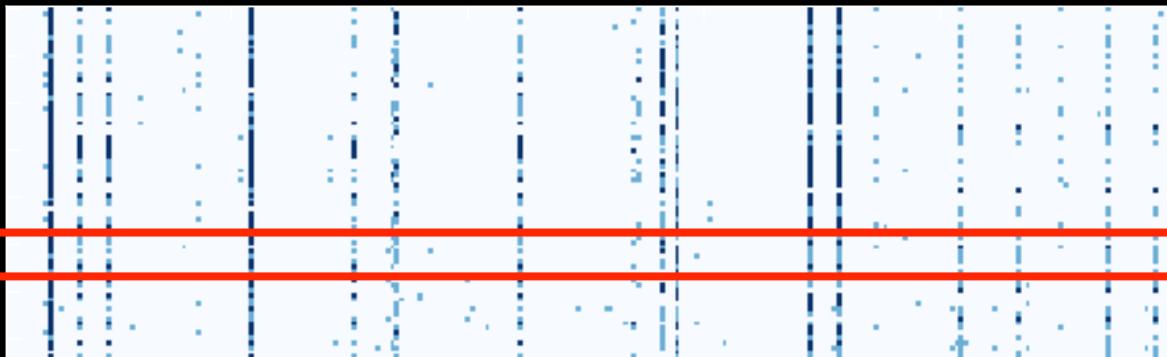


Individuals

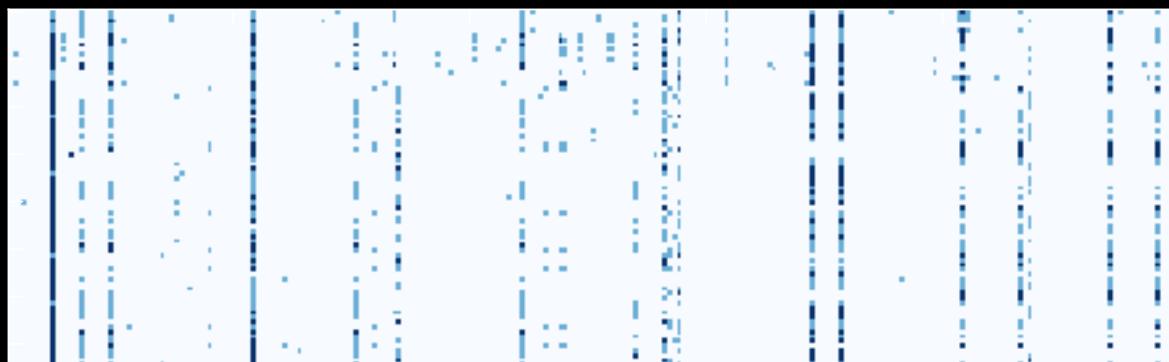
- Homozygous Ref.
- Heterozygous
- Homozygous Alt.

Query: Find the variants that are heterozygous in all affected individuals

Variants



Bad Compression

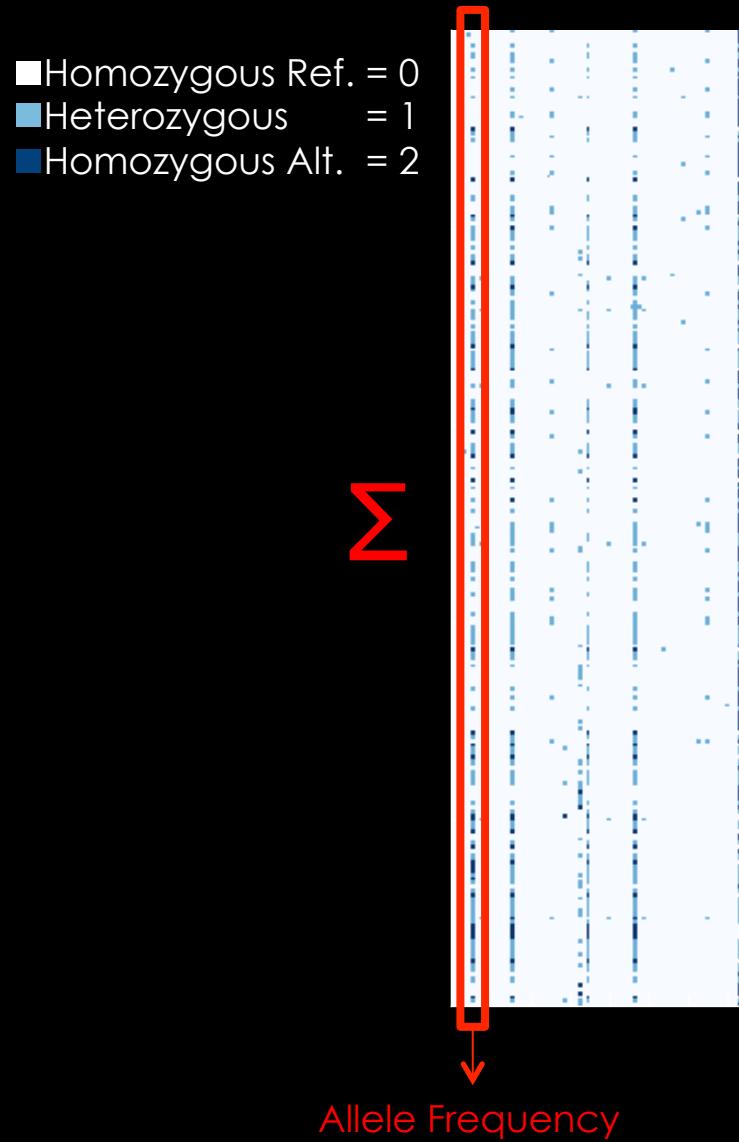


Individuals

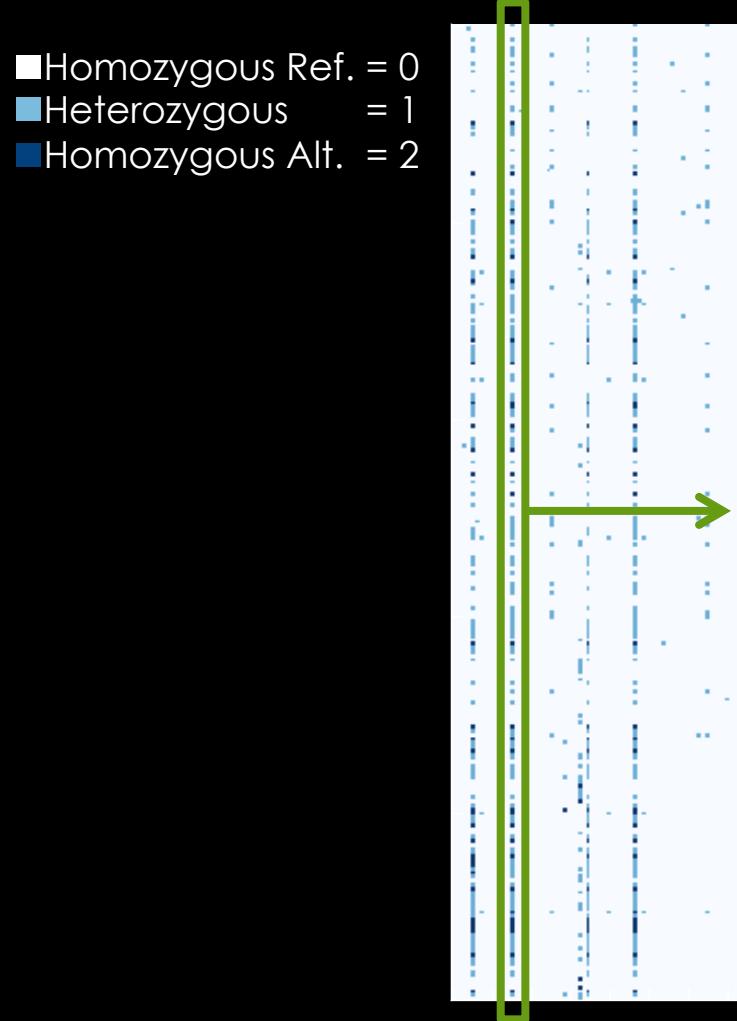


- Homozygous Ref.
- Heterozygous
- Homozygous Alt.

Sort Variants by Allele Frequency

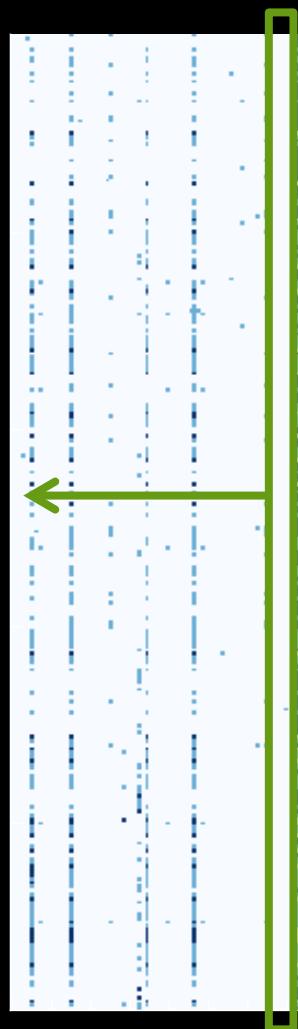


Sort Variants by Allele Frequency



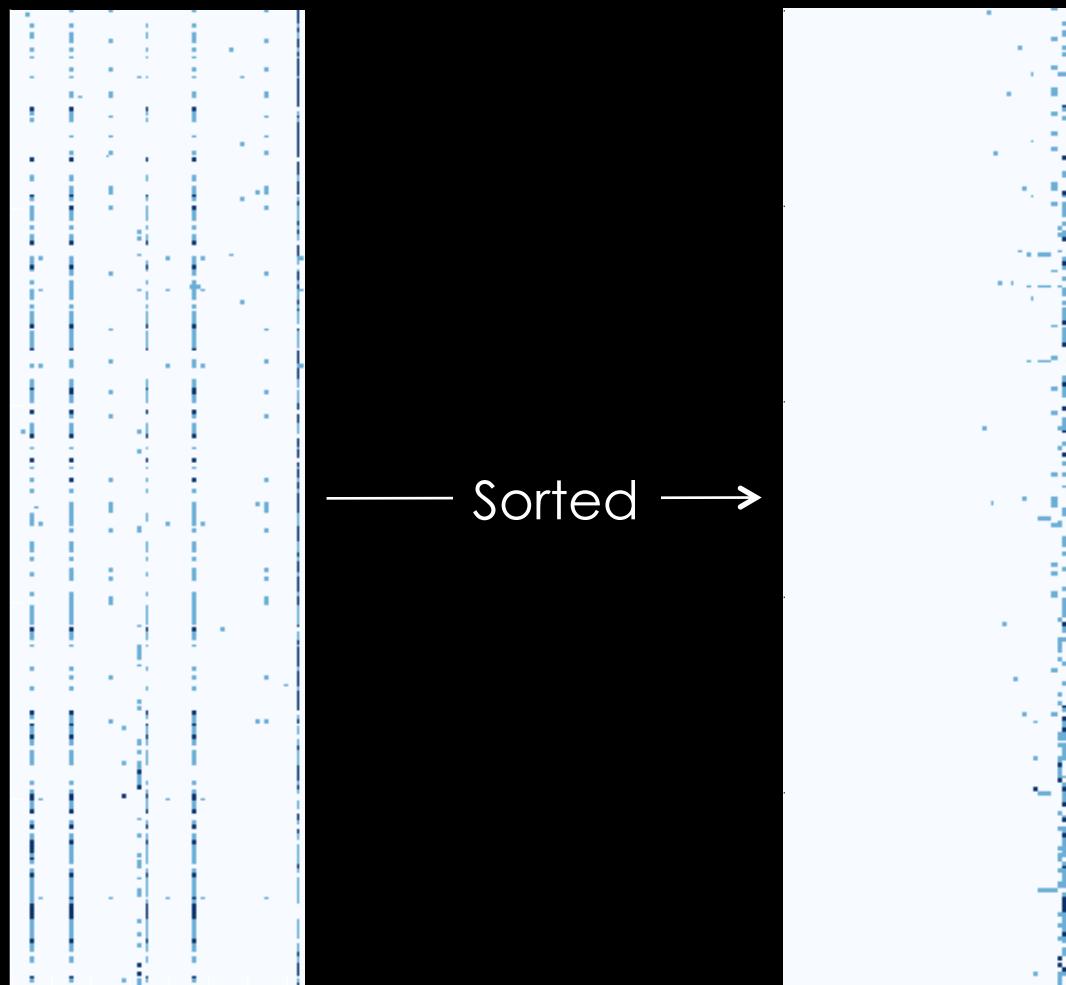
Sort Variants by Allele Frequency

- Homozygous Ref. = 0
- Heterozygous = 1
- Homozygous Alt. = 2



Sort Variants by Allele Frequency

- Homozygous Ref. = 0
- Heterozygous = 1
- Homozygous Alt. = 2



Query: Find the variants that are heterozygous in individuals 1, 2, and 3

1	1 1	0 0	0 0	0 0	0 0	0 0	0 0	0 0	0 0	0 1	0 0	0 0	0 1	0 0	0 0	1 1	0 0
2	0 0	0 0	0 0	1 1	0 1	0 0	0 0	0 0	0 0	0 1	0 0	0 0	0 1	0 0	0 0	1 1	0 0
3	0 1	0 0	0 0	0 0	0 0	0 0	0 0	0 1	0 0	0 1	0 0	0 0	1 1	0 0	0 0	0 0	0 0

Bitmap Index: map records to bit arrays then answer queries using bitwise logical operations

1 1|1 0|0 0|0 0|0 0|0 0|0 0|0 0|0 0|1 0|0 0|0 0|1 0|0 0|0 1|1 0|0

Bitmap Index: map records to bit arrays then answer queries using bitwise logical operations

	1	1 1	0 0	0 0	0 0	0 0	0 0	0 0	0 0	0 1	0 0	0 0	0 1	0 0	0 0	1 1	0 0
Homo Ref	0	1	1	1	1	1	1	1	1	0	1	1	0	1	1	0	1

Bitmap Index: map records to bit arrays then answer queries using bitwise logical operations

1	1 1	0 0	0 0	0 0	0 0	0 0	0 0	0 0	0 0	0 1	0 0	0 0	0 1	0 0	0 0	1 1	0 0
Homo Ref	0	1	1	1	1	1	1	1	0	1	1	0	0	1	1	0	1
Het	0	0	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0

Bitmap Index: map records to bit arrays then answer queries using bitwise logical operations

Bitmap Index: map records to bit arrays then answer queries using bitwise logical operations

Bitmap Index: map records to bit arrays then answer queries using bitwise logical operations

Query: Find the variants that are heterozygous in individuals 1, 2, and 3

Query: Find the variants that are heterozygous in individuals 1, 2, and 3

	Homo Ref	0	1	1	1	1	1	1	1	0	1	1	0	1	1	0	1
1	Het	0	0	0	0	0	0	0	0	1	0	0	1	0	0	0	0
	Homo Alt	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
	Unknown	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	Homo Ref	1	1	1	0	0	1	1	1	0	1	1	0	1	1	0	1
2	Het	0	0	0	0	1	0	0	0	1	0	0	1	0	0	0	0
	Homo Alt	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0
	Unknown	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	Homo Ref	0	1	1	1	1	1	1	0	1	1	1	0	1	1	0	1
3	Het	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
	Homo Alt	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
	Unknown	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	AND	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0

Genotype Filters

```
"HET"  
"HOM_REF"  
"HOM_ALT"
```

Summary Statistics

```
"count(HET)"  
"pct(HET HOM_ALT)"  
"maf()"
```

Statistic Filters

```
"pct(HET HOM_ALT) <= 0.01"  
"count(HOM_REF) > 1"  
"maf() < 0.1"
```

VCF

1|1 0|0 0|0 0|0 0|0 0|0 0|0 0|0 0|0 0|1 0|0

24 bits/genotype

Binary

10000000 00000000 01000001 00001000

2 bits/genotype

Bit Map

01111111	01101101
00000000	10010000
10000000	00000010
00000000	00000000

4 bits/genotype

Run-Length Encoding Complicates* Logical Operations

	00000000	00000111	11111111	11110000
OR	11111111	10000000	00000000	00000000
<hr/>				
	11111111	10000111	11111111	11111111

	13 bits		
	↓		
	00001101	10001111	00000100
OR	10001001	10010111	
<hr/>			
	9 bits		

Succinct Data Structure

- Near-optimal compression
- Allow operations without inflation

*likely, but not proven to be impractical

Word Aligned Hybrid [Wu, TODS 2006]

Run Length Encoding

10010000 → value : length (bits)

1 : 1	0 : 2	1 : 1	0 : 4
10000001	00000010	10000001	00000100

Word Aligned Hybrid

Literal → type : uncompressed value

1001000 → 0101000

Fill → type : value : length (words)

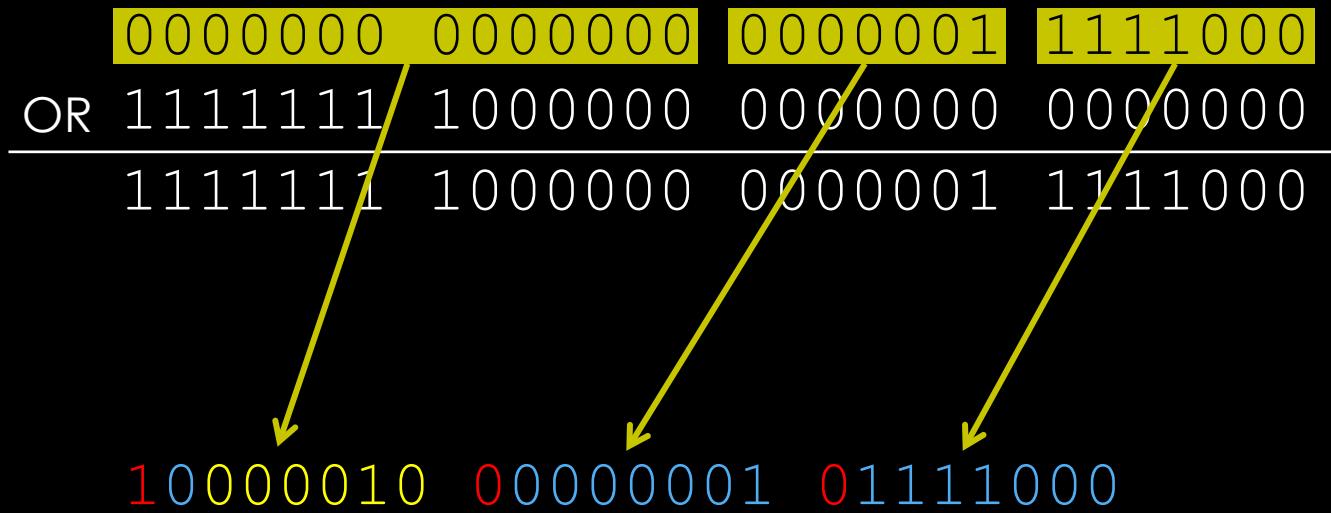
0000000 0000000 → 1000010

1111111 1111111 1111111 → 1100011

Word Aligned Hybrid [Wu, TODS 2006]

	0000000	0000000	0000001	1111000
OR	1111111	1000000	0000000	0000000
<hr/>				
	1111111	1000000	0000001	1111000

Word Aligned Hybrid [Wu, TODS 2006]



Word Aligned Hybrid [Wu, TODS 2006]



Word Aligned Hybrid [Wu, TODS 2006]

	0000000	0000000	0000001	1111000
OR	1111111	1000000	0000000	0000000
<hr/>				
	1111111	1000000	0000001	1111000

10000010	00000001	01111000
01111111	01000000	10000010

Word Aligned Hybrid [Wu, TODS 2006]

	0000000	0000000	0000001	1111000
OR	1111111	1000000	0000000	0000000
<hr/>				
	1111111	1000000	0000001	1111000

2 words
↓

10000010	00000001	01111000
01111111	01000000	10000010

↑
1 word ↑
1 word

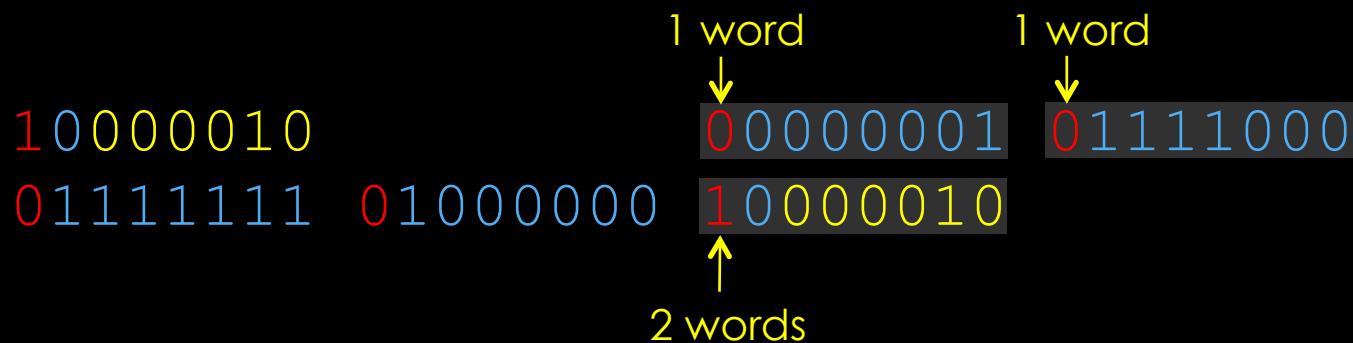
Word Aligned Hybrid [Wu, TODS 2006]

	0000000	0000000	0000001	1111000
OR	1111111	1000000	0000000	0000000
<hr/>				
	1111111	1000000	0000001	1111000

10000010 → 0000001 0111100
01111111 01000000 10000010

Word Aligned Hybrid [Wu, TODS 2006]

	0000000	0000000	0000001	1111000
OR	1111111	1000000	0000000	0000000
<hr/>				
	1111111	1000000	0000001	1111000



Word Aligned Hybrid [Wu, TODS 2006]

