



github.com/ryanlayer/giggle

Search **ALL** Genome Annotations

Ryan M. Layer

ryan.layer@gmail.com

@ryanlayer

University of Utah

quinlanlab.org

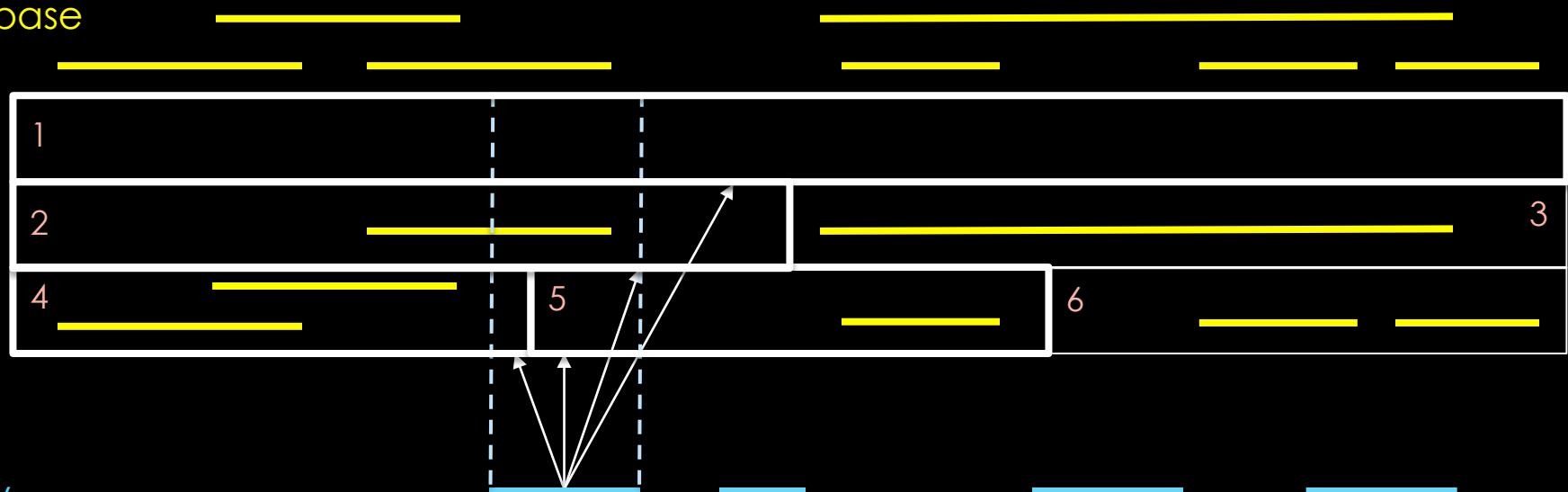
```
graph TD; GFF[GFF] --- NARROWPEAK[NARROWPEAK]; GFF --- VCF[VCF]; BROADPEAK[BROADPEAK] --- BEDPE[BEDPE]; CRAM[CRAM] --- BED[BED]; CRAM --- SAM[SAM]; CRAM --- BCF[BCF]; BED --- SAM; BED --- BCF; BAM[BAM] --- BIGWIG[BIGWIG]; WIG[WIG] --- BIGBED[BIGBED]; BIGWIG --- BEDGRAPH[BEDGRAPH];
```



Hierarchical Binning

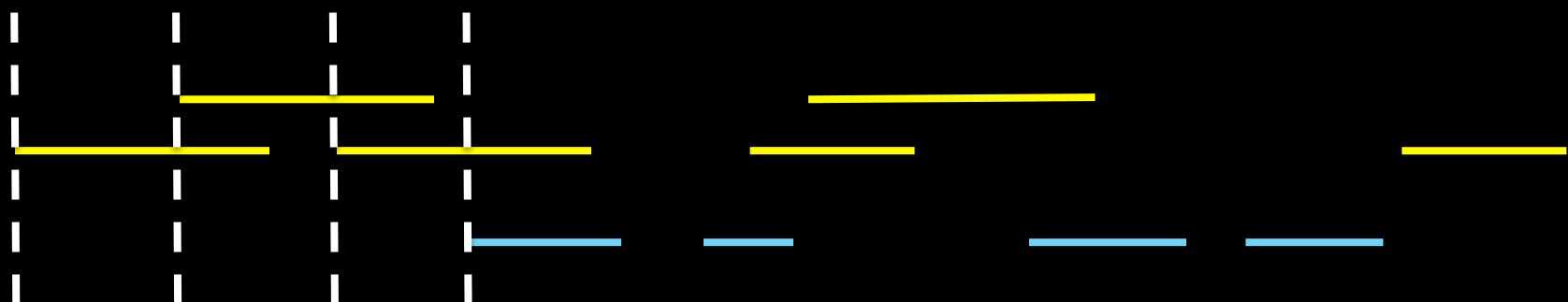
[UCSC Kent 2002, TABIX Li 2011, BEDTOOLS Quinlan 2010]

Database

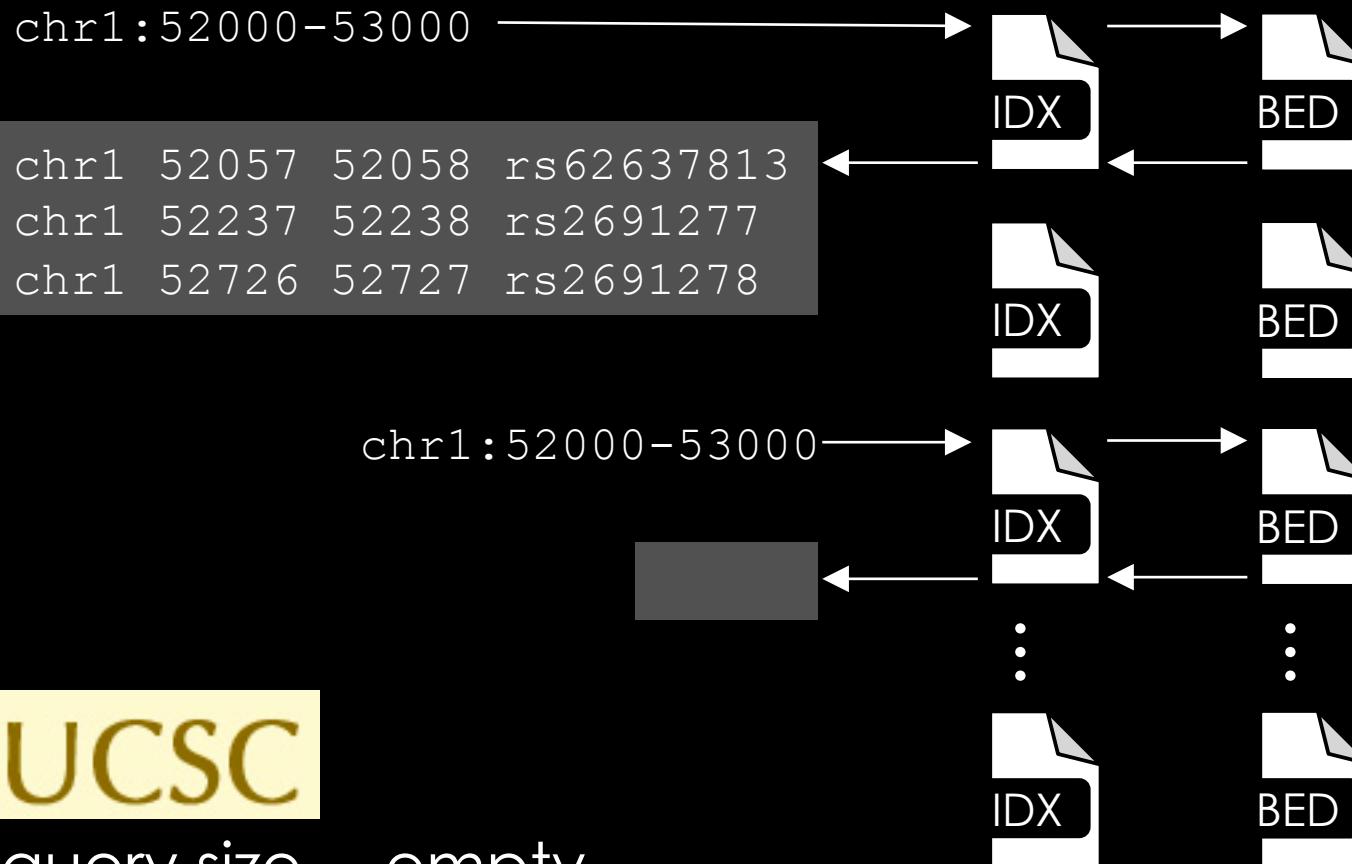


Query

Sweep [BEDTOOLS]



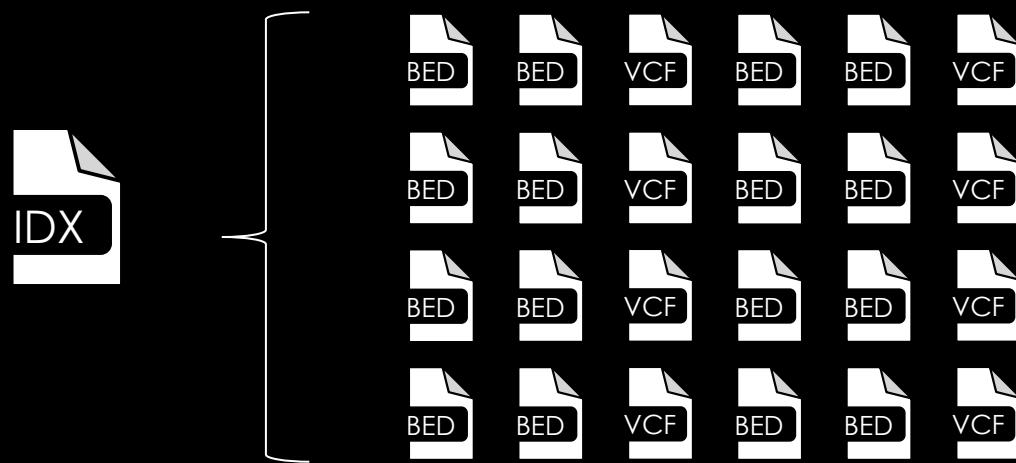
TABIX



<u>query size</u>	<u>empty</u>
100bp	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