	0000000	0000000	0000001	1111000
OR	1111111	1000000	0000000	0000000
<hr/>				
	1111111	1000000	0000001	1111000

	10000010		00000001	01111000
OR	01111111	01000000	10000010	
<hr/>				

Word Aligned Hybrid [Wu, TODS 2006]

	0000000	0000000	0000001	1111000
OR	1111111	1000000	0000000	0000000
<hr/>				
	1111111	1000000	0000001	1111000

	10000010		00000001	01111000
OR	01111111	01000000	10000010	
<hr/>				
	01111111	01000000	00000001	01111000

Word Aligned Hybrid [Wu, TODS 2006]

Smaller Files:

Compresses with run-length & literal values

Faster Queries:

Bitwise logical operations without inflation

GQT v BCFTOOLS & PLINK v1.9

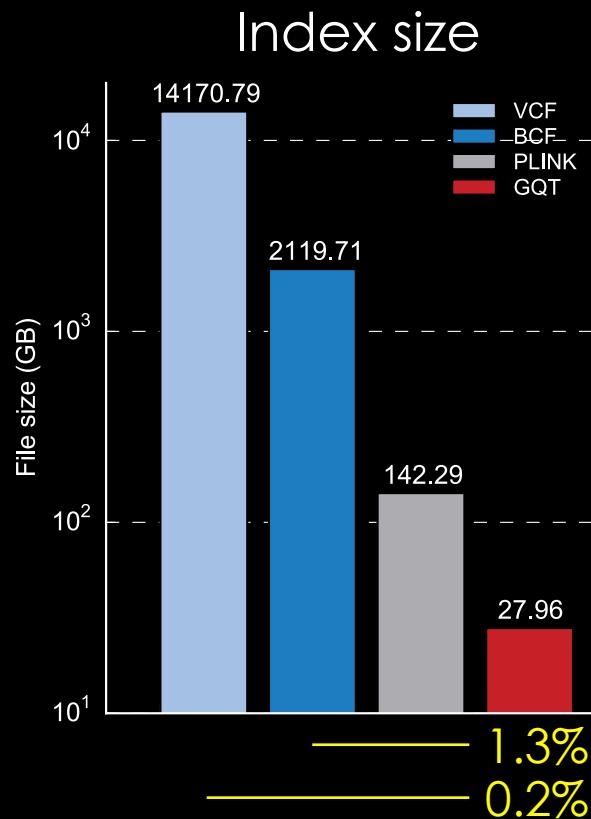
The Exome Aggregation Consortium (ExAC)

60K individual 9.3M variants

Linux v3.13, gcc v4.8.2, 4.0 GHz Intel Core i7 (Haswell)

GQT v BCFTOOLS & PLINK v1.9

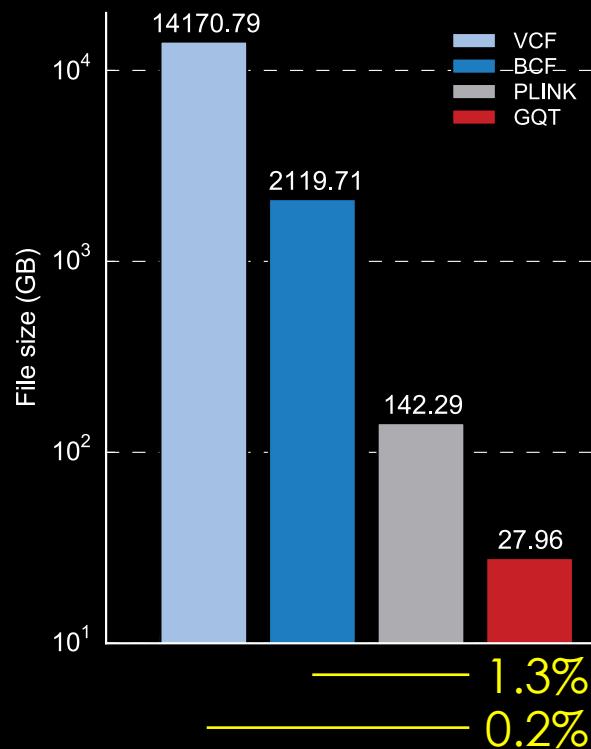
The Exome Aggregation Consortium (ExAC)
60K individual 9.3M variants
Linux v3.13, gcc v4.8.2, 4.0 GHz Intel Core i7 (Haswell)



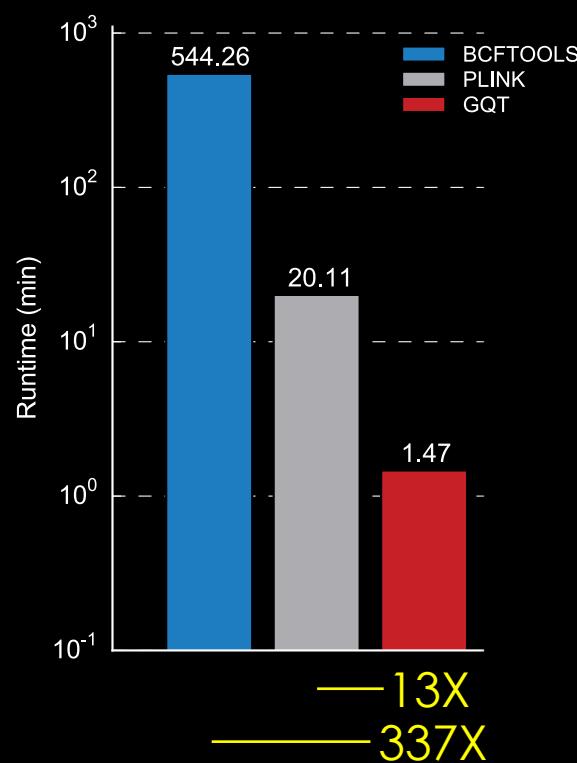
GQT v BCFTOOLS & PLINK v1.9

The Exome Aggregation Consortium (ExAC)
60K individual 9.3M variants
Linux v3.13, gcc v4.8.2, 4.0 GHz Intel Core i7 (Haswell)

Index size



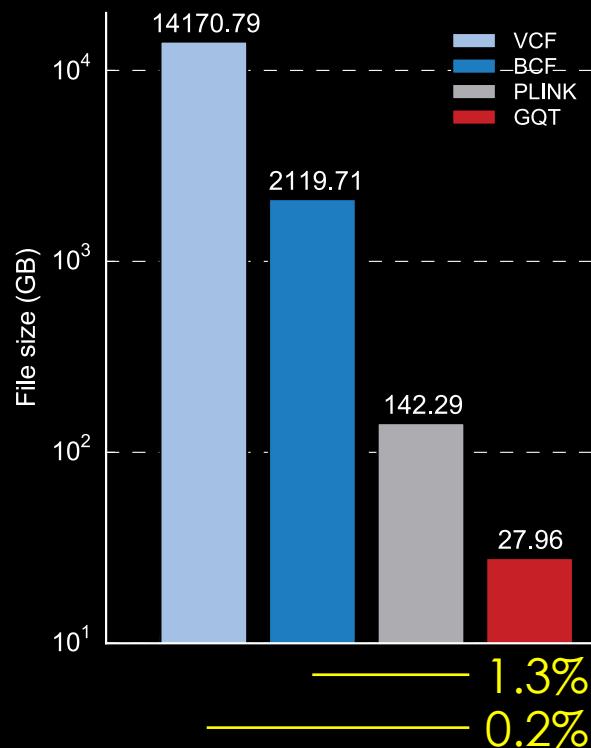
Alt. Allele Count



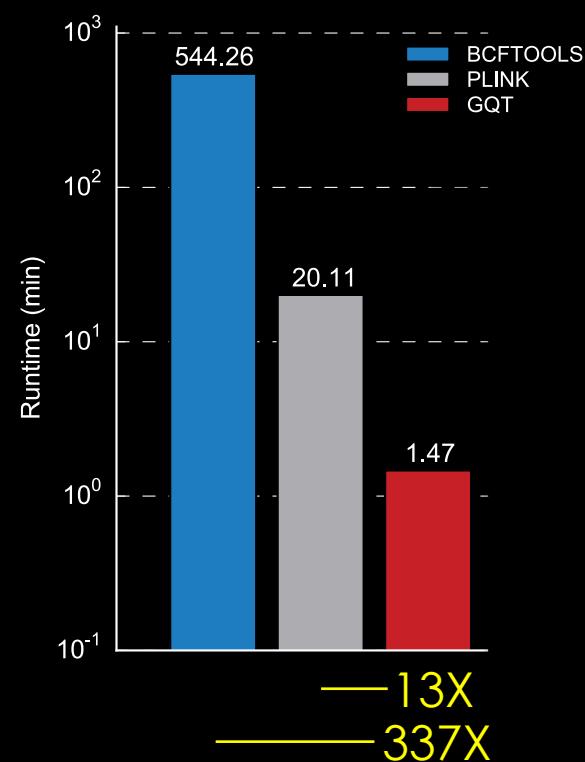
GQT v BCFTOOLS & PLINK v1.9

The Exome Aggregation Consortium (ExAC)
60K individual 9.3M variants
Linux v3.13, gcc v4.8.2, 4.0 GHz Intel Core i7 (Haswell)

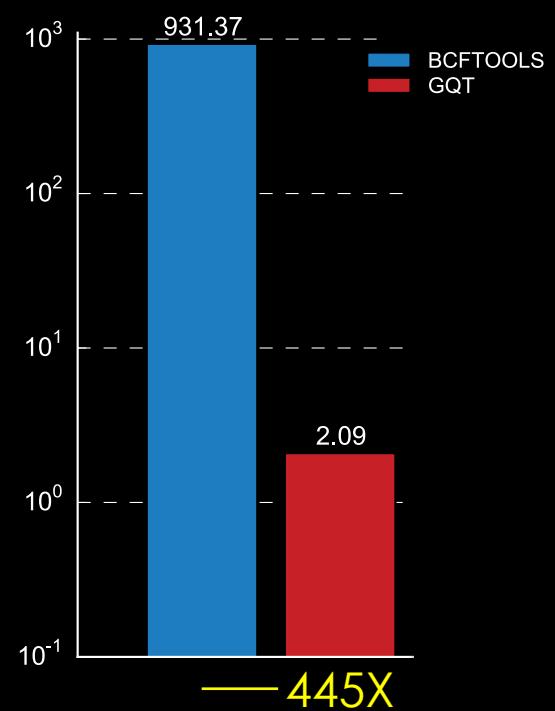
Index size



Alt. Allele Count

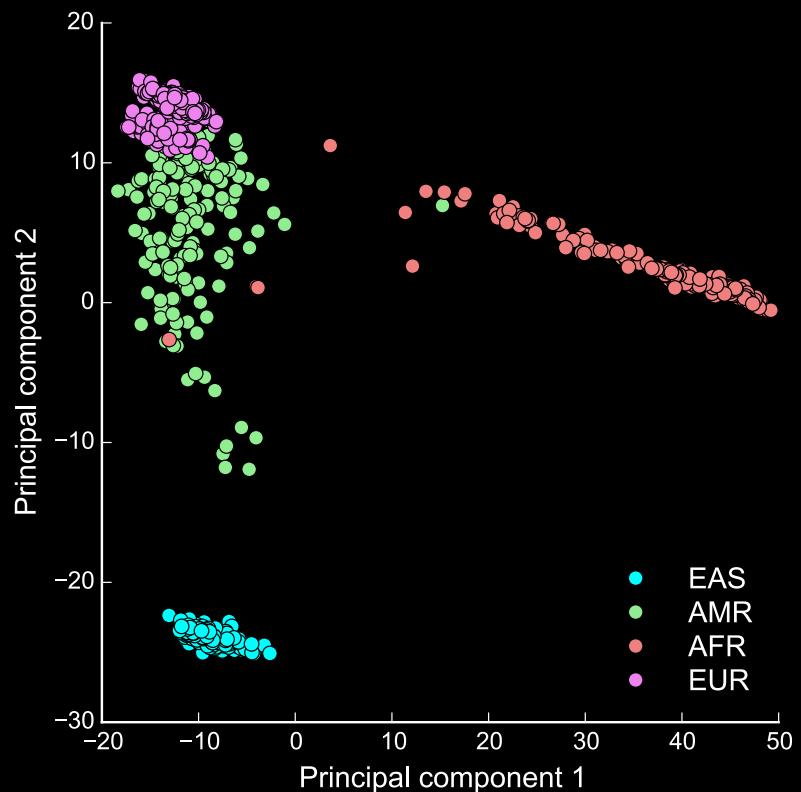
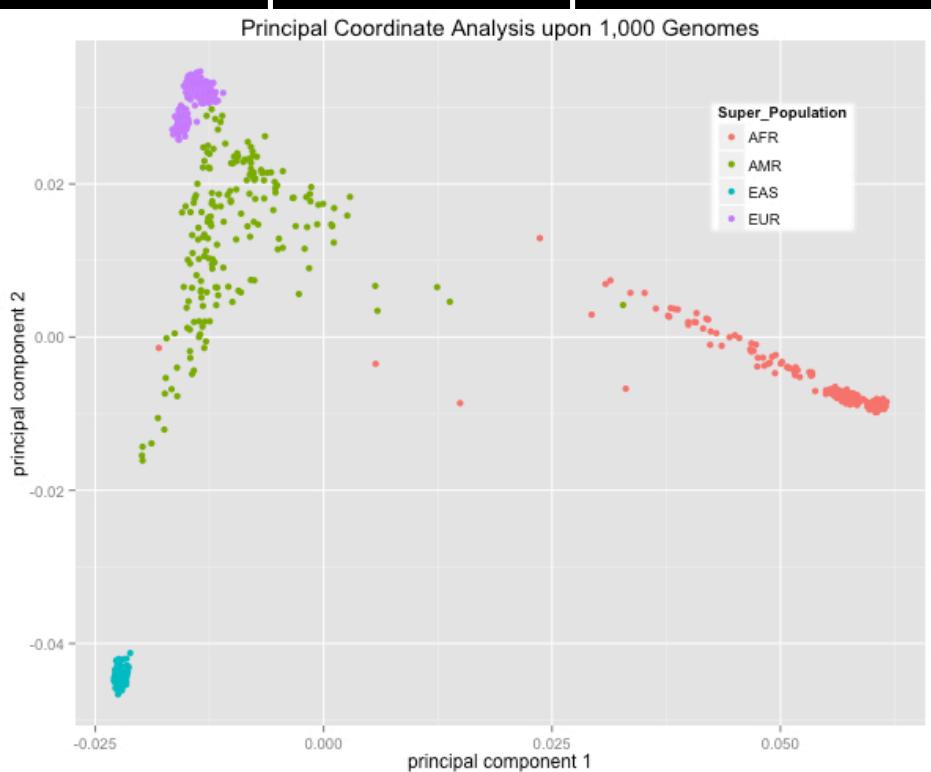


Rare Variant Search



Google Genomics & Apache Spark

GQT



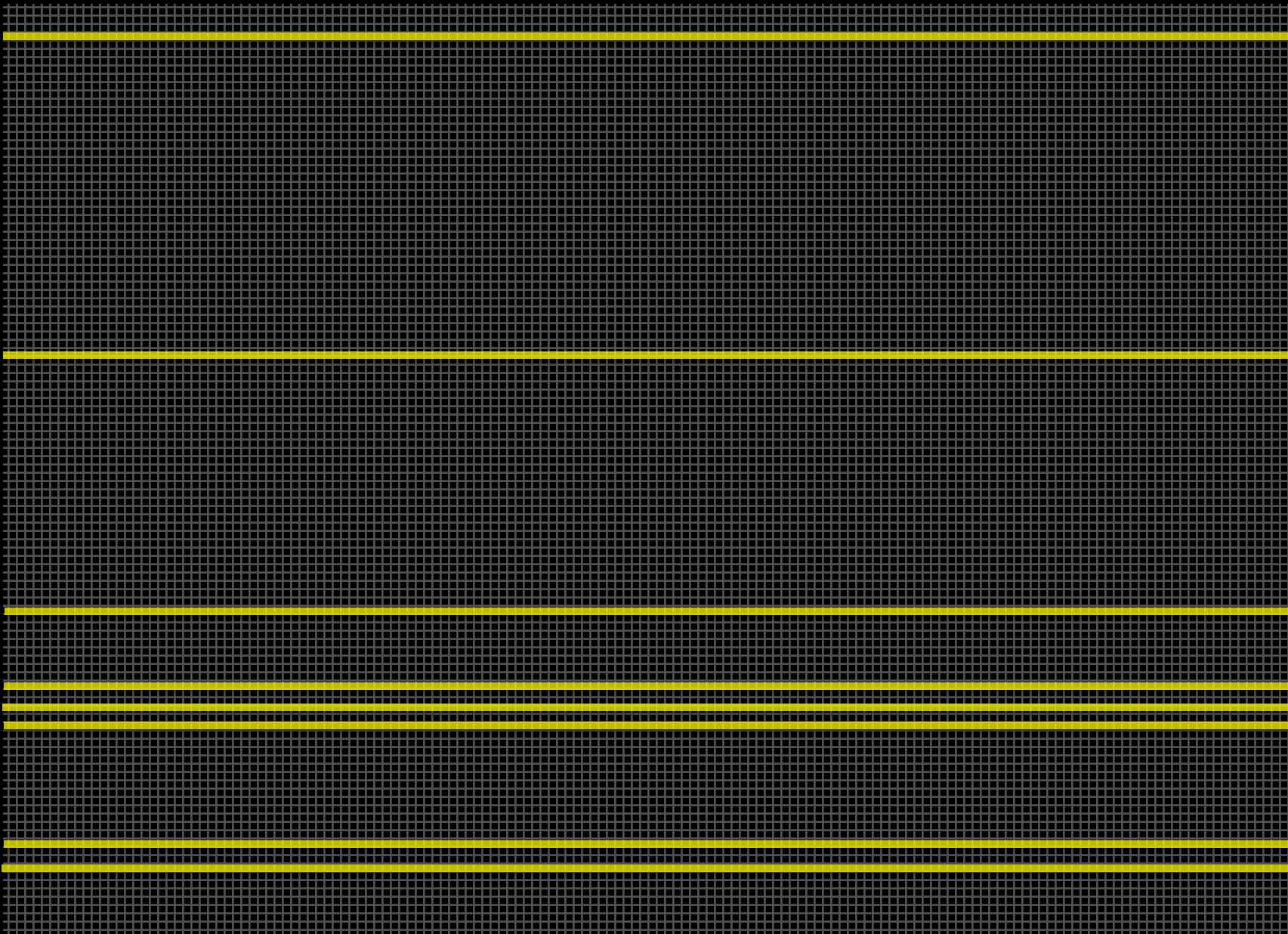
"... about 2 hours
on 60 eight core machines."¹
 $2 \times 60 \times 8 = 960$ compute hours

32 min 31 sec
on 1 core
0.542 compute hours

~1,800X speedup

100,000,000 Variants

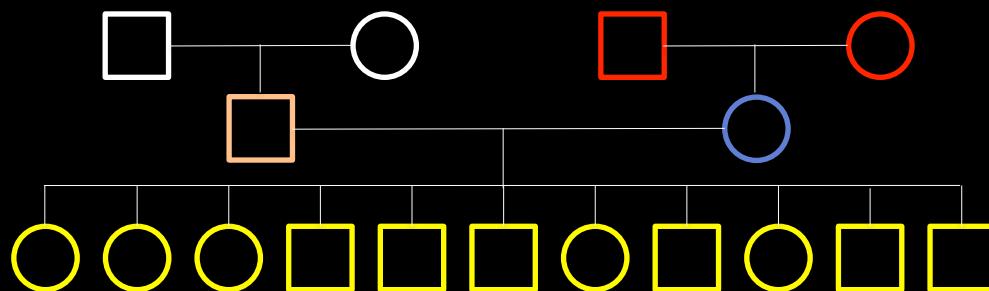
1,000,000 Individuals



Pedigree File (PED)

<u>Family ID</u>	<u>Individual ID</u>	<u>Paternal ID</u>	<u>Maternal ID</u>	<u>Sex</u>	<u>Population</u>	<u>Phenotype</u>
1463	NA12889	0	0	1	CEU	0
1463	NA12890	0	0	2	CEU	0
1463	NA12891	0	0	1	CEU	0
1463	NA12892	0	0	2	CEU	0
1463	NA12877	NA12889	NA12890	1	CEU	0
1463	NA12878	NA12891	NA12892	2	CEU	0
1463	NA12879	NA12877	NA12878	2	CEU	0
1463	NA12880	NA12877	NA12878	2	CEU	0
1463	NA12881	NA12877	NA12878	2	CEU	0
1463	NA12882	NA12877	NA12878	1	CEU	0
1463	NA12883	NA12877	NA12878	1	CEU	0
1463	NA12884	NA12877	NA12878	1	CEU	0
1463	NA12885	NA12877	NA12878	2	CEU	0
1463	NA12886	NA12877	NA12878	1	CEU	0
1463	NA12887	NA12877	NA12878	2	CEU	0
1463	NA12888	NA12877	NA12878	1	CEU	0
1463	NA12893	NA12877	NA12878	1	CEU	0

CEPH/Utah Pedigree 1463



Phenotype Query

Parents

```
"Individual_ID IN ('NA12878', 'NA12877')"
```

Female children

```
"Sex = 2 AND Maternal_ID = 'NA12878'"
```

Affected across the CEU population

```
"Phenotype = 2 AND Population ='CEU'"
```

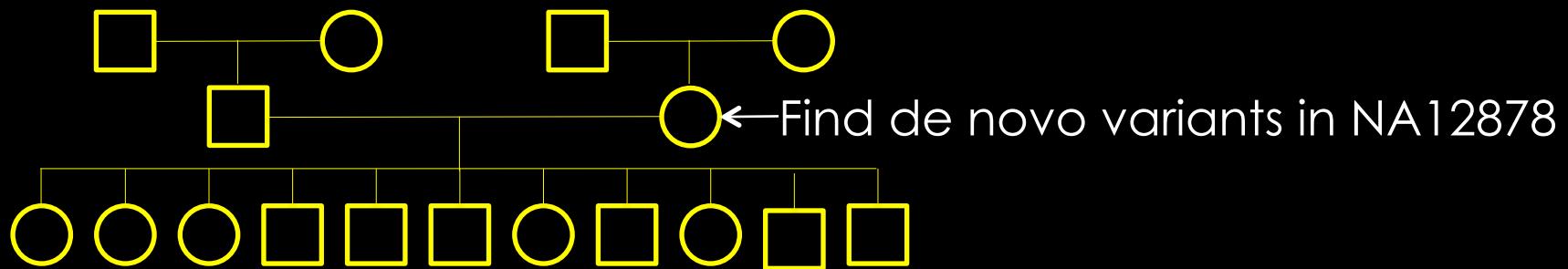
Subphenotypes

```
"LDL > 300"
```

```
"Maternal_infection = 1"
```

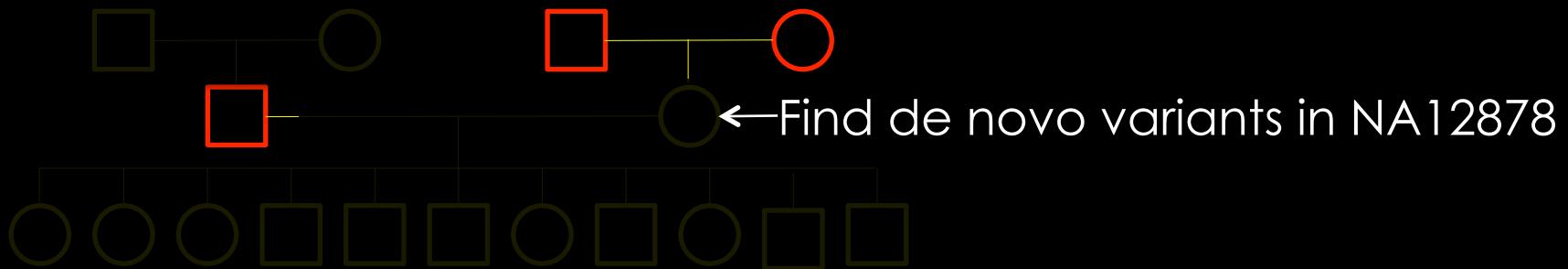
```
"Kidney_failure = 1"
```

Chaining Phenotype & Genotype Queries



```
gqt query -i ceph1463.gqt -d ceph1463.ped.db \
```

Chaining Phenotype & Genotype Queries



```
gqt query -i ceph1463.gqt -d ceph1463.ped.db \  
-p "sample_id in ('NA12891', 'NA12892', 'NA12877')" \  
-g "HOM_REF" \
```

Get all variants that are:

1. homozygous reference in the grandparents & father

Chaining Phenotype & Genotype Queries

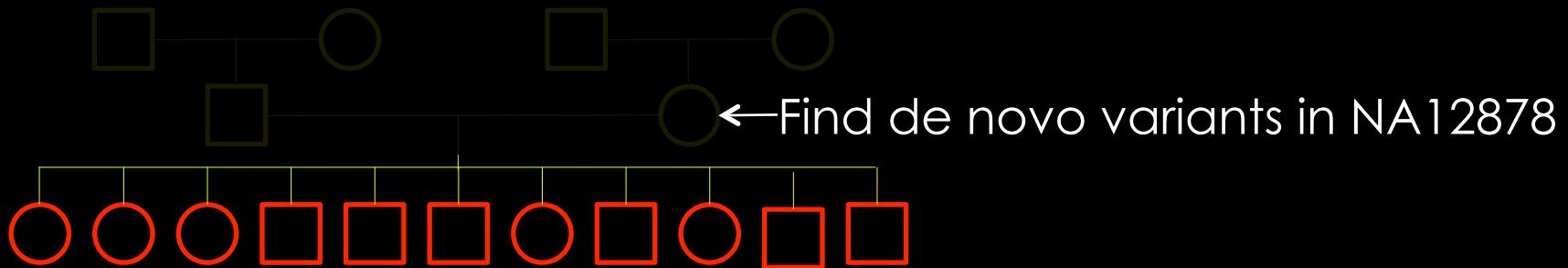


```
gqt query -i ceph1463.gqt -d ceph1463.ped.db \
-p "sample_id in ('NA12891', 'NA12892', 'NA12877')" \
-g "HOM_REF" \
-p "sample_id = 'NA12878'" \
-g "HET" \
```

Get all variants that are:

1. homozygous reference in the grandparents & father
2. heterozygous in the mother

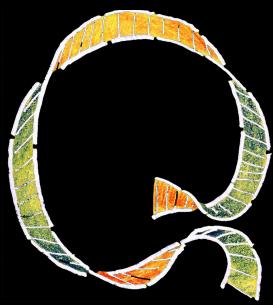
Chaining Phenotype & Genotype Queries



```
gqt query -i ceph1463.gqt -d ceph1463.ped.db \
-p "sample_id in ('NA12891', 'NA12892', 'NA12877')" \
-g "HOM_REF" \
-p "sample_id = 'NA12878'" \
-g "HET" \
-p "maternal_id = 'NA12878'" \
-g "count(HET) >=1"
```

Get all variants that are:

1. homozygous reference in the grandparents & father
2. heterozygous in the mother
3. heterozygous in at least 1 of the children



GENOTYPE QUERY TOOLS

github.com/ryanlayer/gqt

Millions of Individuals

100 Millions of genotypes

Fast Queries

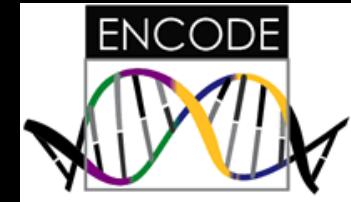
Powerful Interface



300M Intervals



8.5K Samples



7K Samples



2.5K Genomes

ExAC
125K Exomes
15K Genomes



5K Phenotypes
3K Genes



400 Tissues



127 Tissues
55M Intervals



100K Genomes

TOPMed
100K Genomes

CCDG
200K Genomes

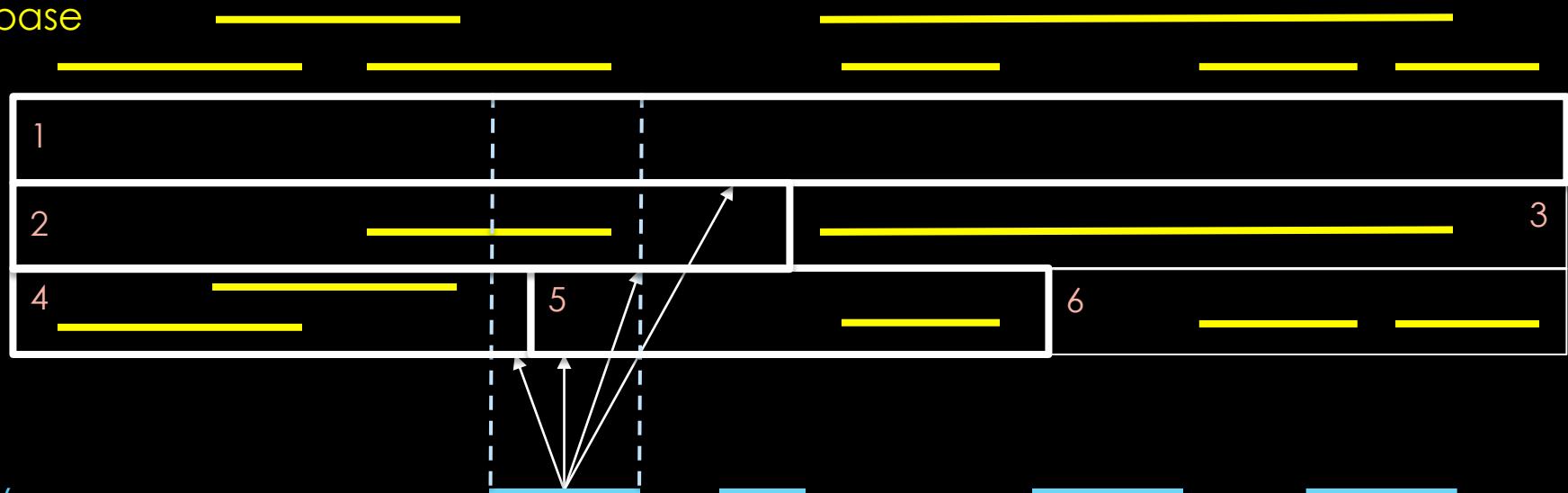
```
graph TD; GFF[GFF] --- NARROWPEAK[NARROWPEAK]; GFF --- VCF[VCF]; BROADPEAK[BROADPEAK] --- BEDPE[BEDPE]; CRAM[CRAM] --- BED[BED]; CRAM --- SAM[SAM]; CRAM --- BCF[BCF]; BED --- SAM; BED --- BCF; BAM[BAM] --- BIGWIG[BIGWIG]; WIG[WIG] --- BIGBED[BIGBED]; BIGWIG --- BEDGRAPH[BEDGRAPH];
```



Hierarchical Binning

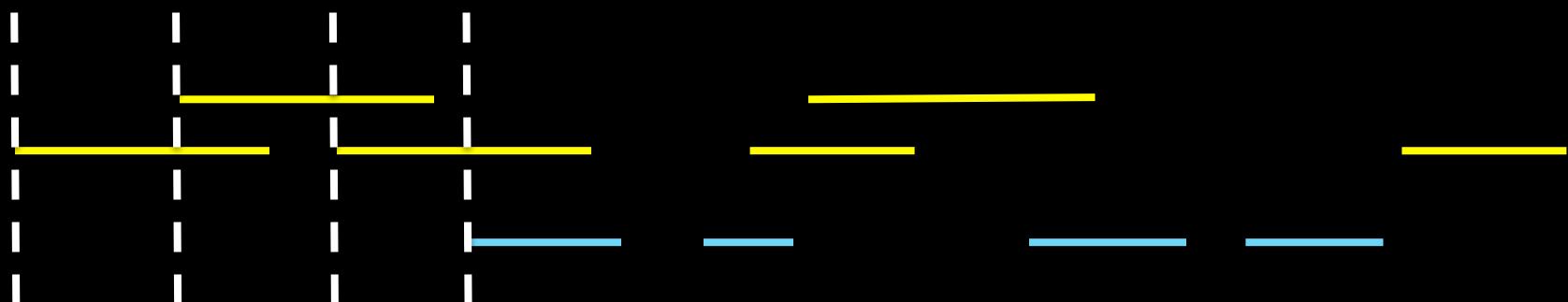
[UCSC Kent 2002, TABIX Li 2011, BEDTOOLS Quinlan 2010]

Database



Query

Sweep [BEDTOOLS]



TABIX

chr1:52000-53000



```
chr1 52057 52058 rs62637813
chr1 52237 52238 rs2691277
chr1 52726 52727 rs2691278
```



chr1:52000-53000



⋮

⋮



query size

100bp

empty

0.73

1000bp

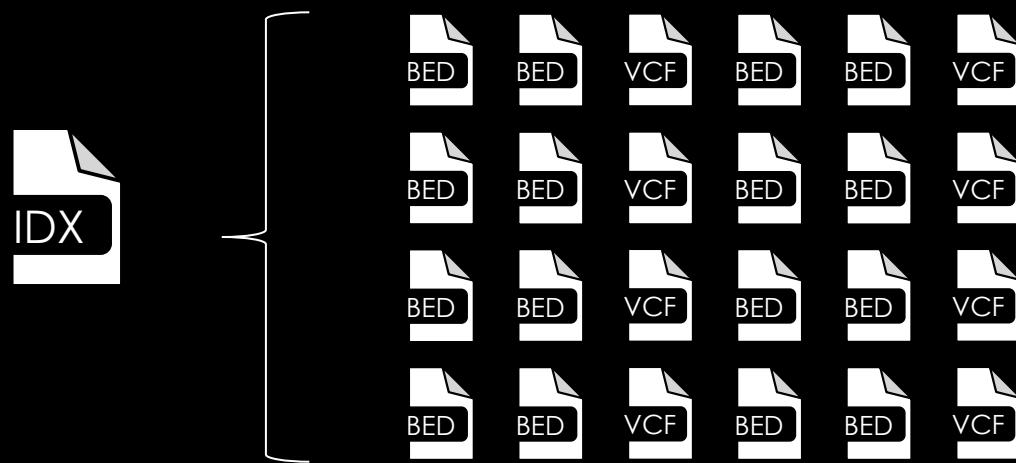
0.67

10000bp

0.58



B+ Tree index of positions and offsets





Query

of intersections

Operation

chr1:52000-53000



12039

Intersections per file

chr1:52000-53000



snp144

13

neandertalMethylation

5

genomicSuperDups

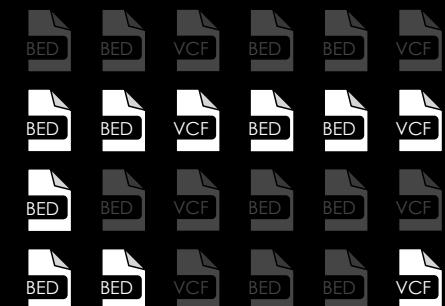
2

Intersecting intervals

chr1:52000-53000



```
#snp144
chr1 52057 52058 rs62637
chr1 52095 52096 rs36775
#neandertalMethylation
chr1 52015 52016 100
chr1 52028 52029 100
```





THE TIME INDEX:
AN ACCESS STRUCTURE FOR TEMPORAL DATA

Ramez Elmasri^{1,*}, Gene T. J. Wu², and Yeong-Joon Kim¹

¹Department of Computer Science, University of Houston, Houston, TX 77204

²Bell Communications Research, 444 Hoes Lane, Piscataway, NJ 08854

ABSTRACT

In this paper, we describe a new indexing technique, the *time index*, for improving the performance of certain classes of temporal queries. The time index can be used to retrieve versions of objects that are valid during a specific time period. It supports the processing of the temporal WHEN operator and temporal aggregate functions efficiently. The time indexing scheme is also extended to improve the performance of the temporal SELECT operator, which retrieves objects that satisfy a certain condition during a specific time period. We will describe the indexing technique, and its search and insertion algorithms. We also describe an algorithm for processing a commonly used temporal JOIN operation. Some results of a simulation for comparing the performance of the time index with other proposed temporal access structures are presented.

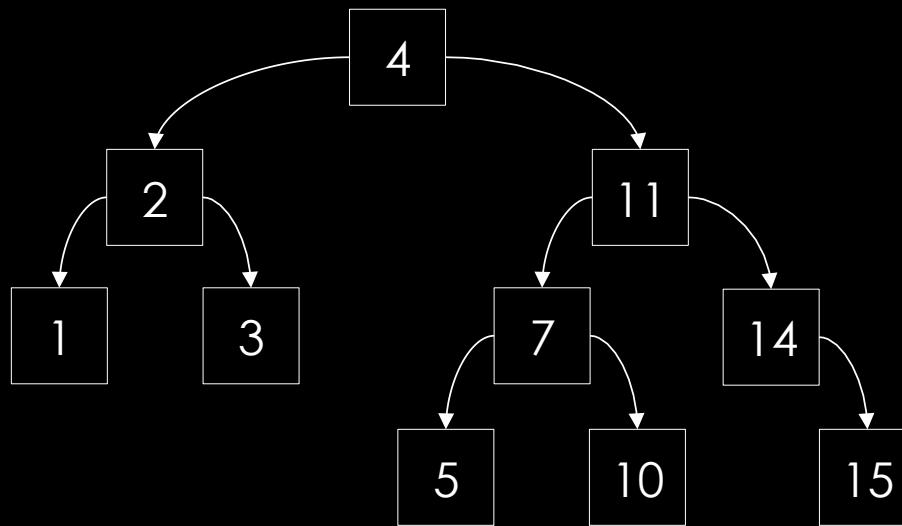
languages [SS87, EW90]. These temporal data models define powerful operations for specifying complex temporal queries. There has been relatively less research in the area of defining efficient storage structures and access paths for temporal data [Lum84, Ahn86, AS88, SG89, GSsu, RS87, KS89, LS89]. These proposals do not discuss indexing schemes for supporting the high-level temporal operators defined in [GY88, EW90]. This paper describes indexing techniques for improving the efficiency of temporal operations, such as *when*, *select*, and *join* [GY88], *temporal selection* and *temporal projection* [EW90], and aggregation functions.

The storage techniques for temporal data proposed in [AS88, Lum84] index or link the versions of each individual object separately. In order to retrieve object versions that are valid during a certain time period, it is necessary to first locate the first (current) version of each object, and then search through the version in-