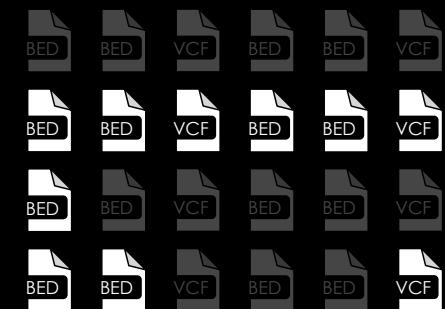
2

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

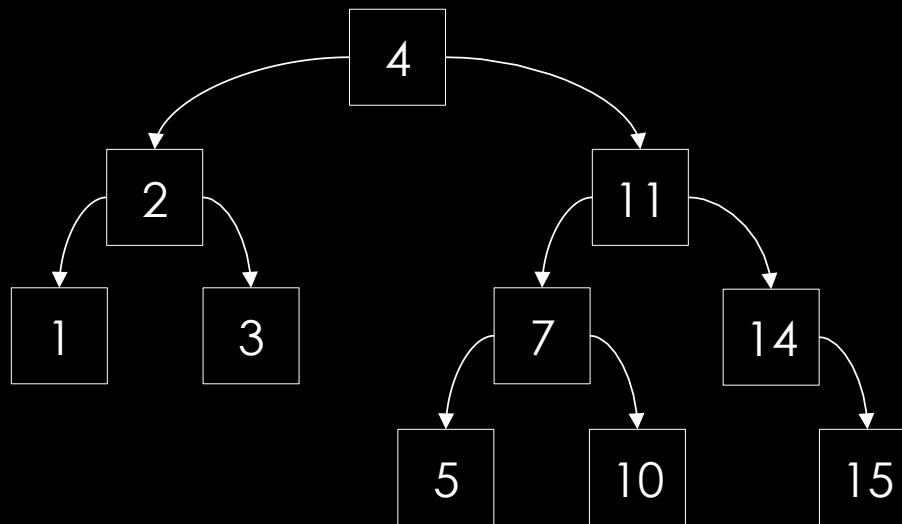
In this paper, we describe a new indexing technique, the *time index*, for improving the performance of certain classes of temporal queries. The time index can be used to retrieve versions of objects that are valid during a specific time period. It supports the processing of the temporal WHEN operator 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.

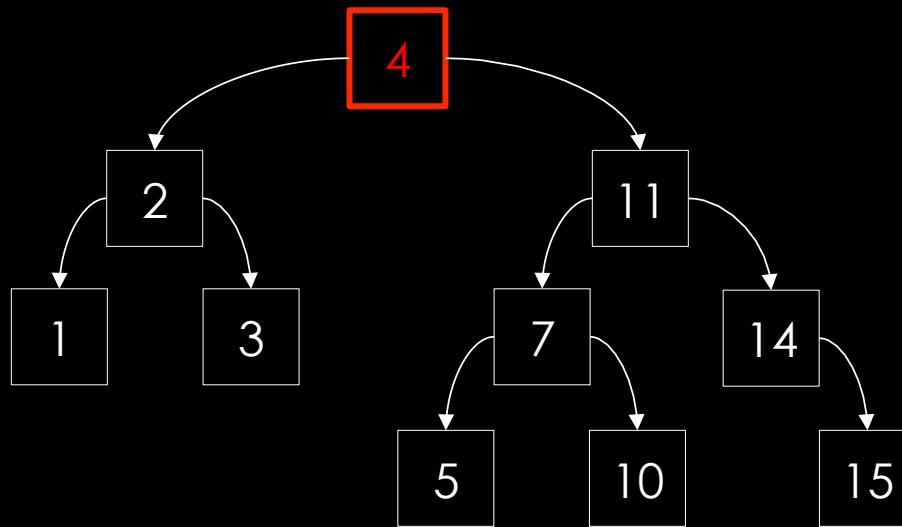
The storage techniques for temporal data proposed in [AS88, Lum84] index or link the versions of each individual object separately. In order to retrieve object versions that are valid during a certain time period, it is necessary to first locate the first (current) version of each object, and then search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

VLDB 1990

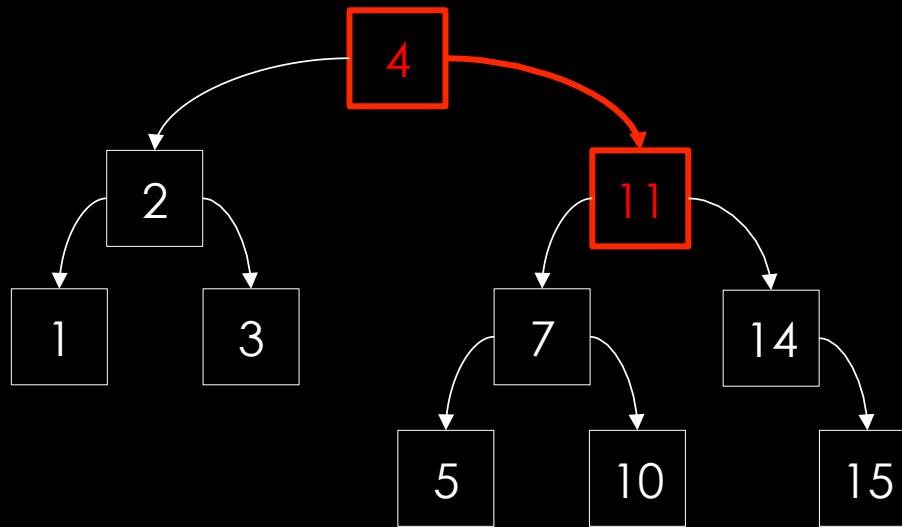
Binary Tree



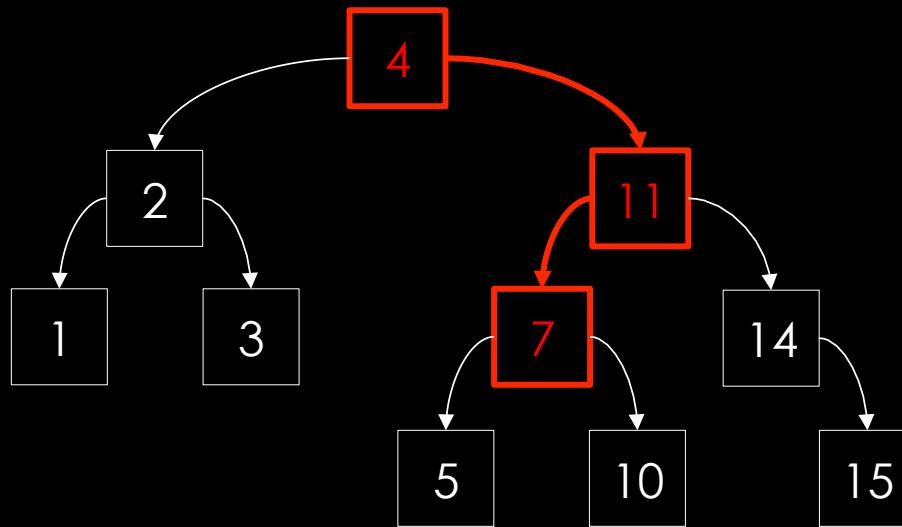
Binary Tree



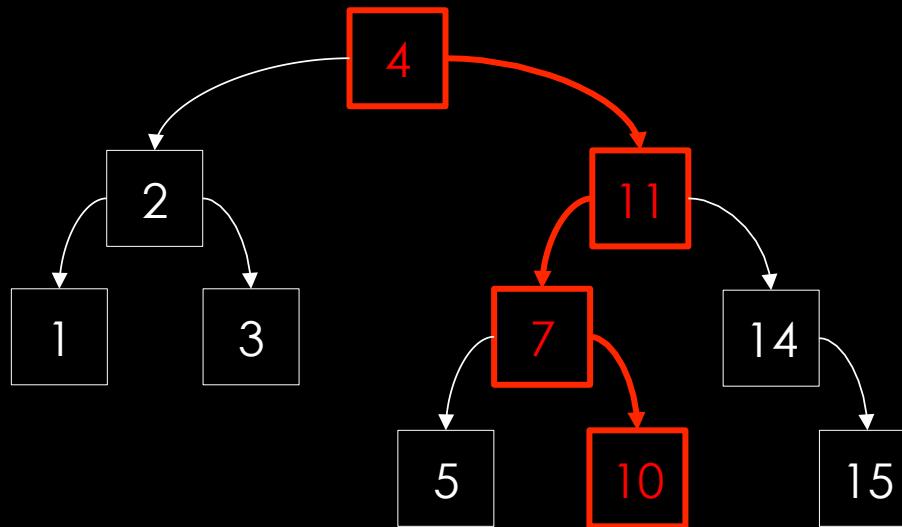
Binary Tree



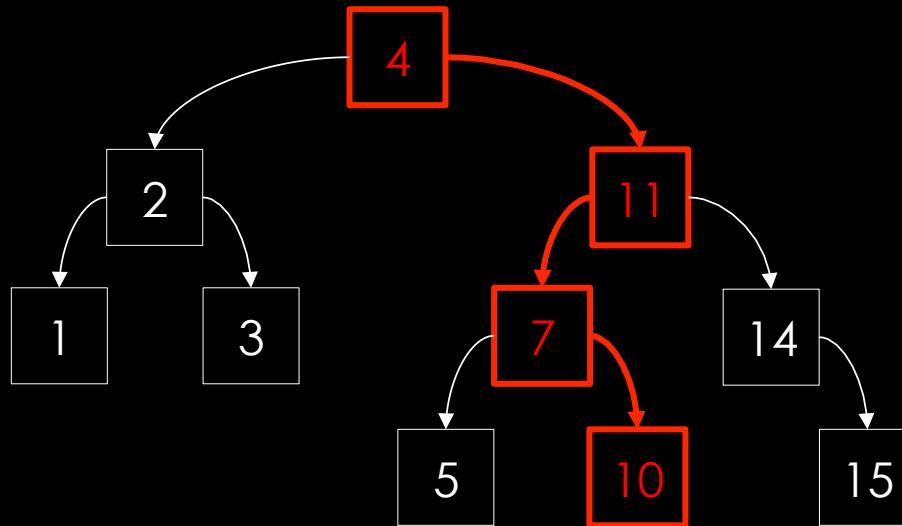
Binary Tree