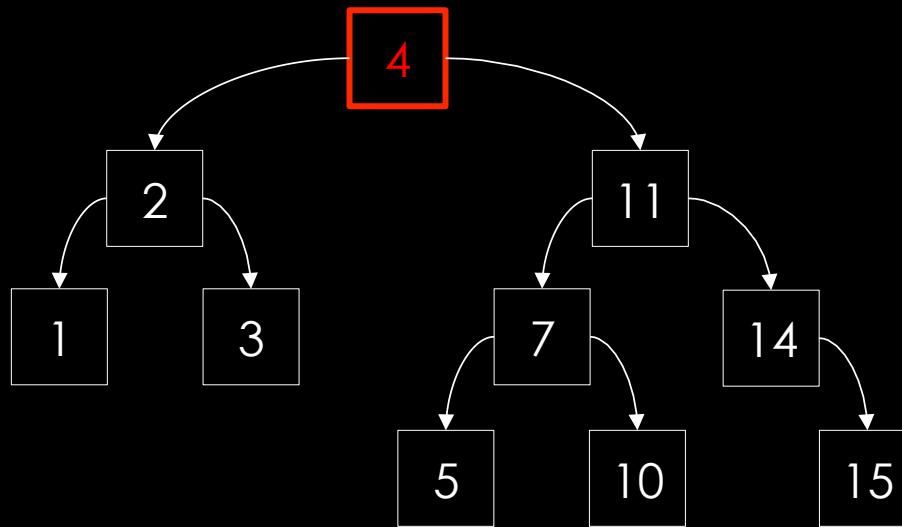
VLDB 1990

B+ Tree: efficient on-disk retrieval
efficient in-order traversal

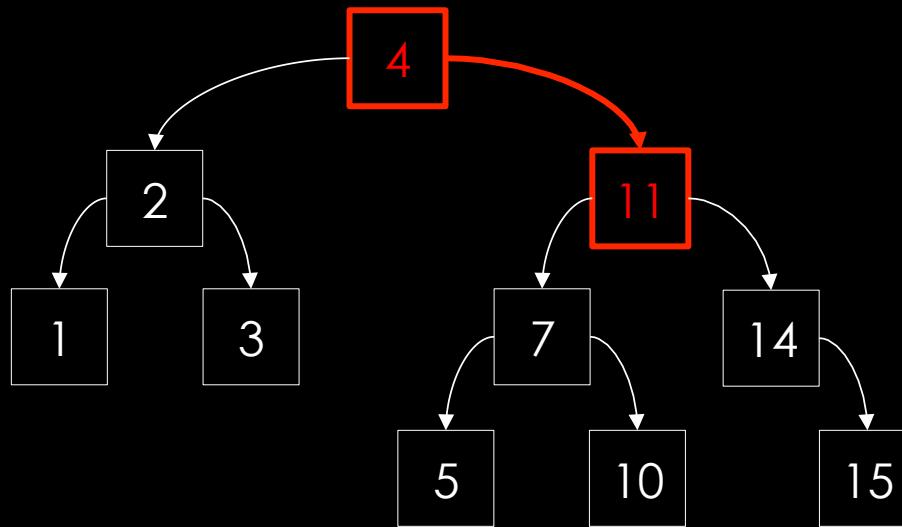
Binary Tree



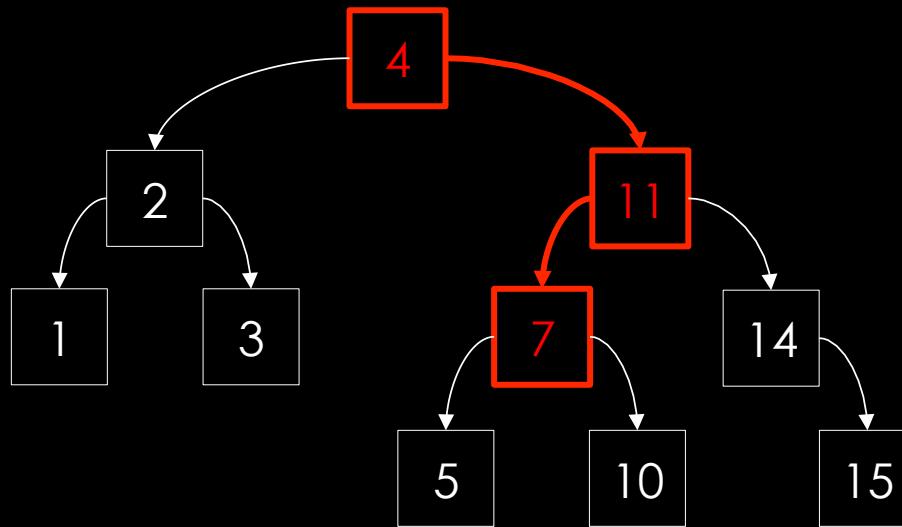
Binary Tree



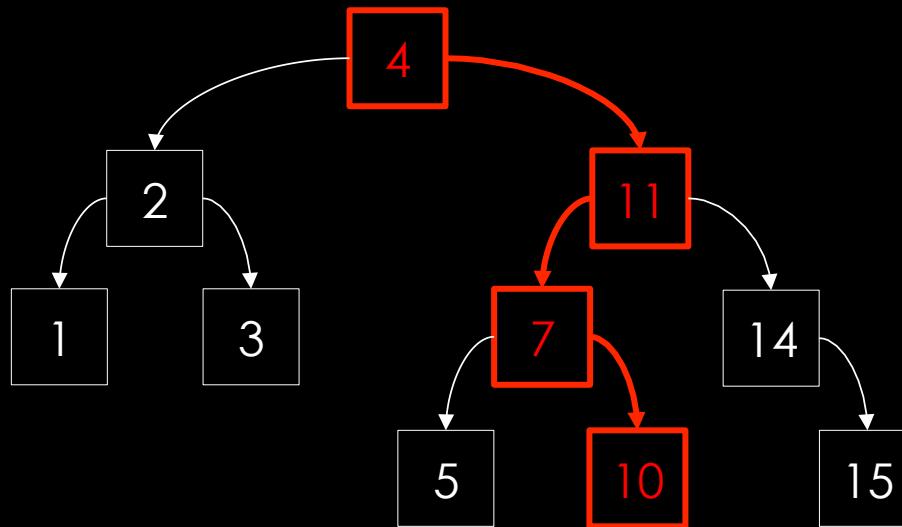
Binary Tree



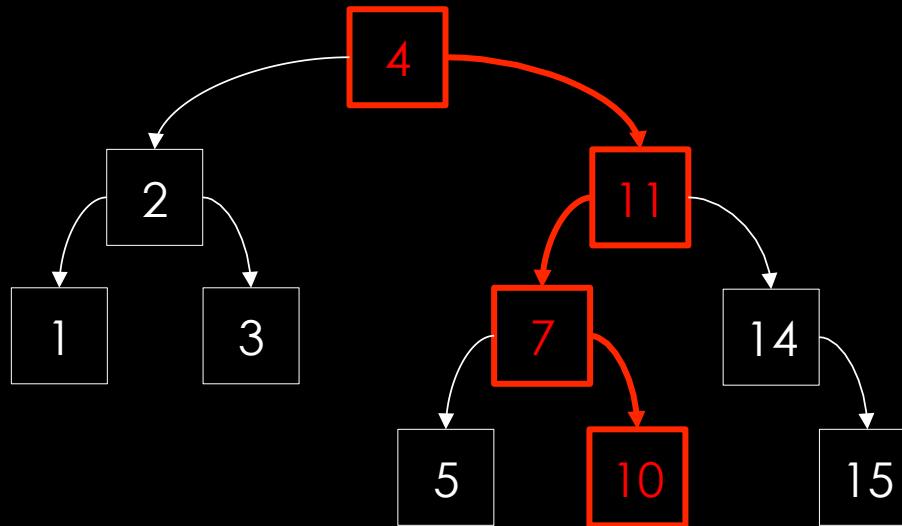
Binary Tree



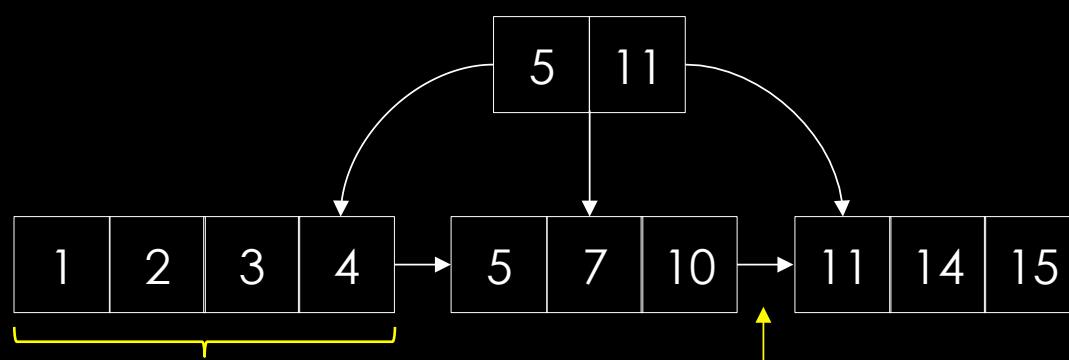
Binary Tree



Binary Tree



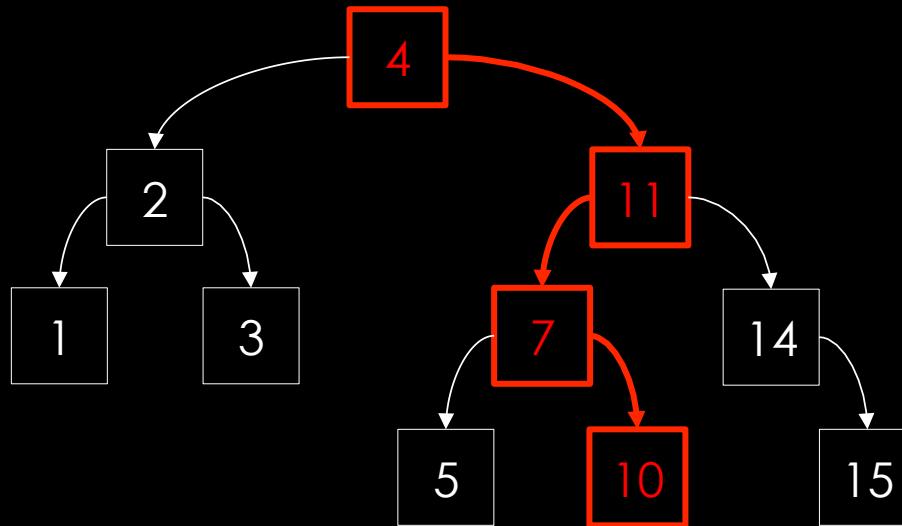
B+ Tree



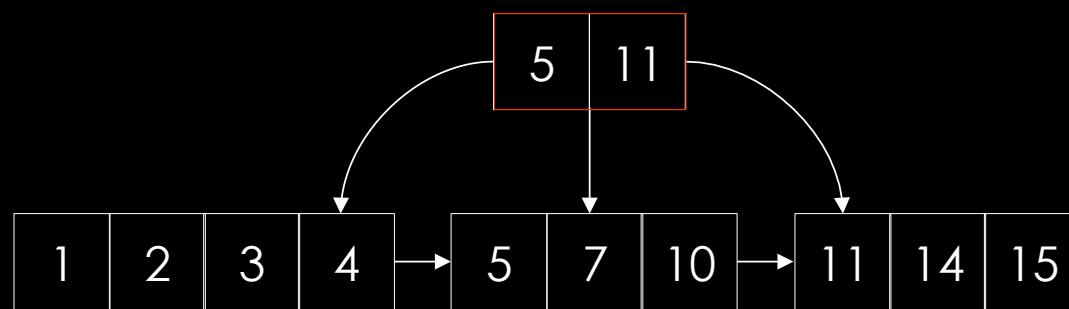
at most N (4) keys per node

linked leaves

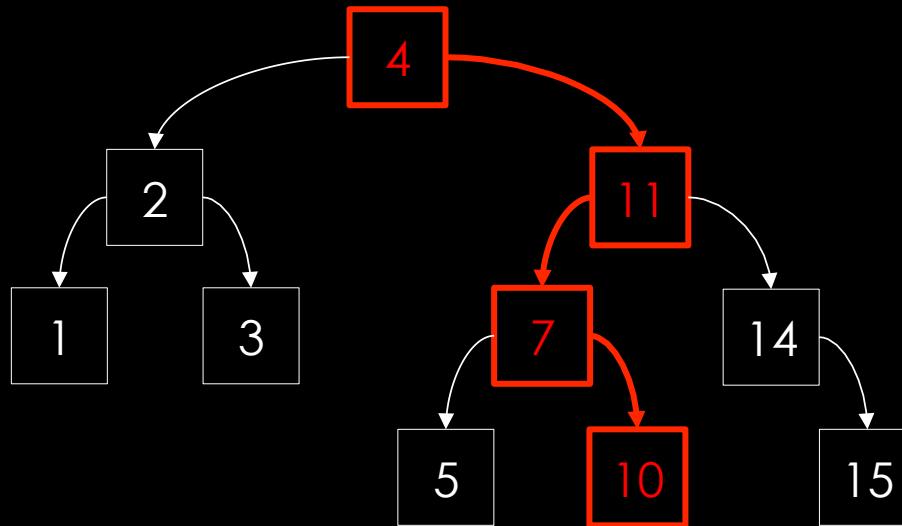
Binary Tree



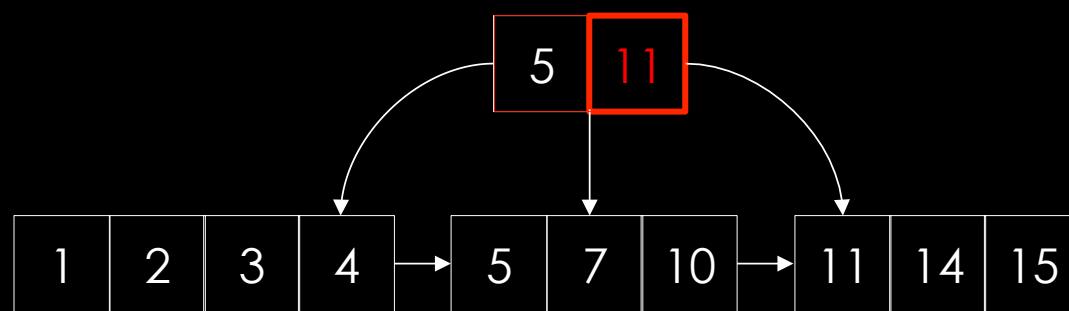
B+ Tree



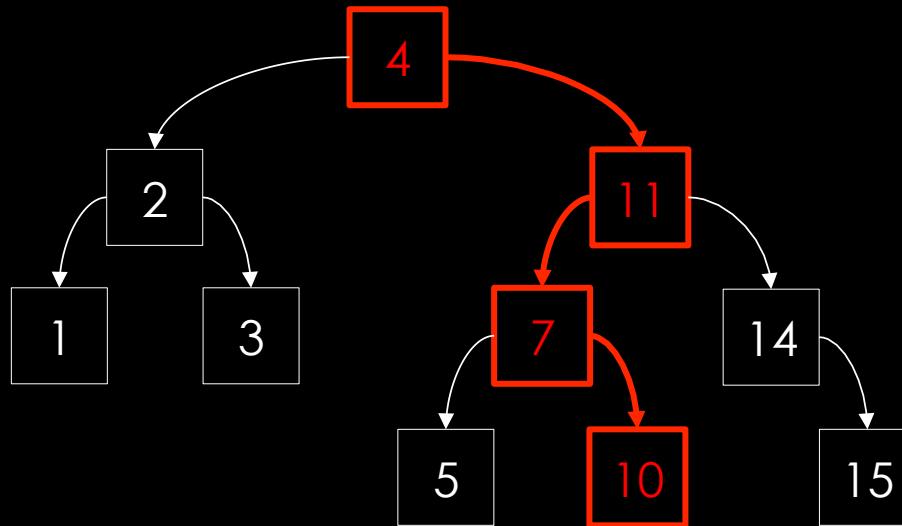
Binary Tree



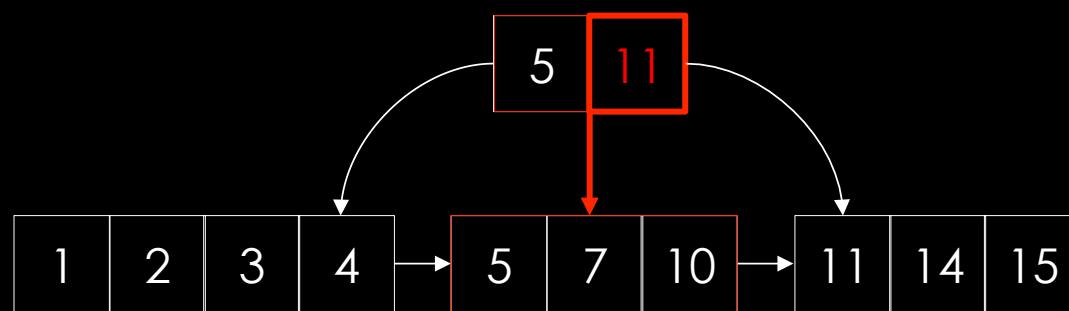
B+ Tree



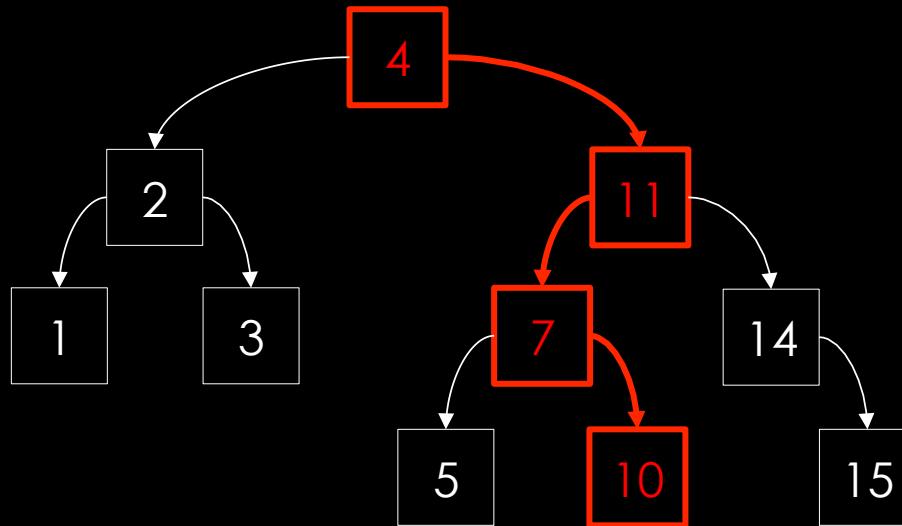
Binary Tree



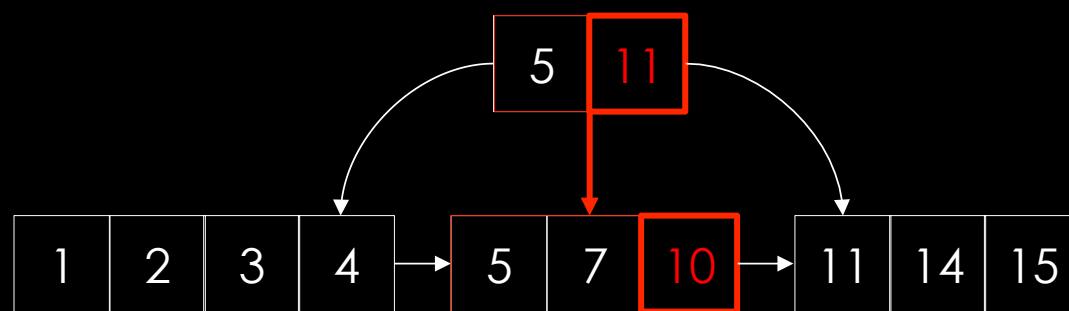
B+ Tree



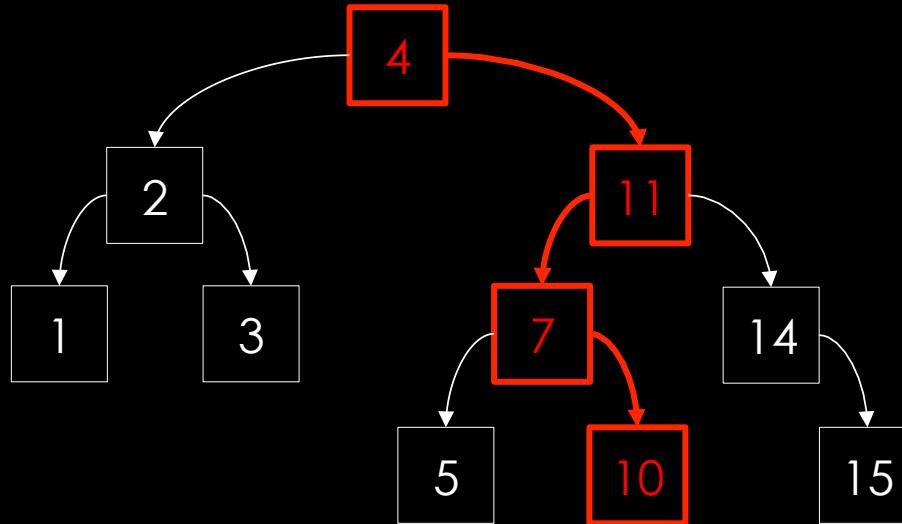
Binary Tree



B+ Tree



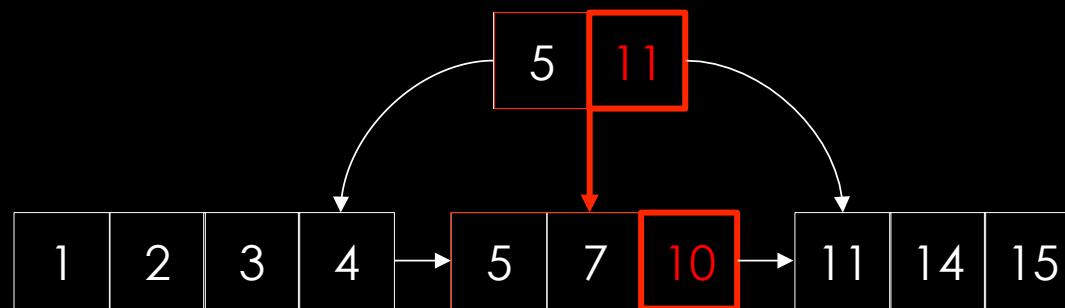
Binary Tree



Disk layout:



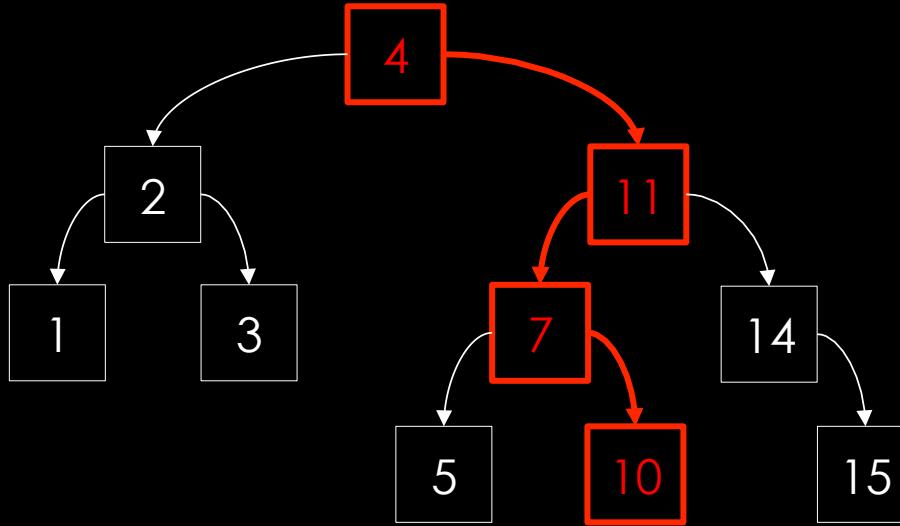
B+ Tree



Disk layout:



Binary Tree



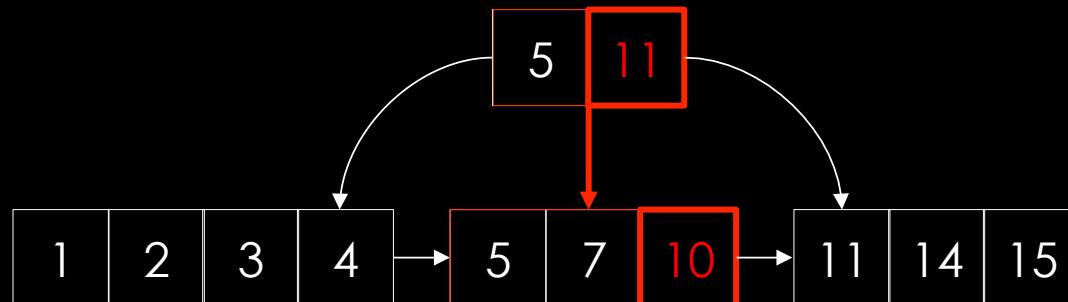
4 disk reads:



SSD disk speed
100,000 IOPS

~2,000,000X
slower than CPU

B+ Tree



2 disk reads:



1 2 3 4 5 6 7 8 9 10 11 12 13 14



1 2 3 4 5 6 7 8 9 10 11 12 13 14



```
index(X, start, end)
insert +X at start
insert -X at end + 1
append X at spanned L's
```

1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

A3

B

B1

B2

C

C2

index(A1, 1, 9)

MAX_KEYS = 4

1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A3

A1

B

B1

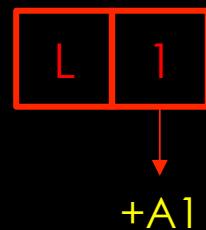
B2

C

C2

C1

index(A1, 1, 9)



+A1

1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A3

A1

B

B1

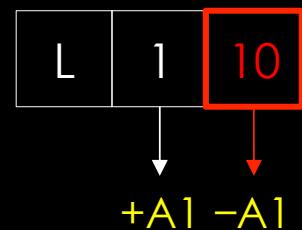
B2

C

C2

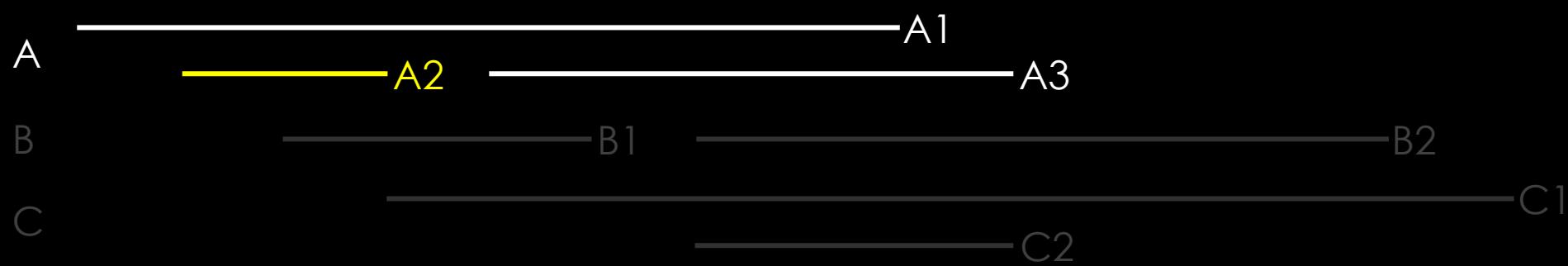
C1

index(A1, 1, 9)