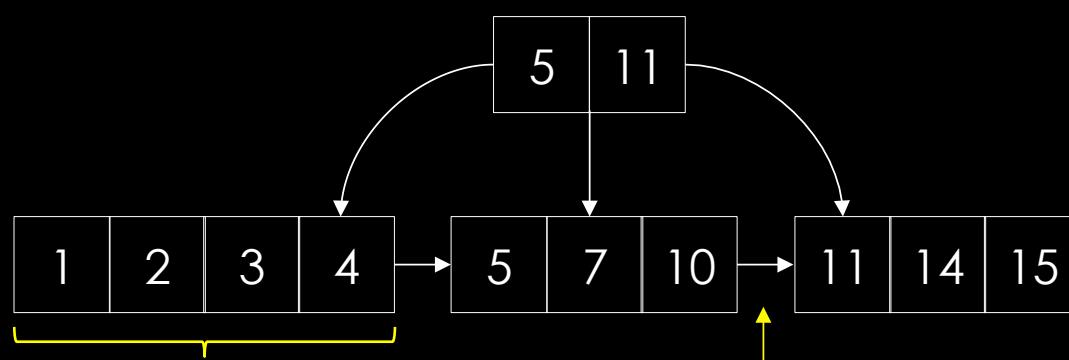
Binary Tree



Binary Tree



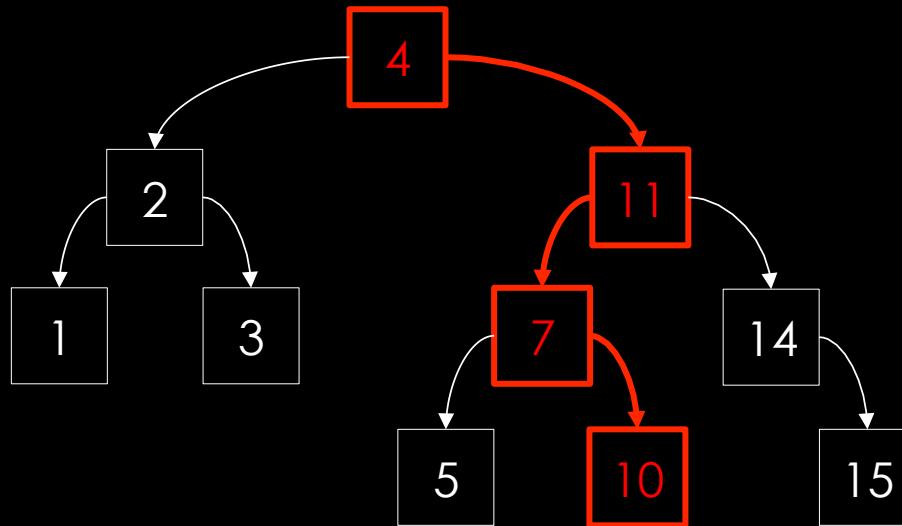
B+ Tree



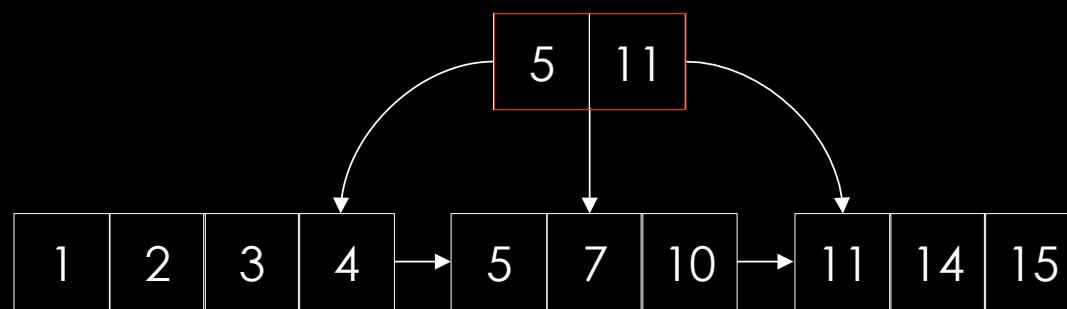
at most N (4) keys per node

linked leaves

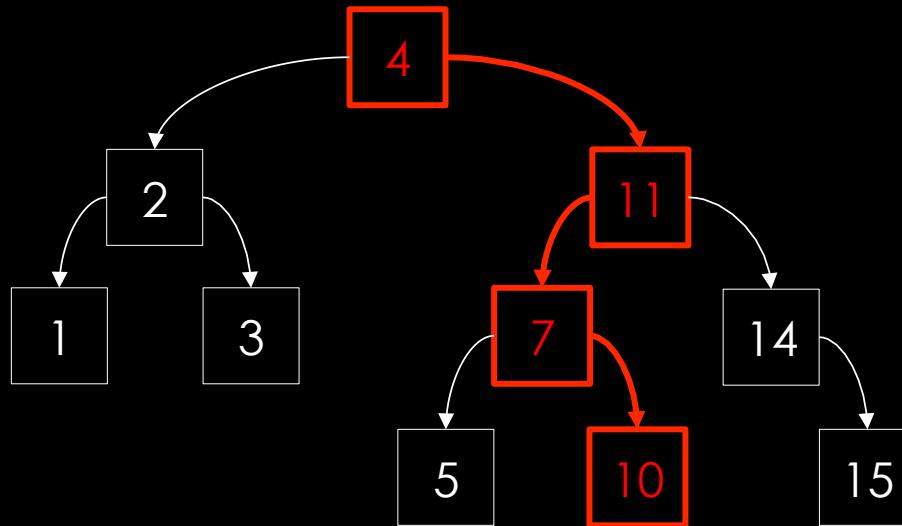
Binary Tree



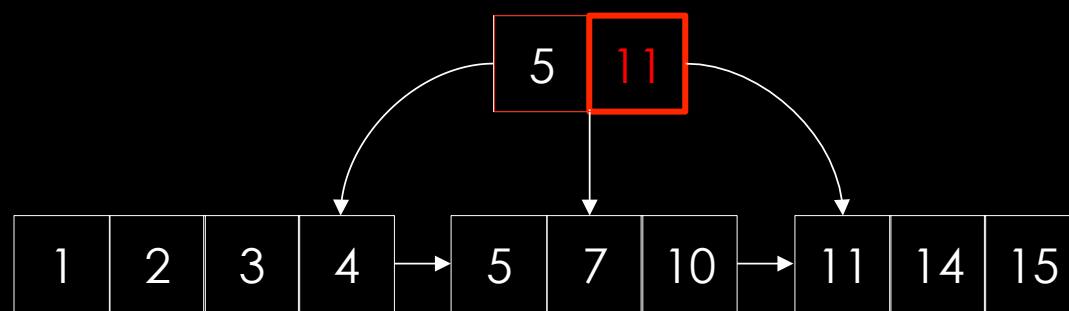
B+ Tree



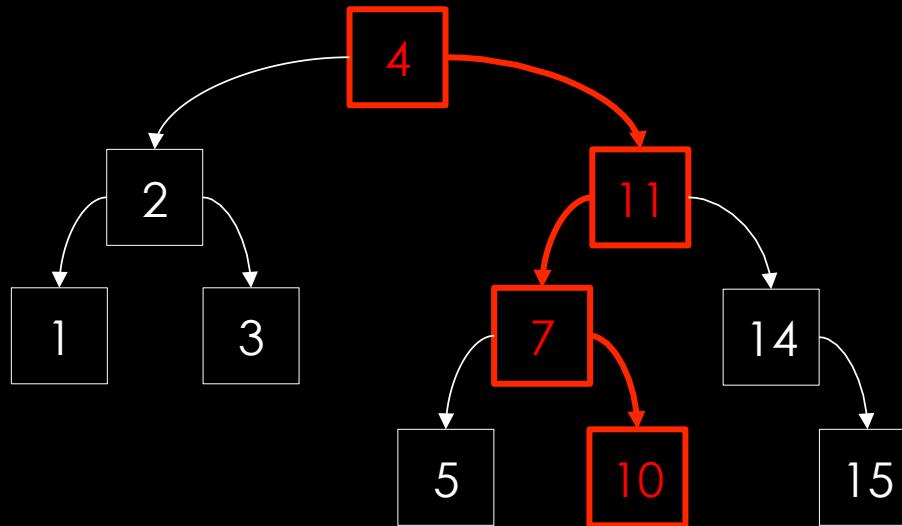
Binary Tree



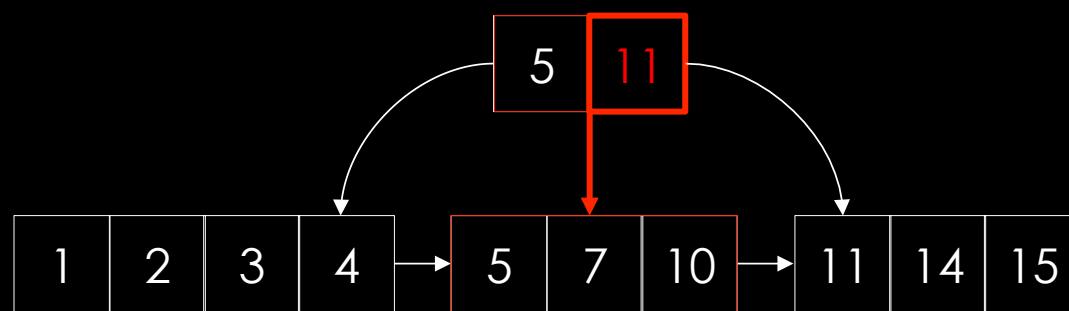
B+ Tree



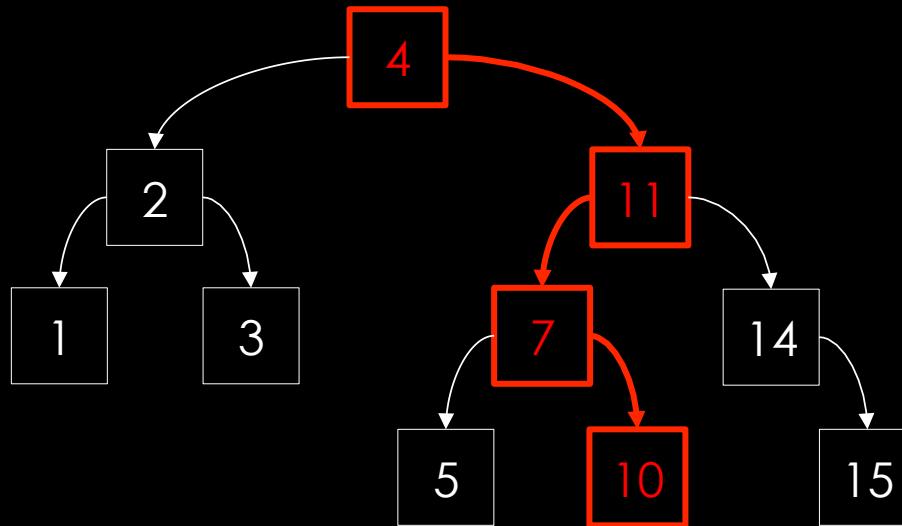
Binary Tree



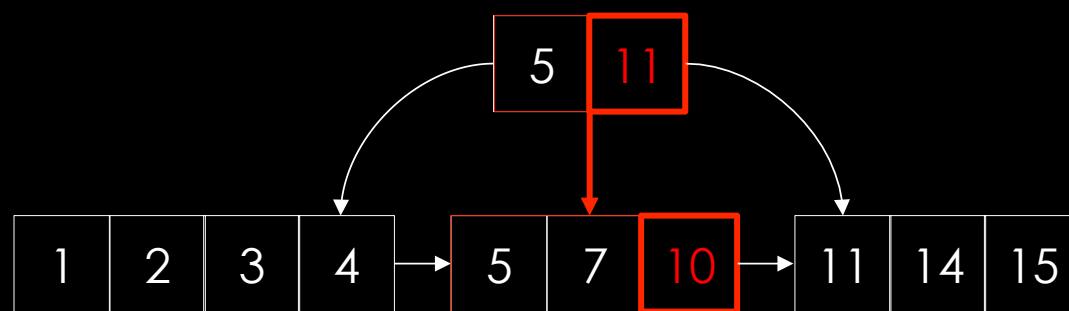