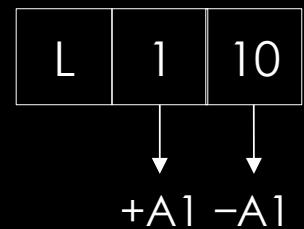


+A1 -A1

1 2 3 4 5 6 7 8 9 10 11 12 13 14



index(A2, 2, 4)



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

B1

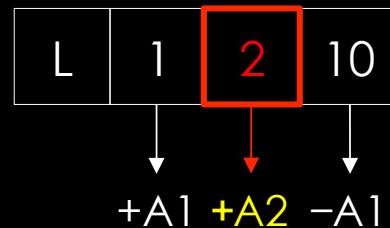
A3

B2

C

C2

index(A2, 2, 4)



+A1 +A2 -A1

1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

B1

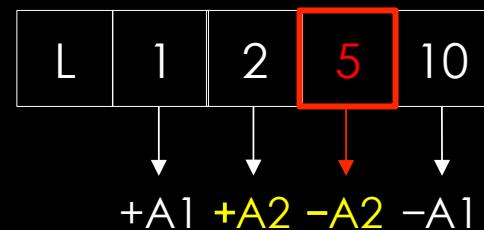
A3

B2

C

C2

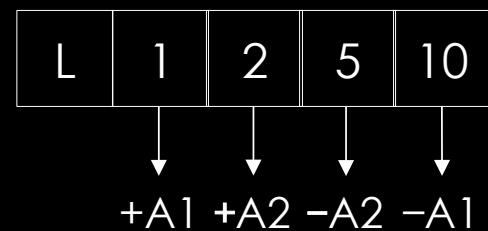
index(A2, 2, 4)



1 2 3 4 5 6 7 8 9 10 11 12 13 14



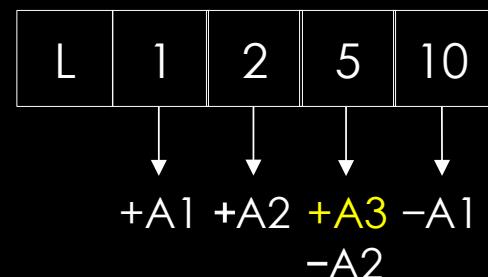
index(A3, 5, 10)



1 2 3 4 5 6 7 8 9 10 11 12 13 14



index(A3, 5, 10)



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

B1

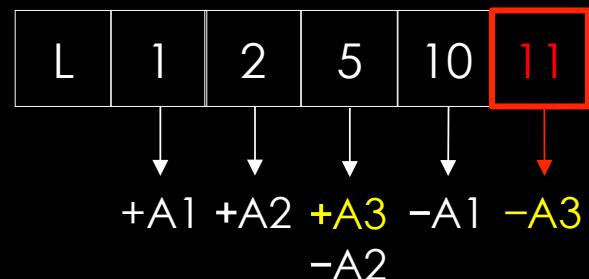
A3

B2

C

C2

index(A3, 5, 10)

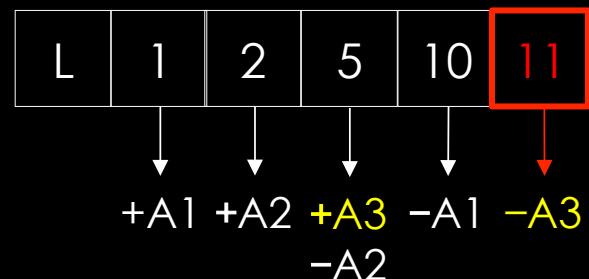


1 2 3 4 5 6 7 8 9 10 11 12 13 14

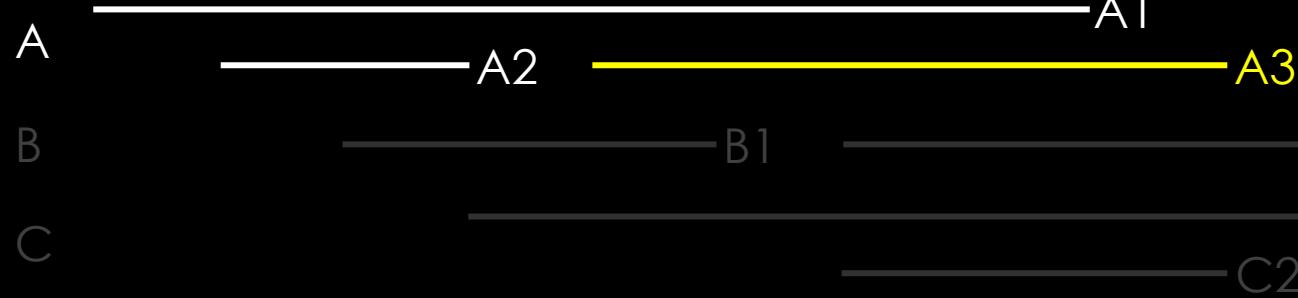


index(A3, 5, 10)

— Too Many Keys —

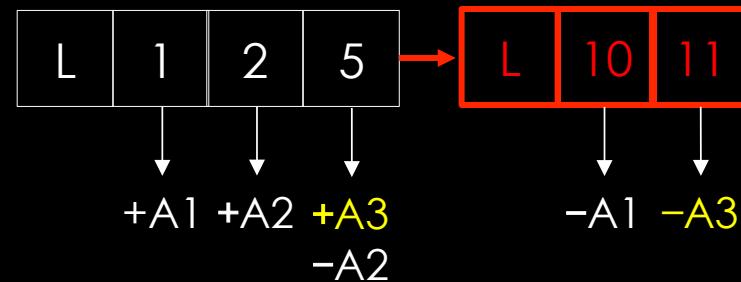


1 2 3 4 5 6 7 8 9 10 11 12 13 14



index(A3, 5, 10)

Split Leaf



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

B1

A3

B2

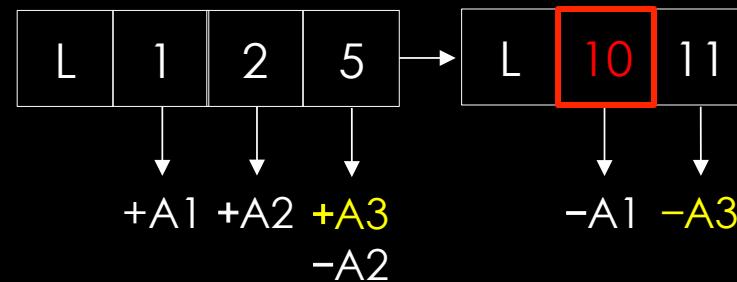
C

C2

C1

index(A3, 5, 10)

Promote a New Root



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

B1

A3

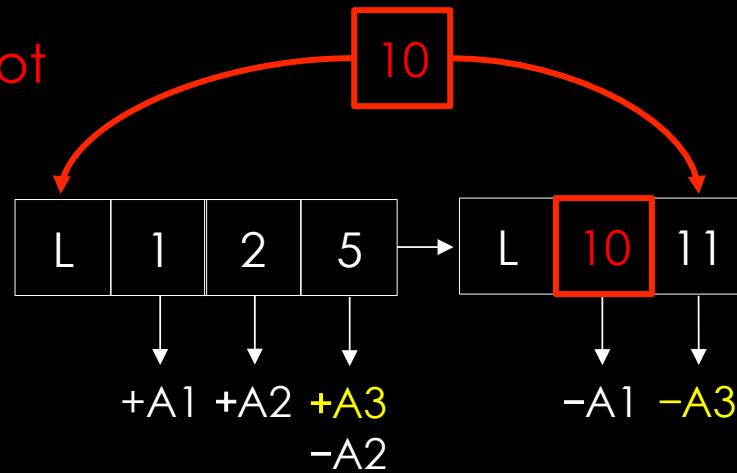
B2

C

C2

index(A3, 3, 6)

Promote a New Root



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

B1

A3

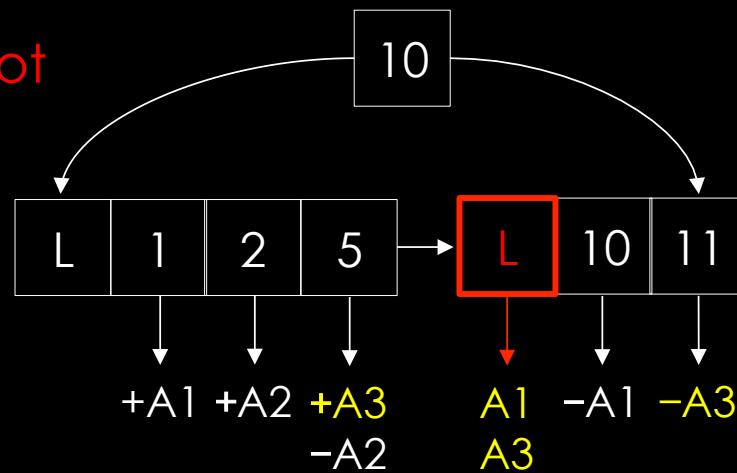
B2

C

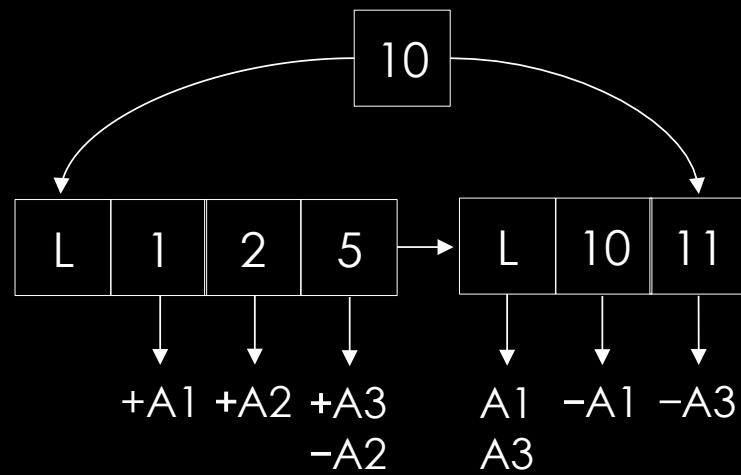
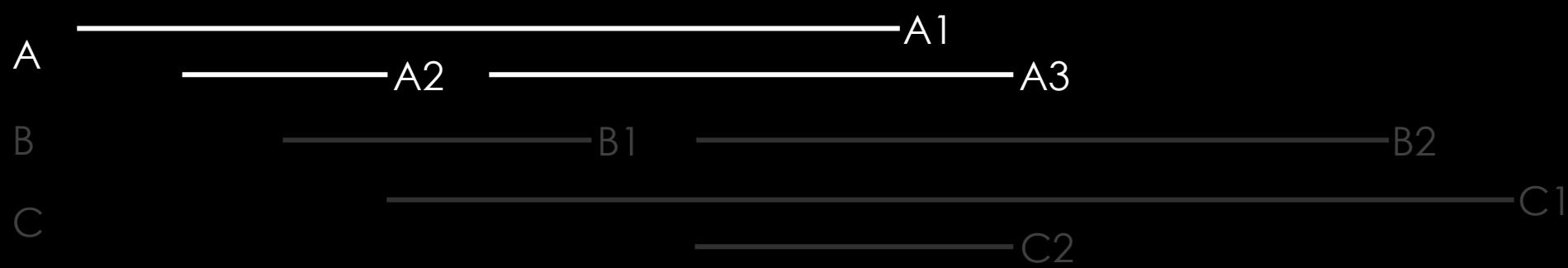
C2

index(A3, 3, 6)

Promote a New Root



1 2 3 4 5 6 7 8 9 10 11 12 13 14



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

A3

B

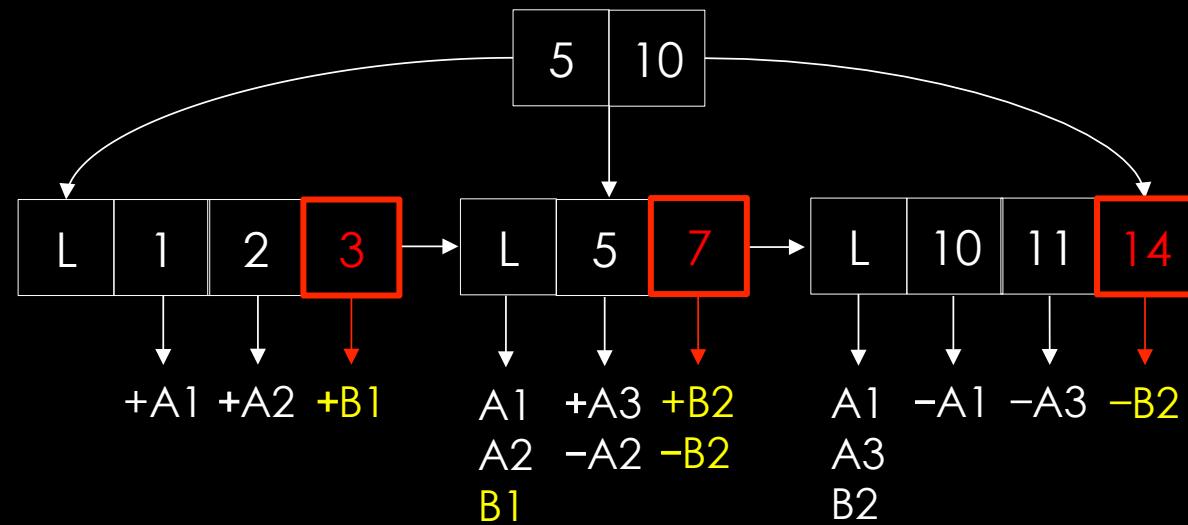
B1

B2

C

C1

C2



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

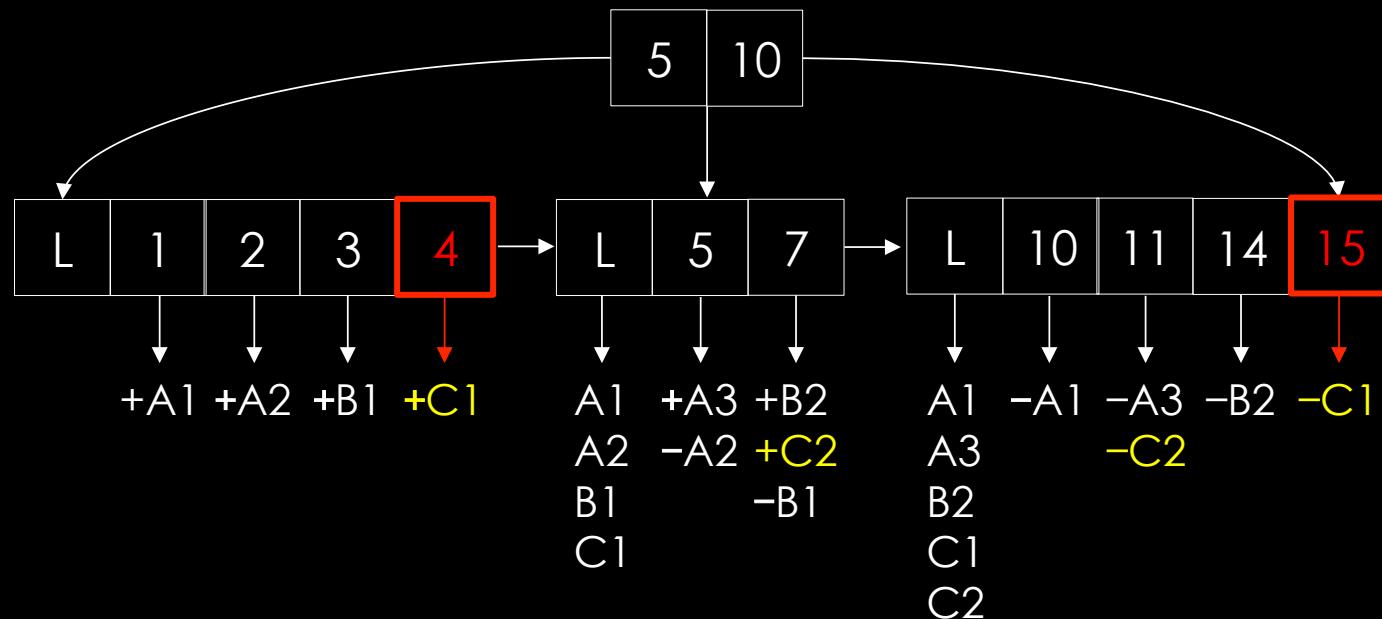
B1

B2

C

C2

C1

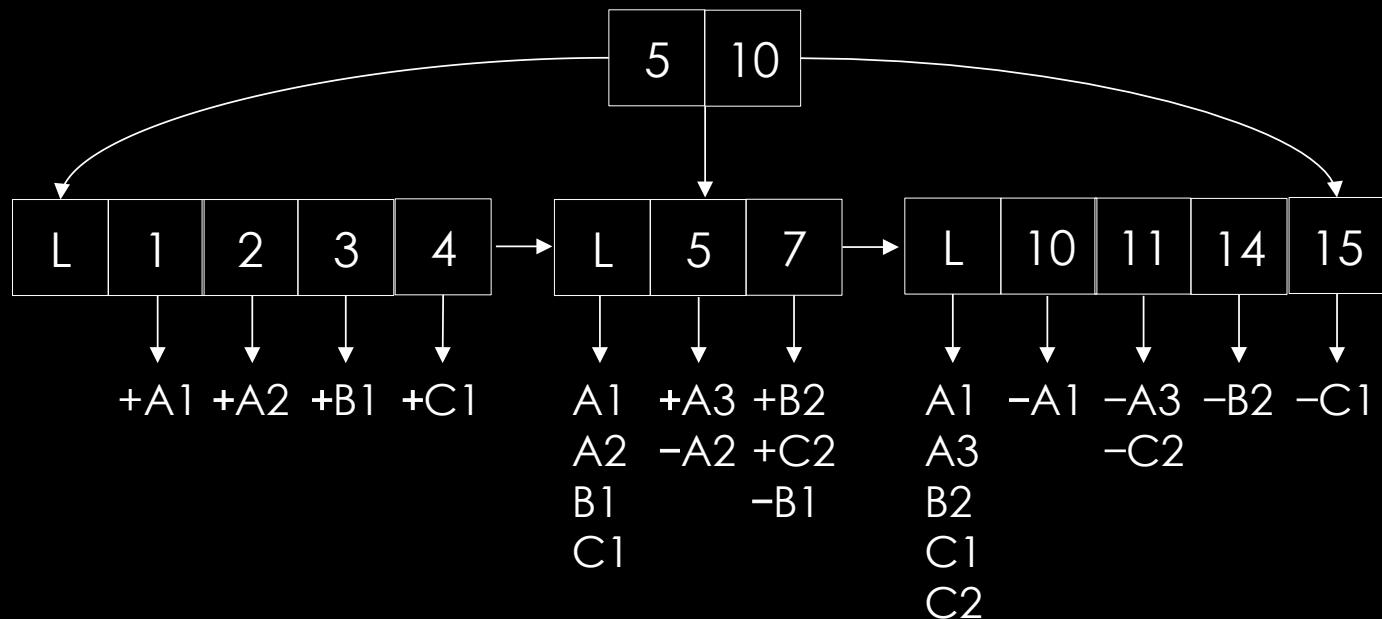


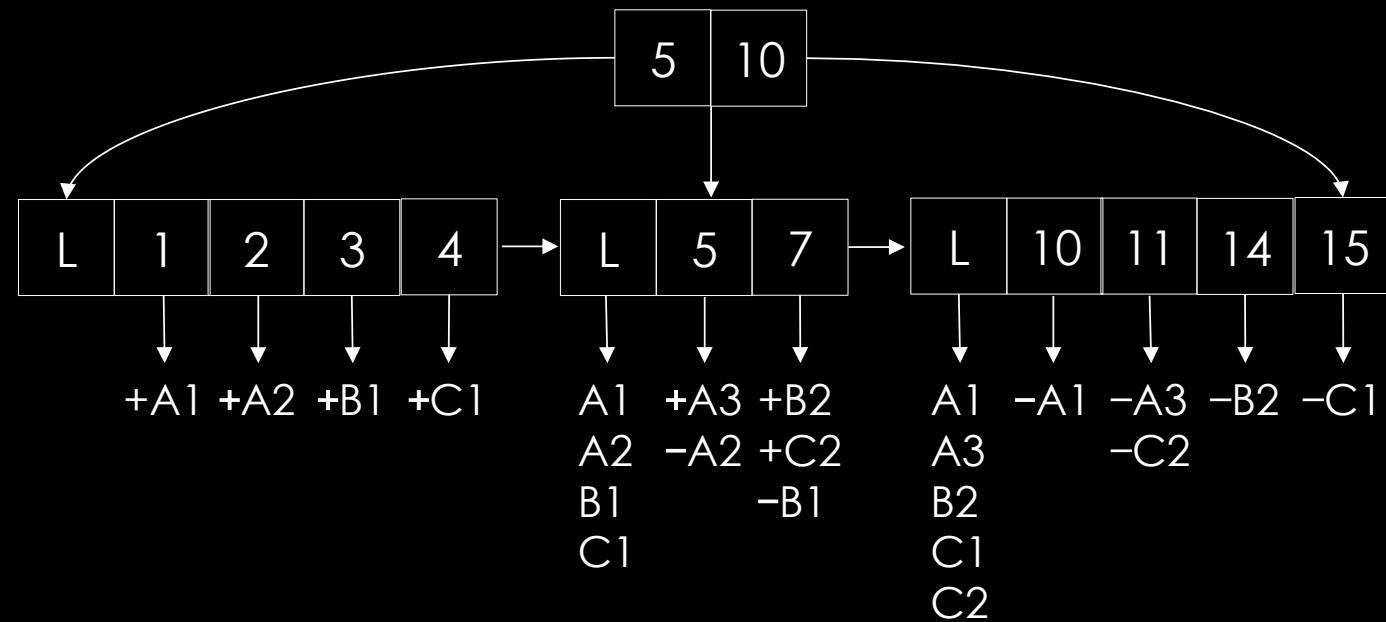
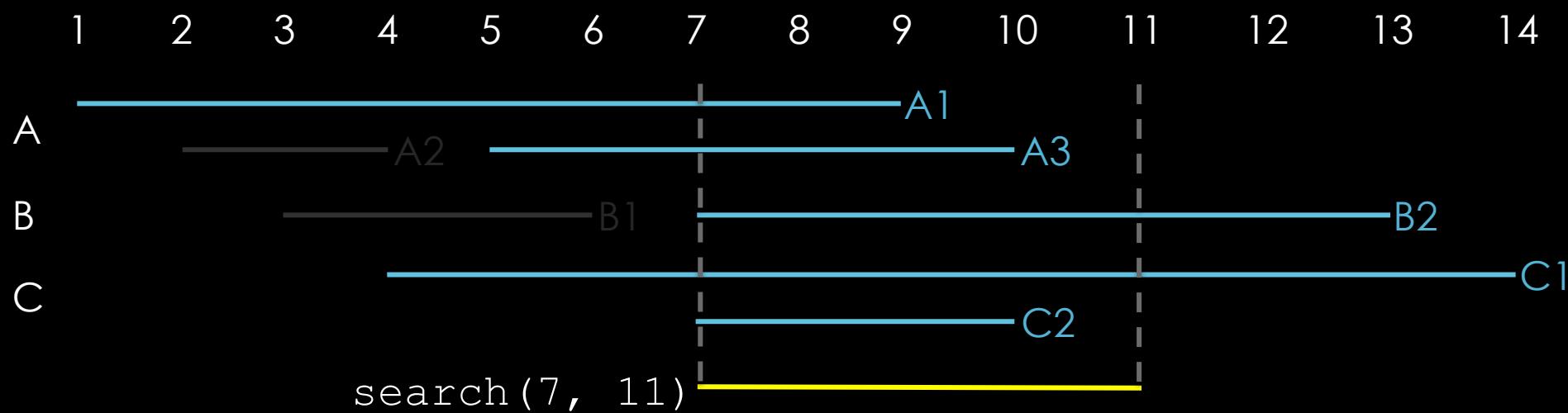
1 2 3 4 5 6 7 8 9 10 11 12 13 14



```
search(start, end)
    I = []
    from = find(start)
    to = find(end);
    I += from.leaf.L
    for node upto from:
        I += node.starts
        I -= node.ends
    for node in [from.next, to]:
        I += node.starts
```