B+ Tree



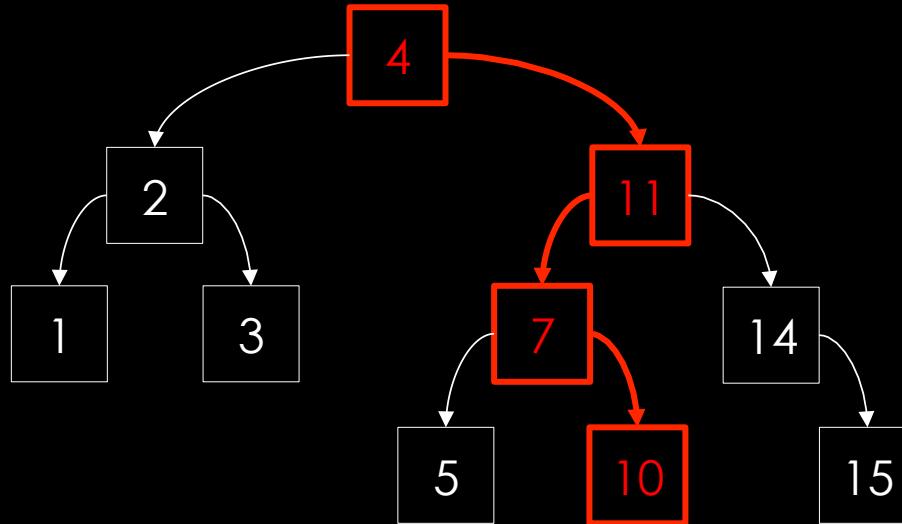
Binary Tree



B+ Tree



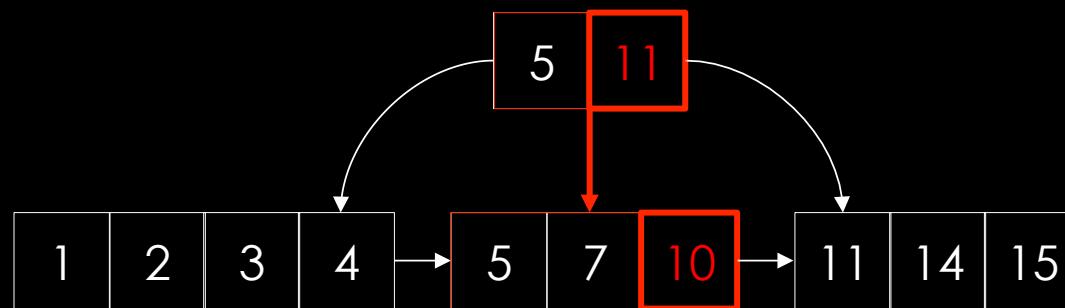
Binary Tree



Disk layout:



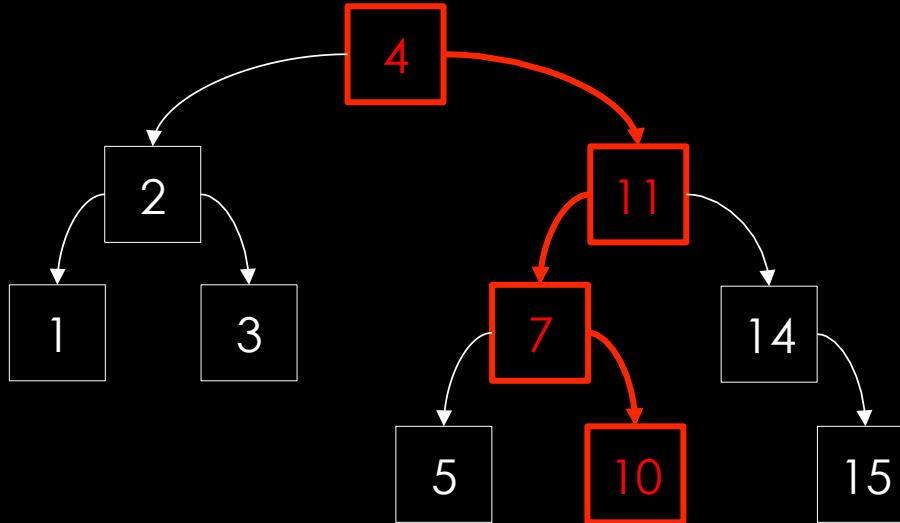
B+ Tree



Disk layout:



Binary Tree