1 2 3 4 5 6 7 8 9 10 11 12 13 14





1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

B

C

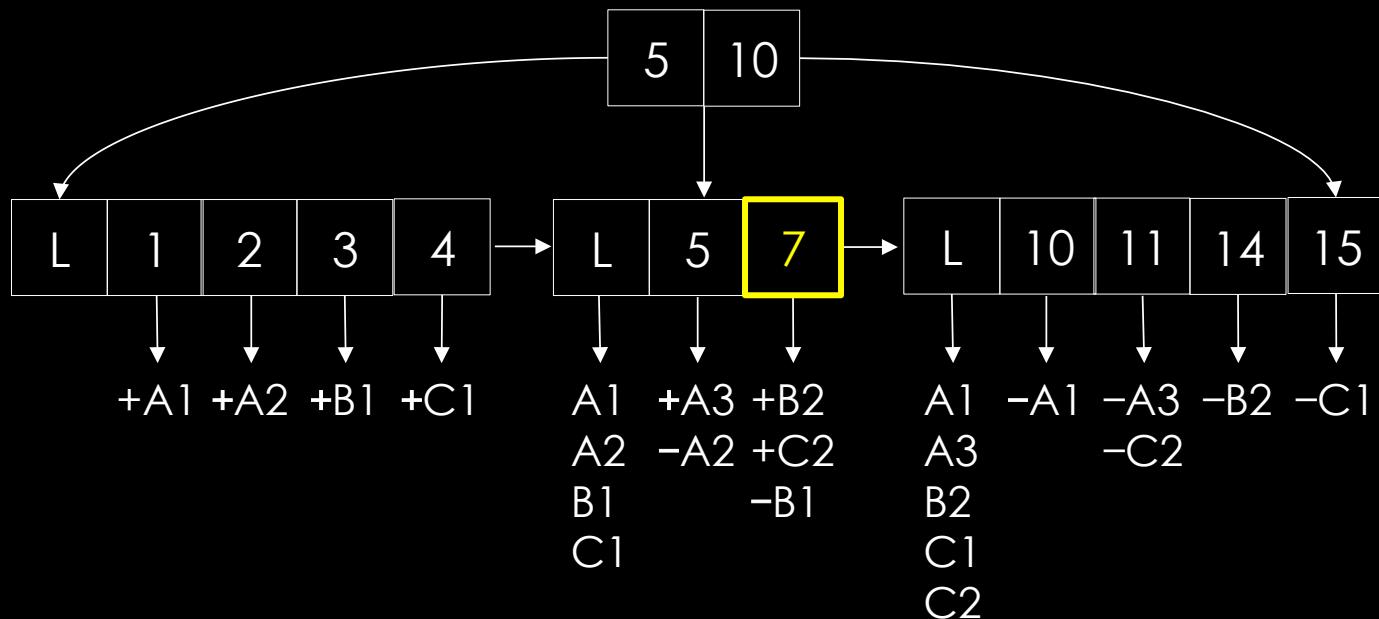
A1
A2
A3

B1
B2

C1
C2

search(7, 11)

I = []



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

B

C

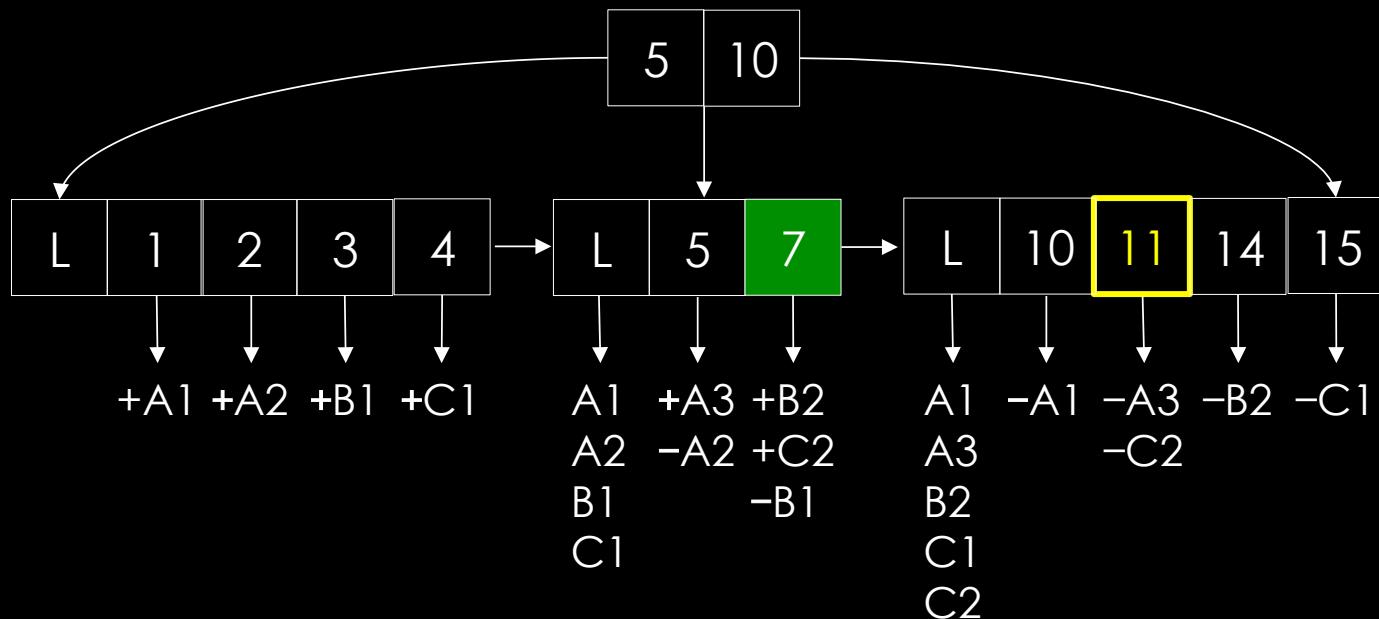
A1
A2
A3

B1
B2

C1
C2

search(7, 11)

I = []



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

B

C

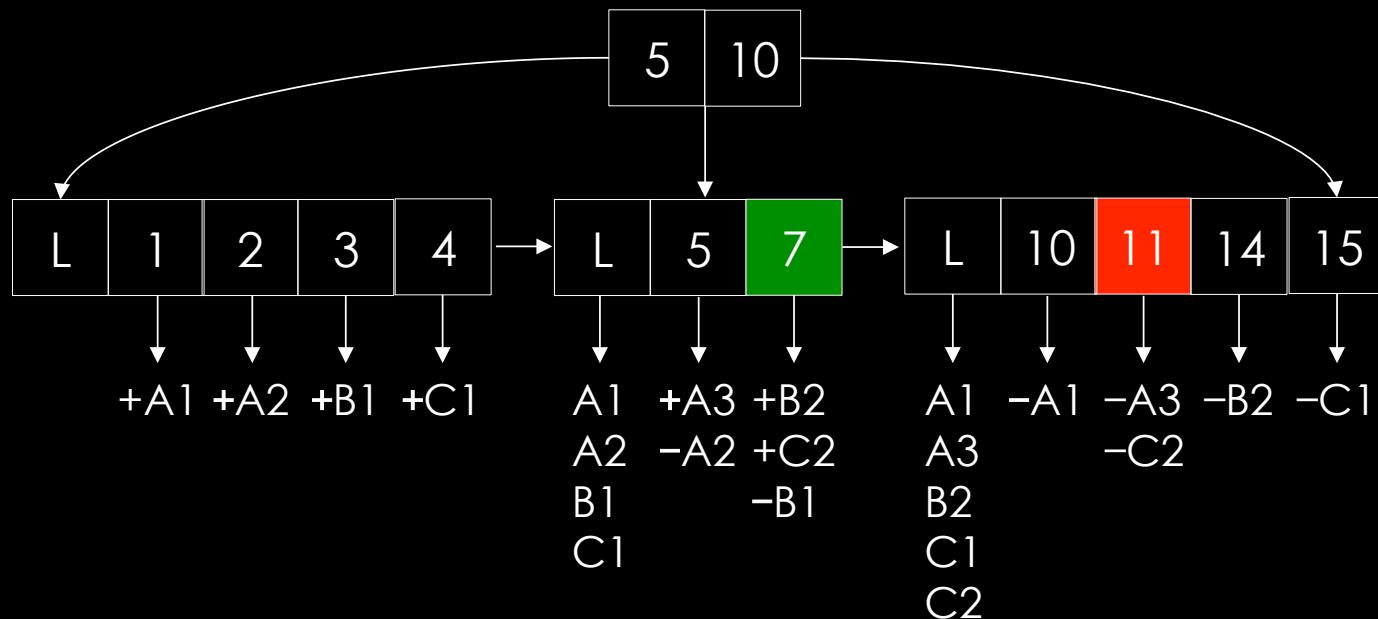
A1
A2
A3

B1
B2

C1
C2

search(7, 11)

I = []



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

B

C

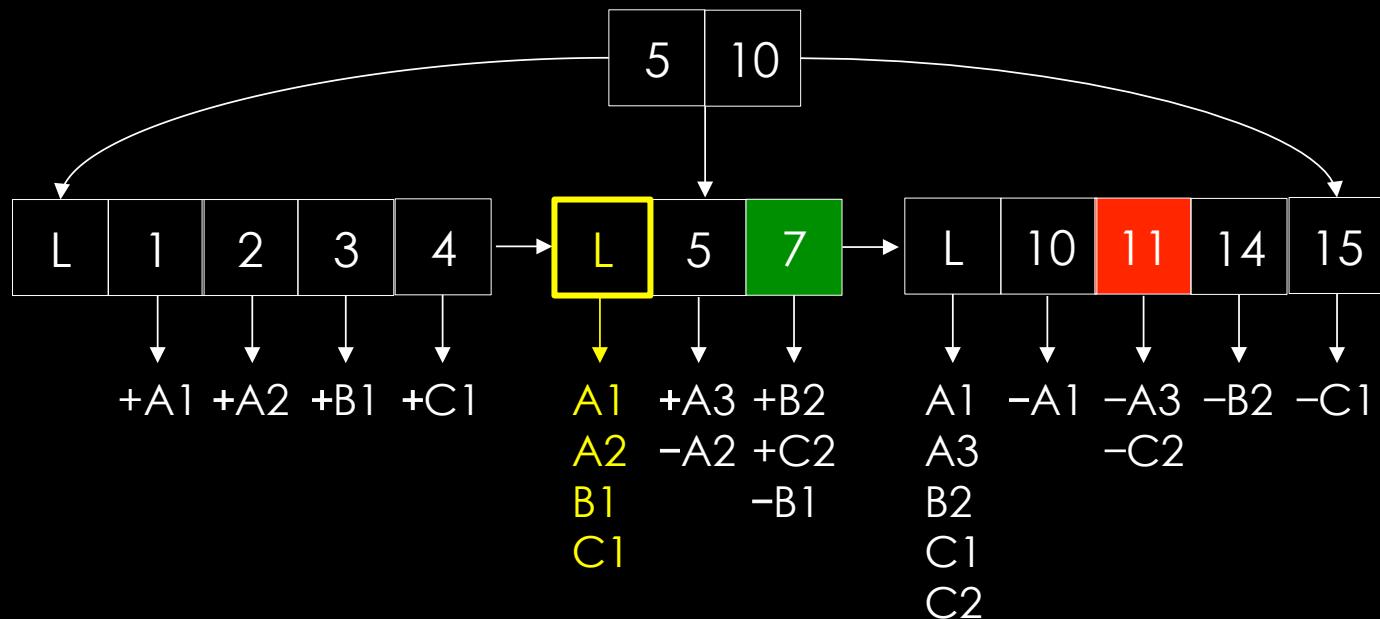
A1
A2
A3

B1
B2

C1
C2

search(7, 11)

I = [A1 A2 B1 C1]



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A3

B

B1

B2

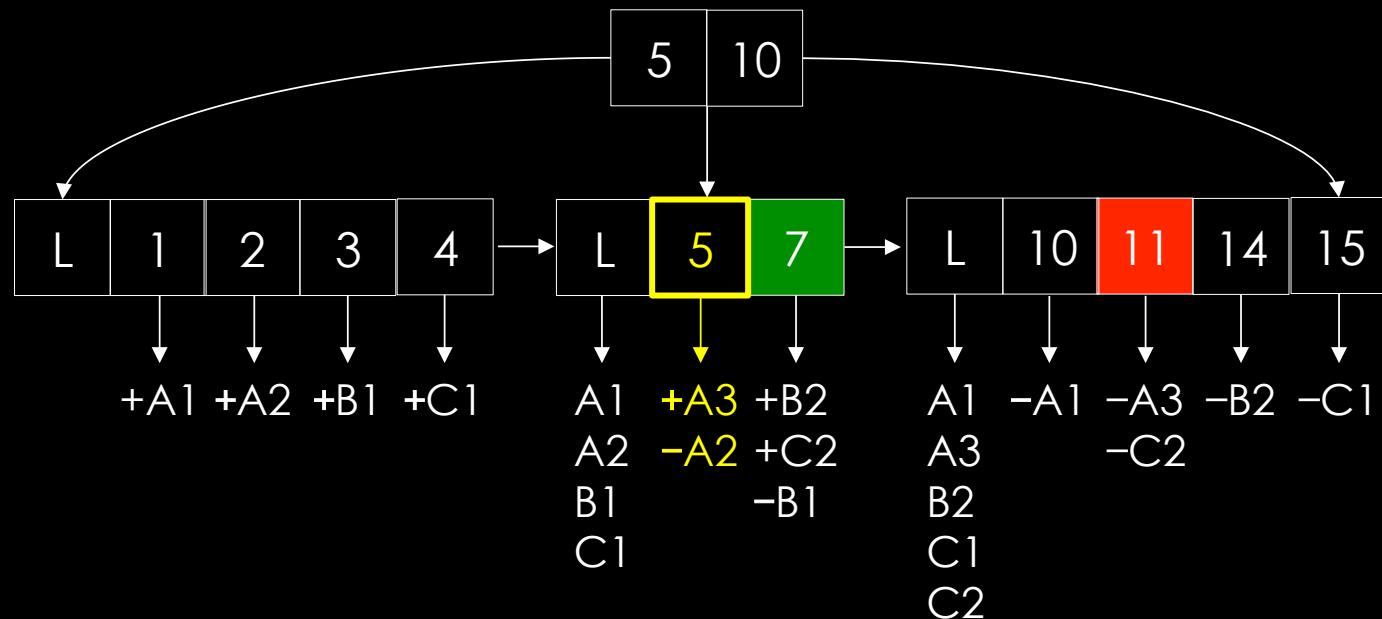
C

C2

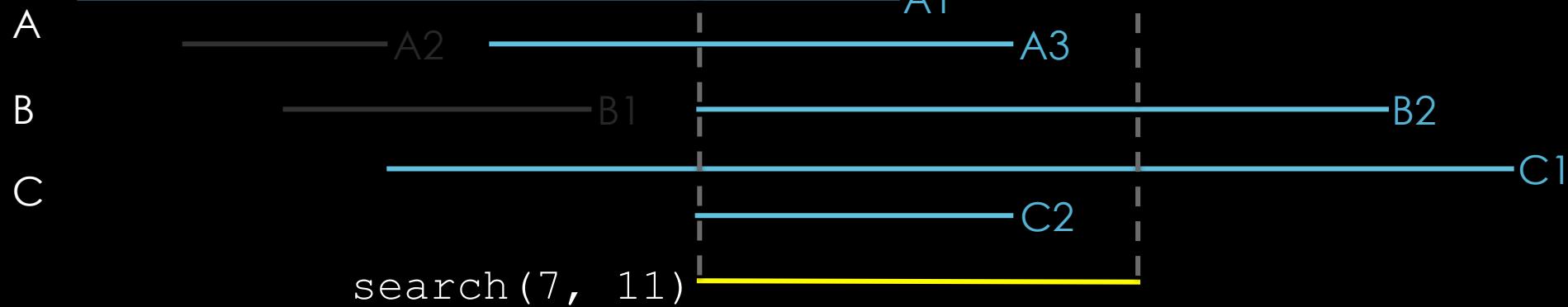
C1

search(7, 11)

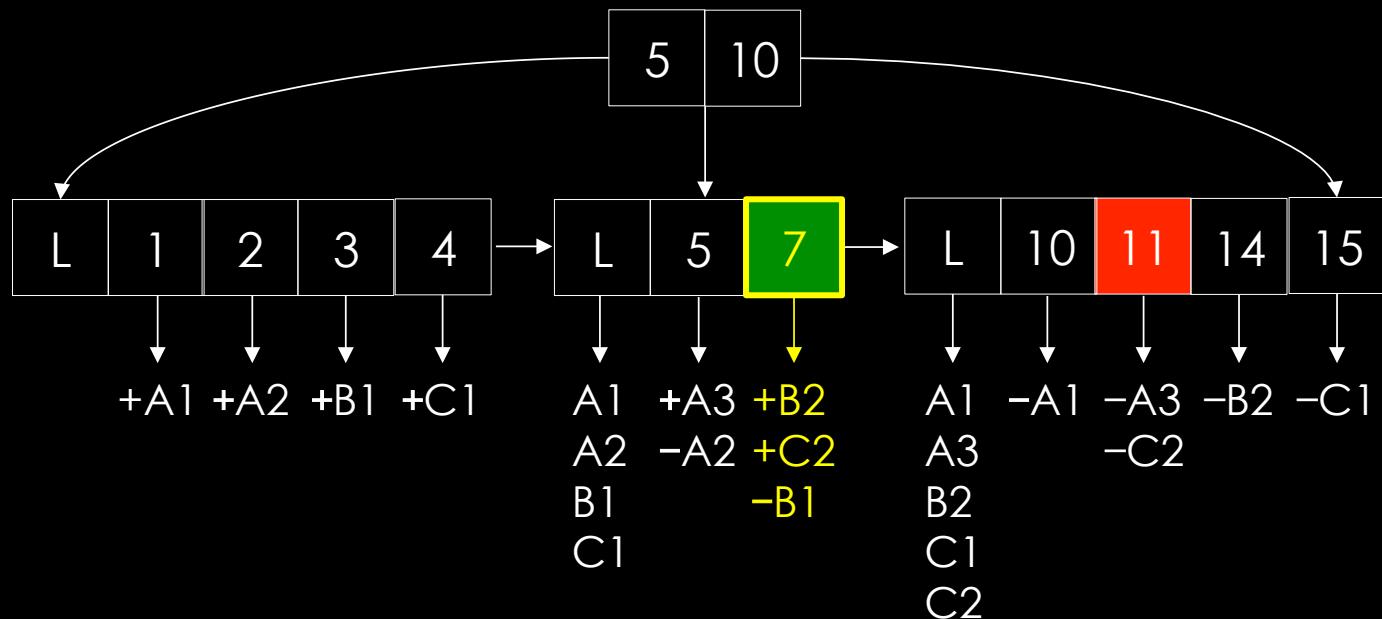
$$I = [A_1 \ A_2 \ B_1 \ C_1] + [A_3] - [A_2] = [A_1 \ A_3 \ B_1 \ C_1]$$

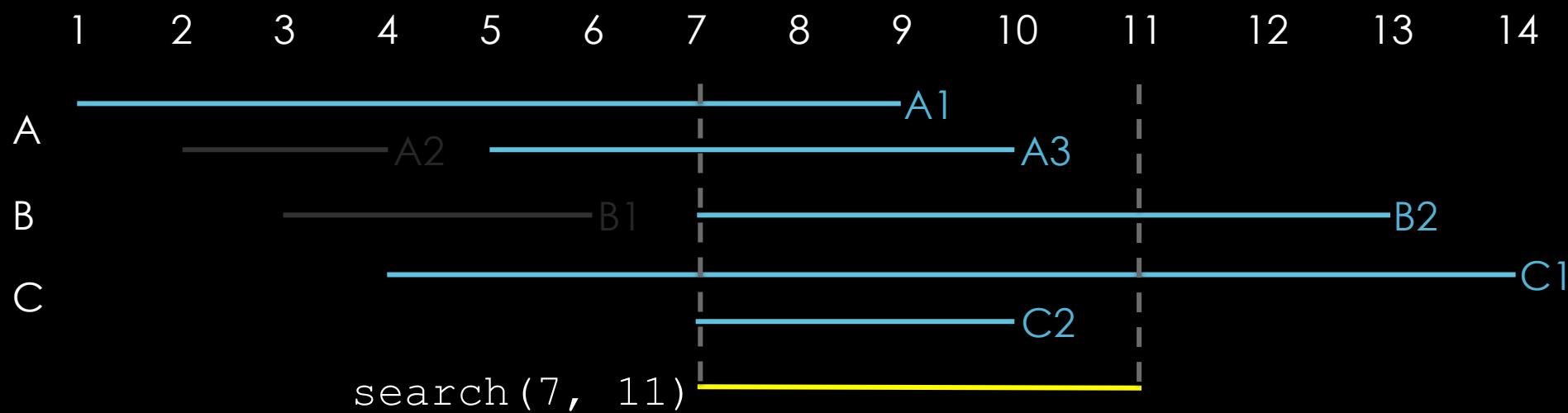


1 2 3 4 5 6 7 8 9 10 11 12 13 14

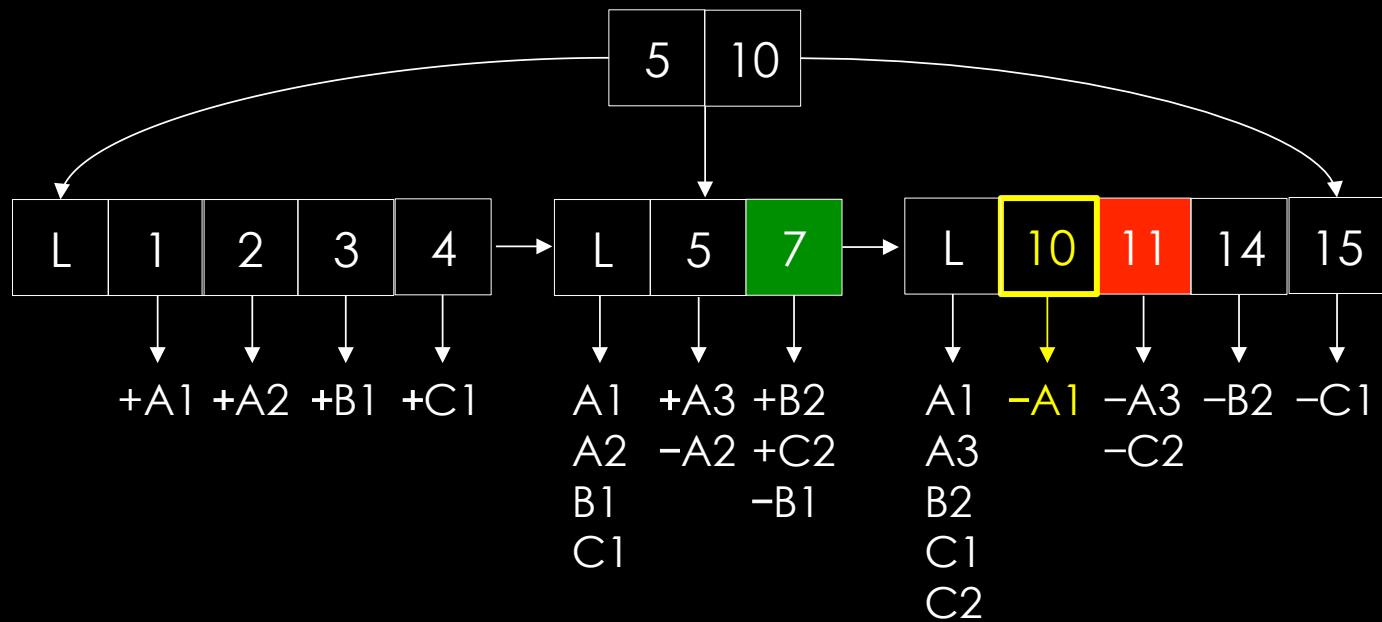


$$I = [A1 \ A3 \ B1 \ C1] + [B2 \ C2] - [B1] = [A1 \ A3 \ B2 \ C1 \ C3]$$





$$I = [A_1 \ A_3 \ B_2 \ C_1 \ C_3]$$



1 2 3 4 5 6 7 8 9 10 11 12 13 14

A

A2

A1

B

B1

A3

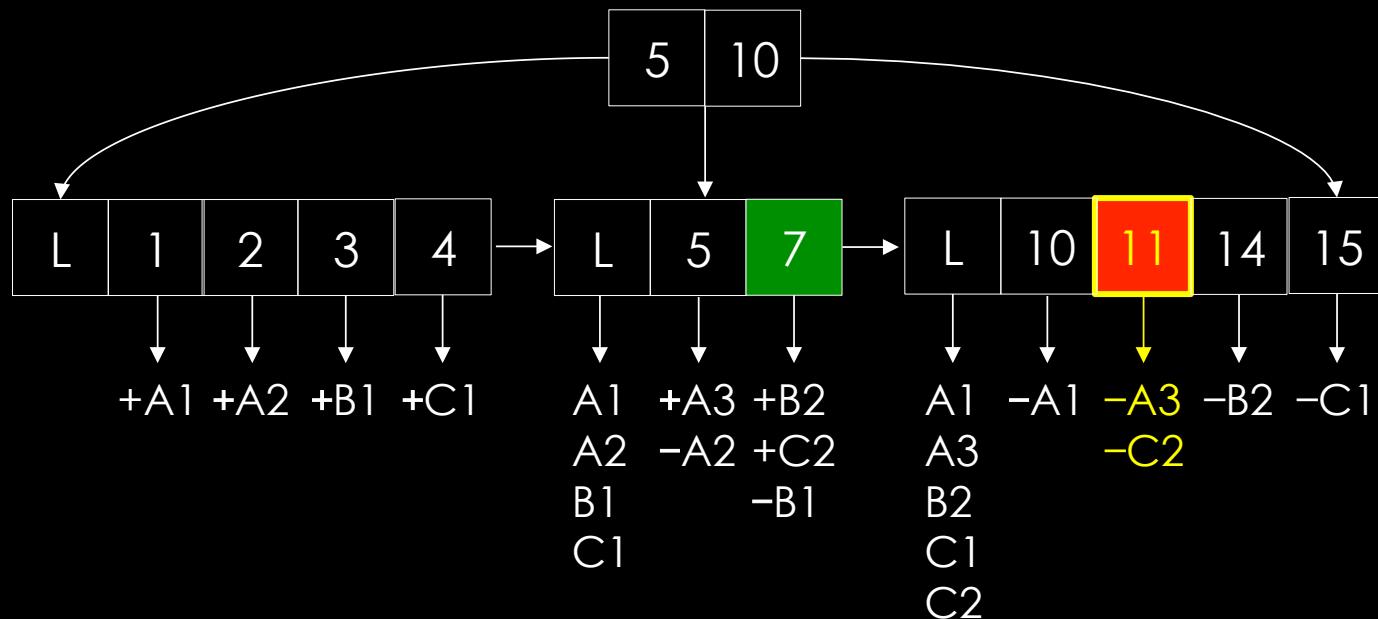
B2

C

C2

search(7, 11)

I = [A1 A3 B2 C1 C3]



Count intersections for 10 - 1M 100bp intervals



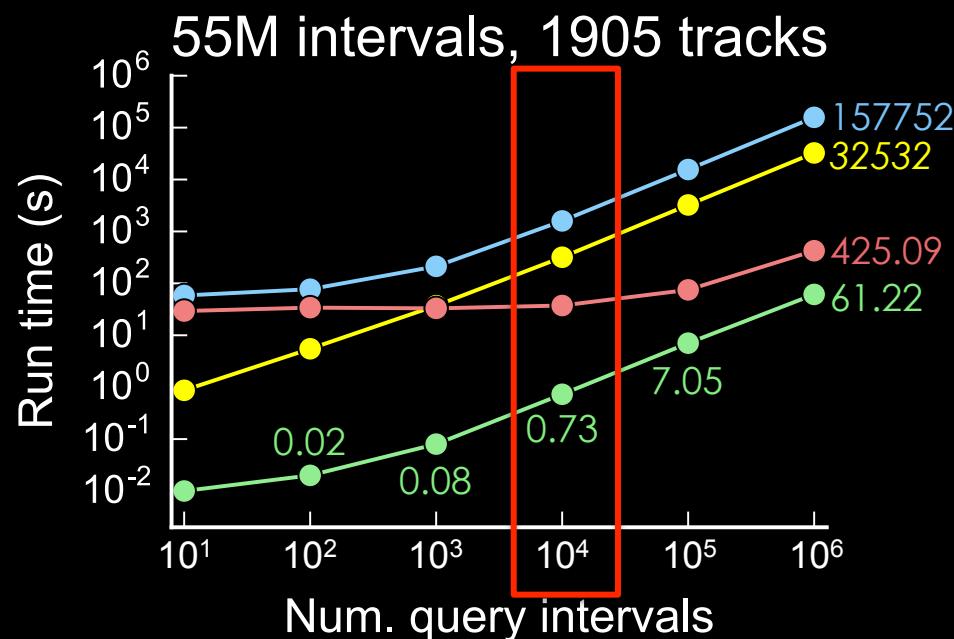
127 reference epigenomes

15 genomic states [CHROMHMM Ernst 2012]

1905 tracks

2.6GHz Intel Xeon CPU, 20MB cache

SQLITE3 w/ UCSC binning
TABIX
BEDTOOLS sorted
GIGGLE



2576X 531X 6.9X



my_chipseq.bed



All **Encode** 1KG GTEx ExAC TOPMed UCSC

DNase-seq of K562

Homo sapiens, adult 53 year
Lab: John Stamatoyannopoulos, UW
Project: ENCODE

ChIP-seq of forebrain

Mus musculus, embryonic 10.5 day
Target: H3K9me3
Lab: Bing Ren, UCSD
Project: ENCODE

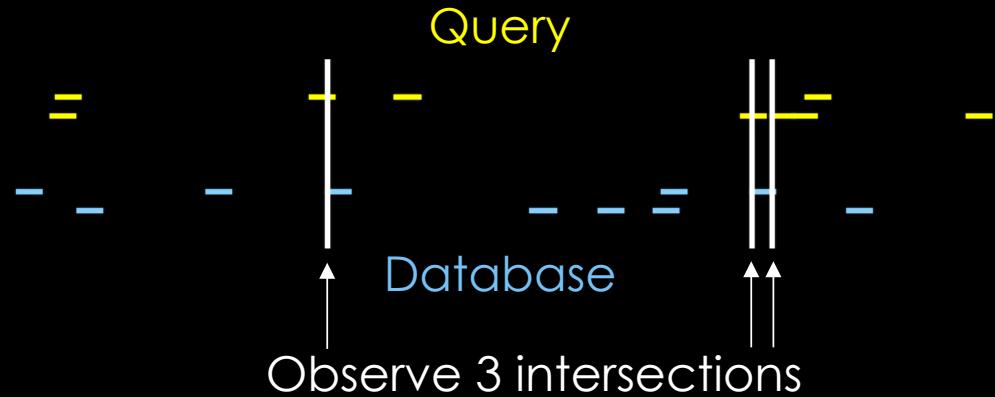
ChIP-seq of heart

Mus musculus, embryonic 10.5 day
Target: H3K27me3
Lab: Bing Ren, UCSD
Project: ENCODE

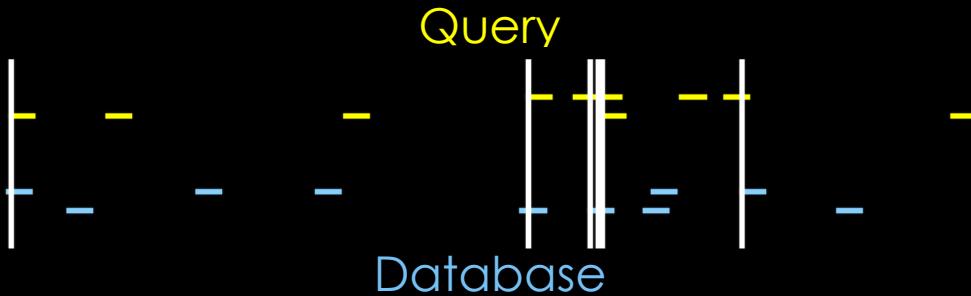
ChIP-seq of embryonic facial prominence

Mus musculus, embryonic 10.5 day
Target: H3K27me3
Lab: Bing Ren, UCSD
Project: ENCODE

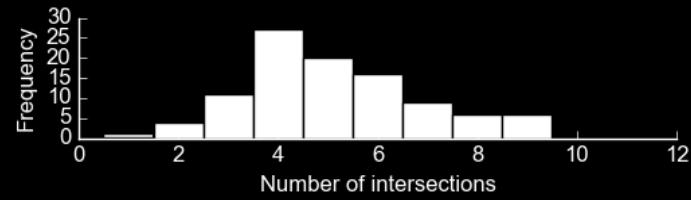
Monte Carlo simulations



Monte Carlo simulations

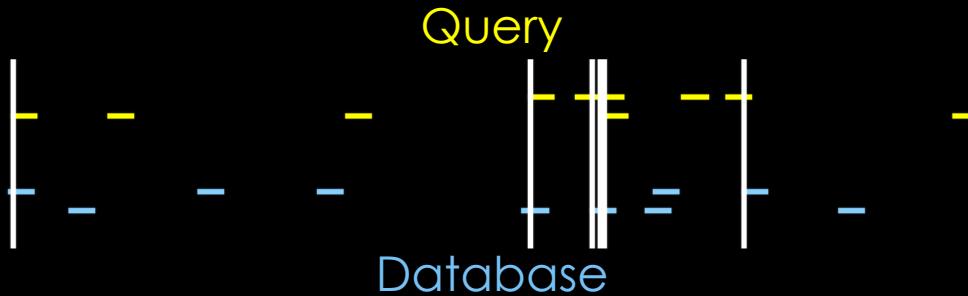


Observe 3 intersections

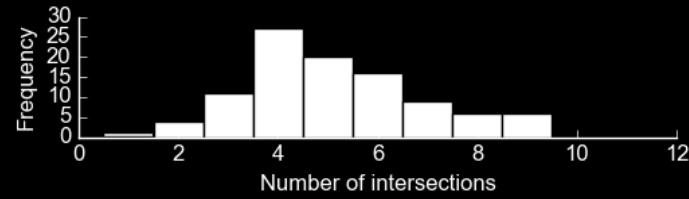


Expect ~4
0.75 "enrichment"

Monte Carlo simulations



Observe 3 intersections



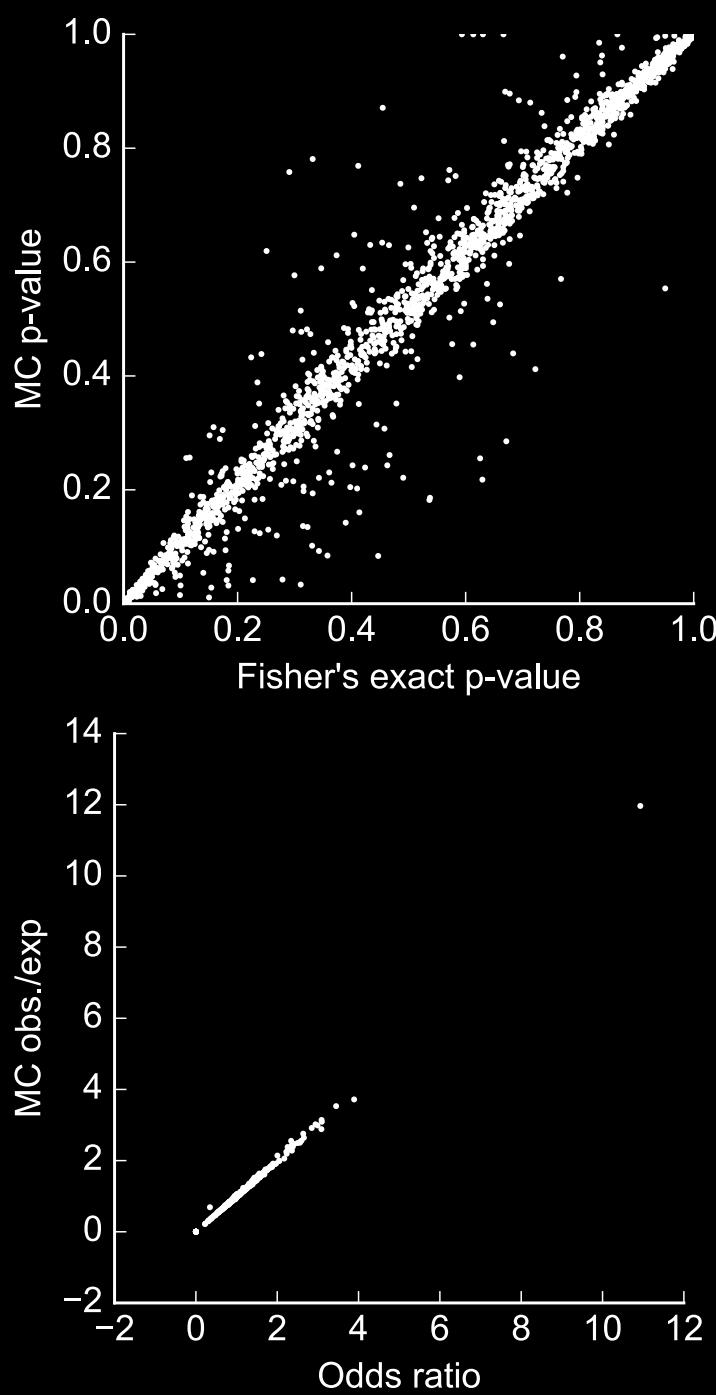
Expect ~4
0.75 "enrichment"

Contingency table

	In target	Not in target
In query	$ Q \cap T $	$ Q - Q \cap T $
Not in query	$ T - Q \cap T $	$\frac{\text{Genome size}}{\mu(Q) + \mu(T)} - (Q + T - Q \cap T)$

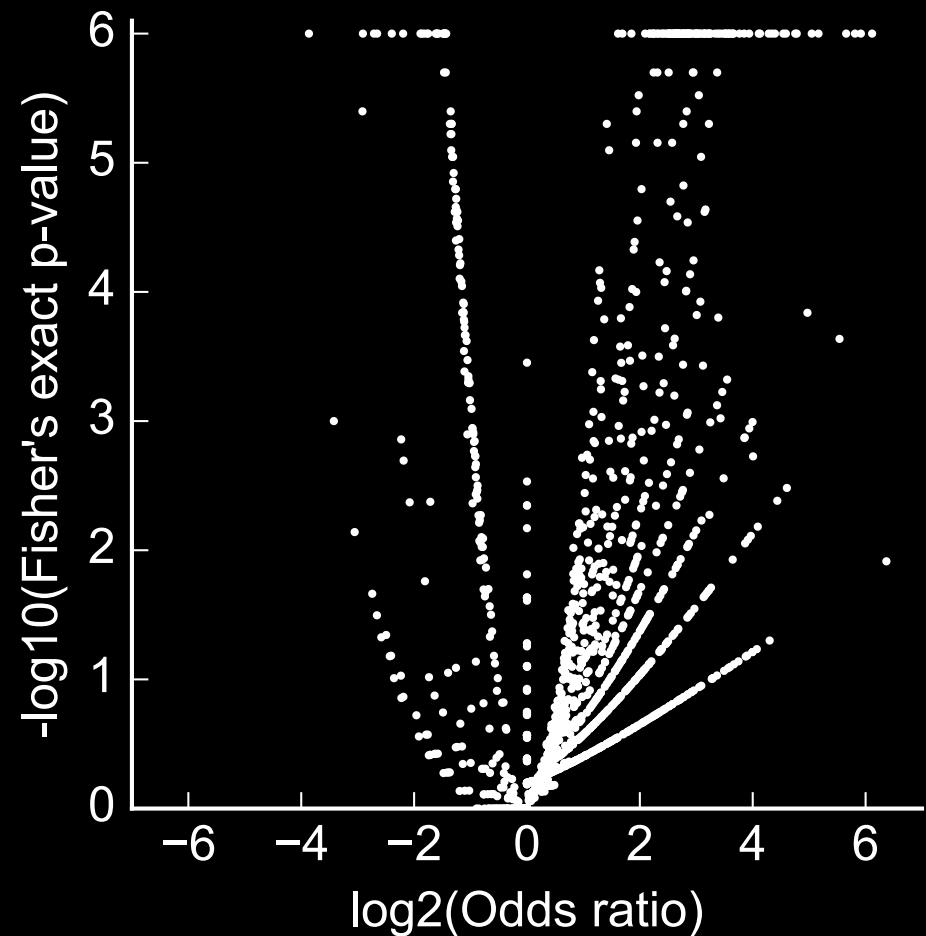


Brent Pedersen



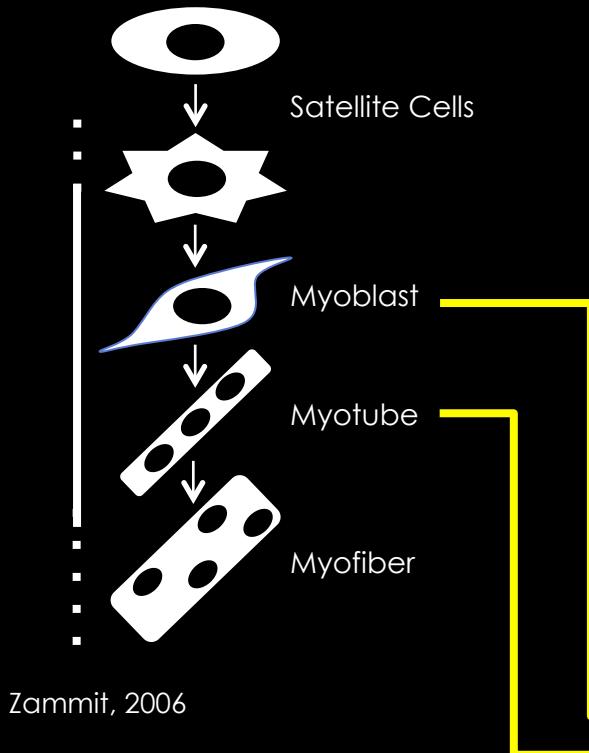
False Positive:
Small fold change, Small variance

False Negative:
Large fold change, Large variance



A novel significance score for gene selection and ranking, Xiao 2014

MyoD

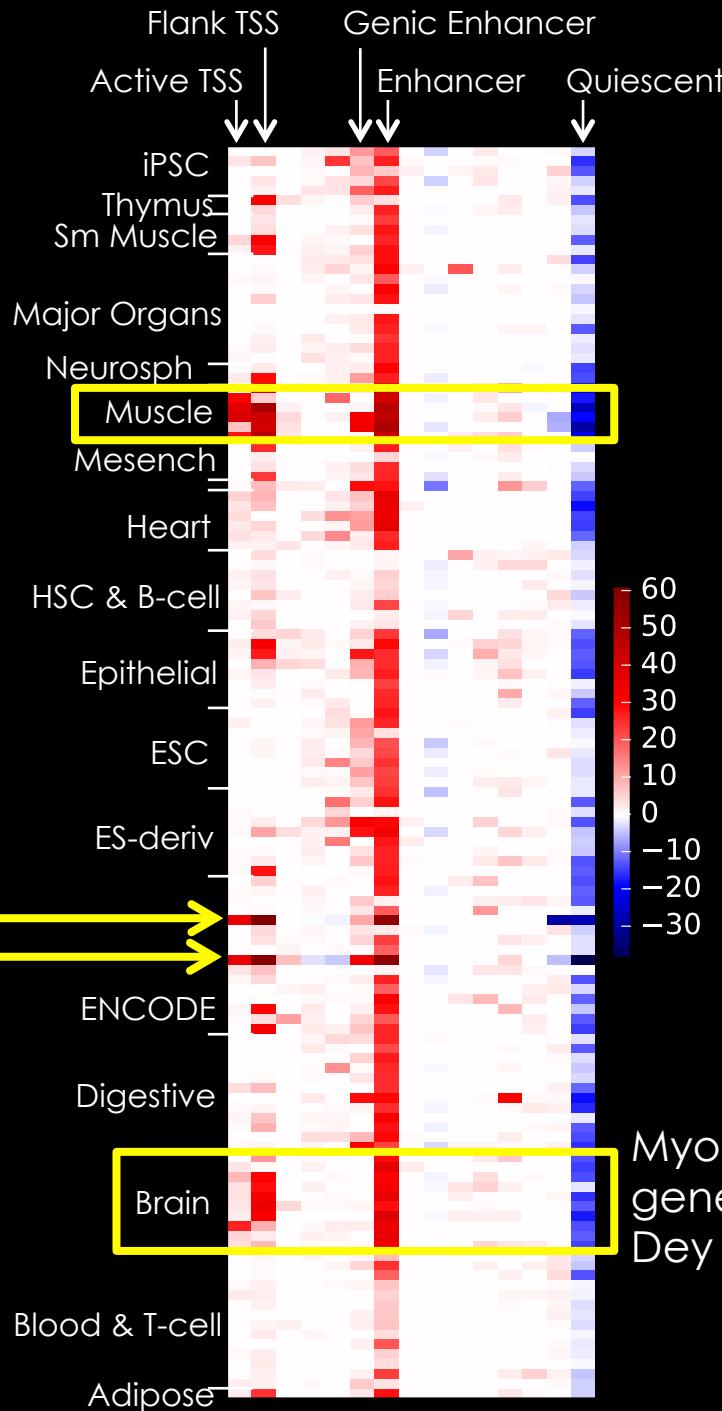


ChIP-Seq MacQuarrie, 2013

+



127 tissues
15 genomic states





Centers for Common
Disease Genomics

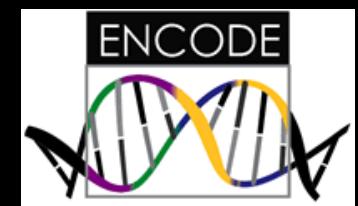
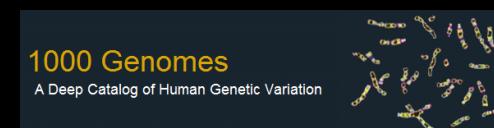
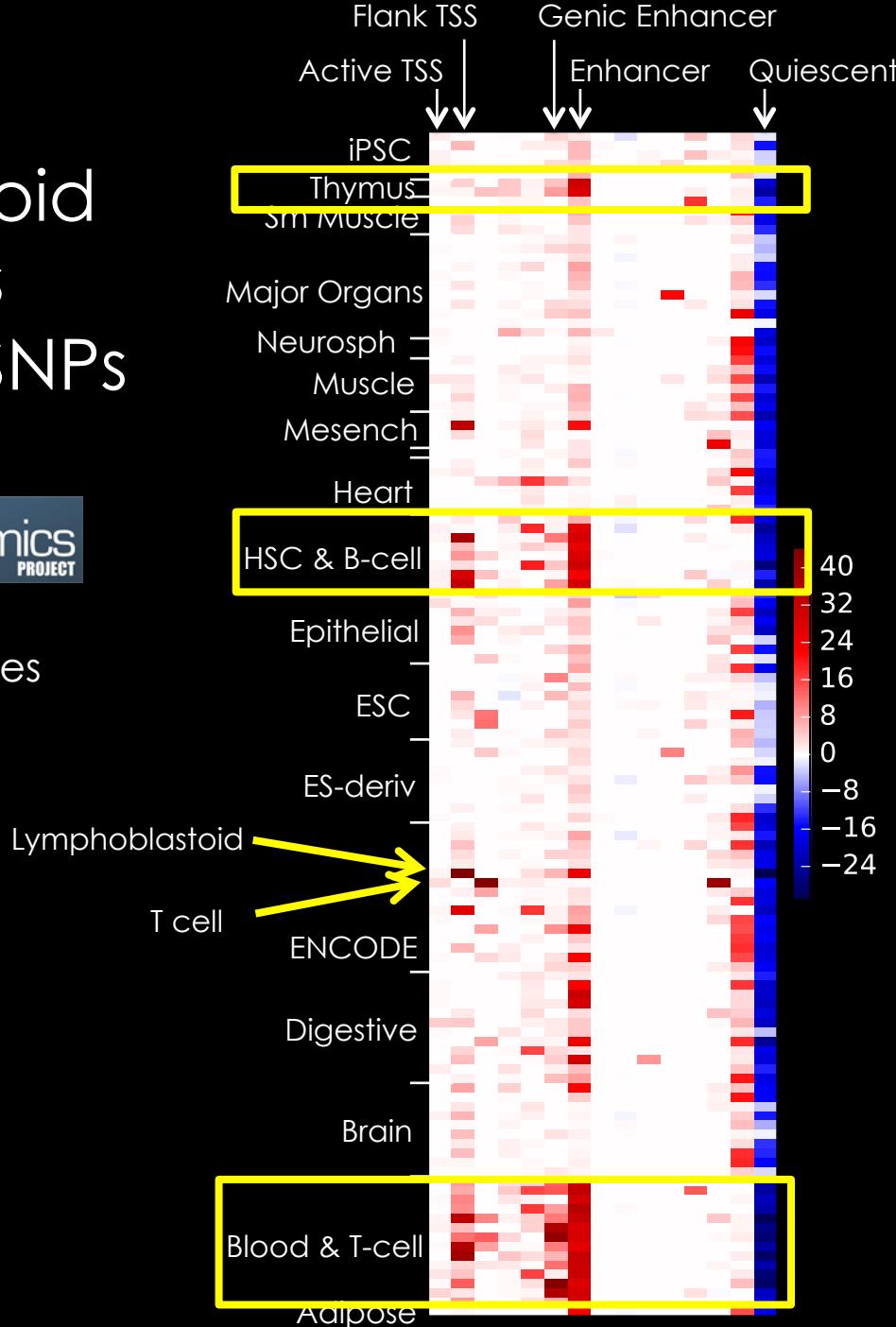


Rheumatoid Arthritis 84 GWAS SNPs

+

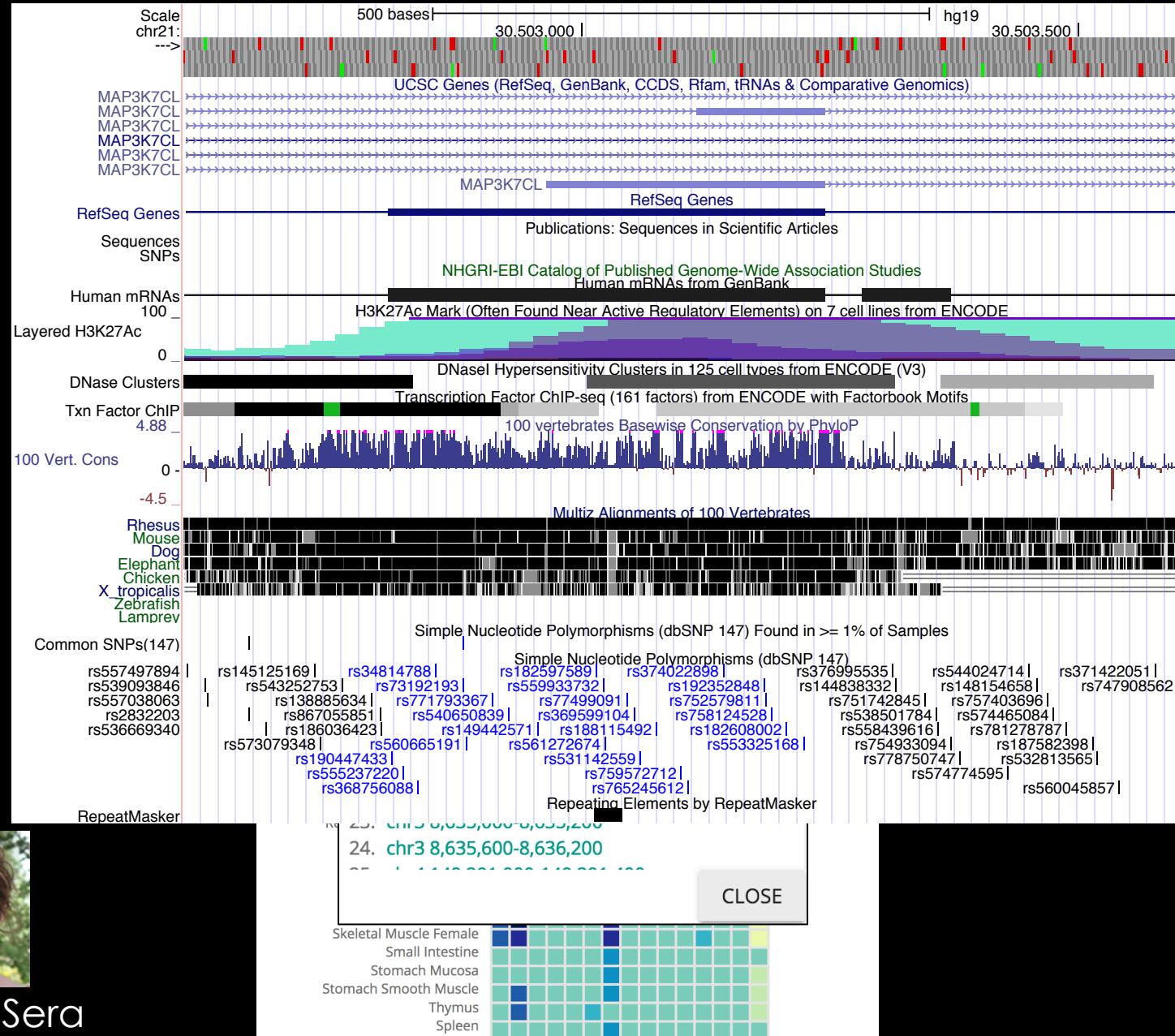


127 tissues
15 genomic states



ClinVar

GGTATGGGGCCAAGAGATATATCT
CGGCTGTCATCACTTAGACCTCAC
TGGGCATAAAAGTCAGGGCAGAGC
TGCATCTGACTCCTGAGGAGAAGT
TGGTATCAAGGTTACAAGACAGGT
GACTCTCTCTGCCTATTGGTCTAT



Tonya Di Sera



GIGGLE

github.com/ryanlayer/giggle

Thousands of Annotations

100 Millions of Intervals

Fast Queries

High Quality Enrichment Score

Reference:

AGCGTGCATTGCCGCGGTCCAACGTGTCAATAGACAACATTCTCGTCGAAATACACAGTTCTGGGAGTGATGTGGCTAGGGAATCACTG



CGTGCATTG
ATTGGCGCG

GOATTGCCCG
ATTGCCCGCG

GOATTGCCCG
ATTGCCCGCGG

GENOTYPE QUERY TOOLS



C C

G A

C C

C A

G C

CATTCTCGTCG
ACAACATTCTC

TGGGAGTGAT
TTTCTGGGA

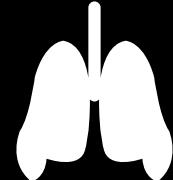
GATGTGGCT

A
TTTCTGGGA



CAATAGACAACATT

TCGAAATACACA



SFARI

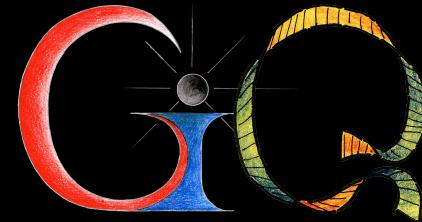
gatk

70M SNVs

519 

lumpy

100K SVs



UCSC

Not in problematic regions (SEGDPUs)

ExAC ← 1000 Genomes A Deep Catalog of Human Genetic Variation ← ClinVar ← OMIM ← Not previously observed

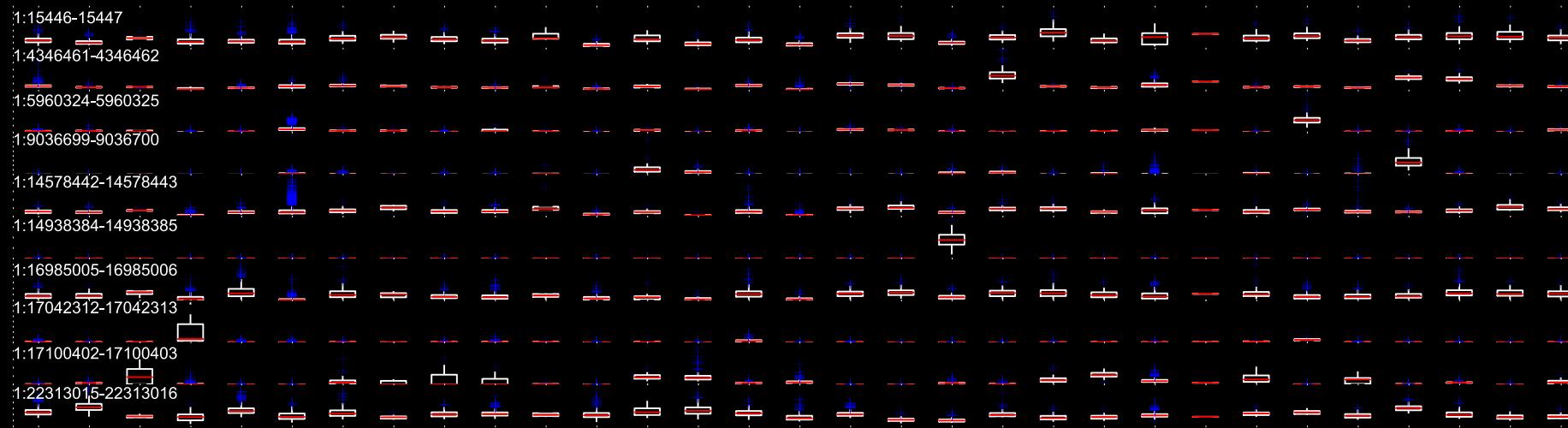
 GTEx

de novo in affected

Gene expression

Tissues

RPKM



SFARI

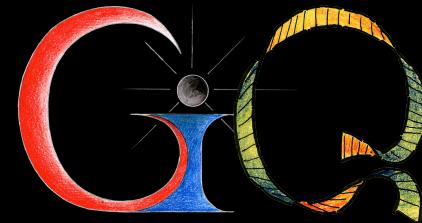
gatk

70M SNVs

519

lumpy

100K SVs



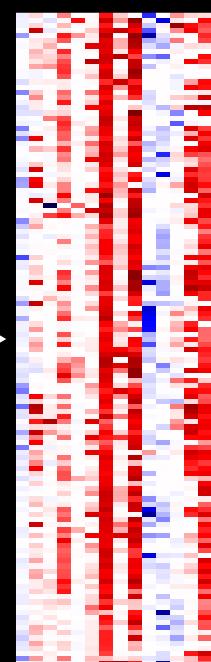
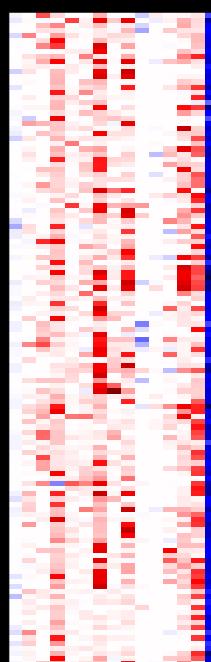
UCSC

Not in problematic regions (SEG DUPs)

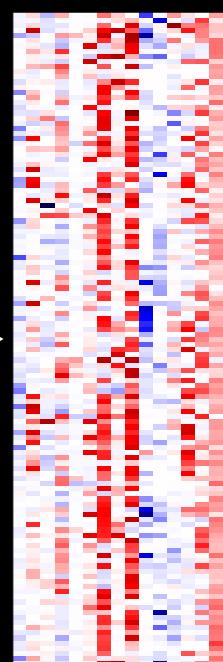
ExAC ← 1000 Genomes A Deep Catalog of Human Genetic Variation ← ClinVar ← OMIM ← Not previously observed

de novo in
affected

de novo in
unaffected



difference



SFARI

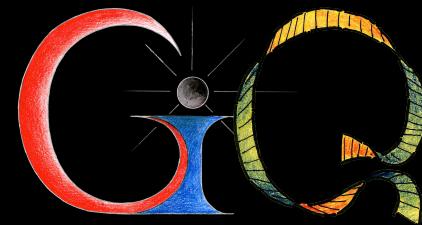
gatk

70M SNVs

519

lumpy

100K SVs



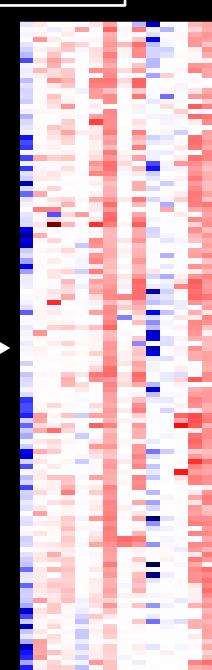
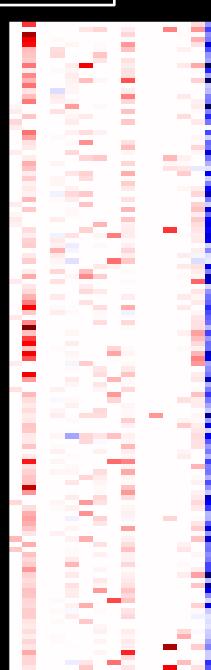
UCSC

Not in problematic regions (SEG DUPs)

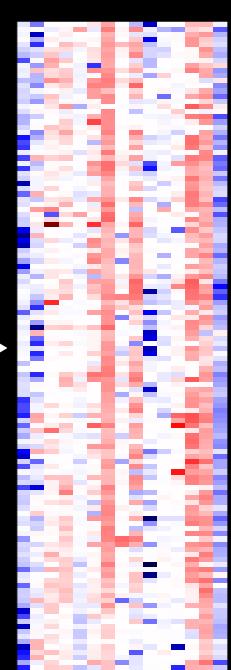
ExAC ← 1000 Genomes ← ClinVar ← OMIM ← Not previously observed

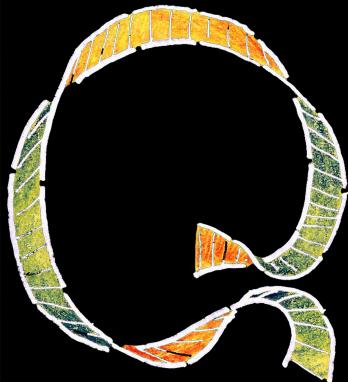
de novo in
affected AFR

de novo in
unaffected AFR

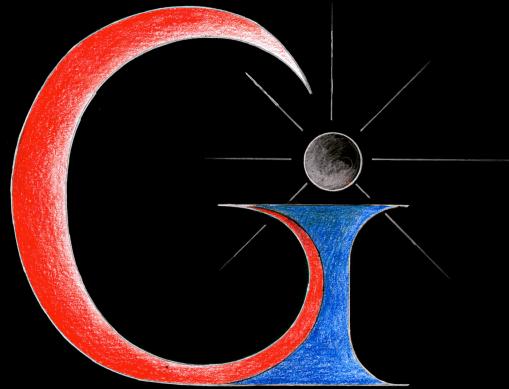


difference





&



GENOTYPE QUERY TOOLS

GIGGLE

github.com/ryanlayer/giggle

Questions?

Aaron Quinlan, Brent

layer@colorado.edu, ryanlayer@gmail.com

@ryanlayer

Aaron Quinlan, Brent



USTAR Center for
Genetic Discovery