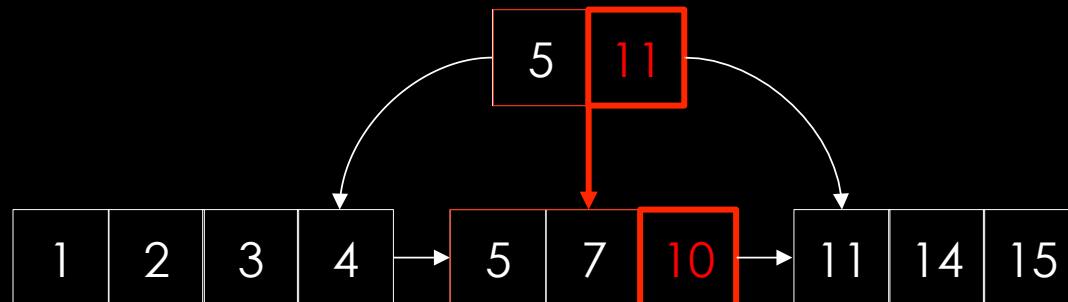
4 disk reads:



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

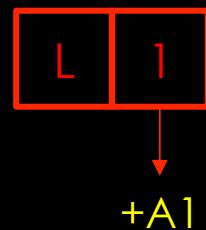
B2

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

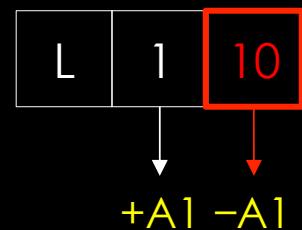
B2

C

C2

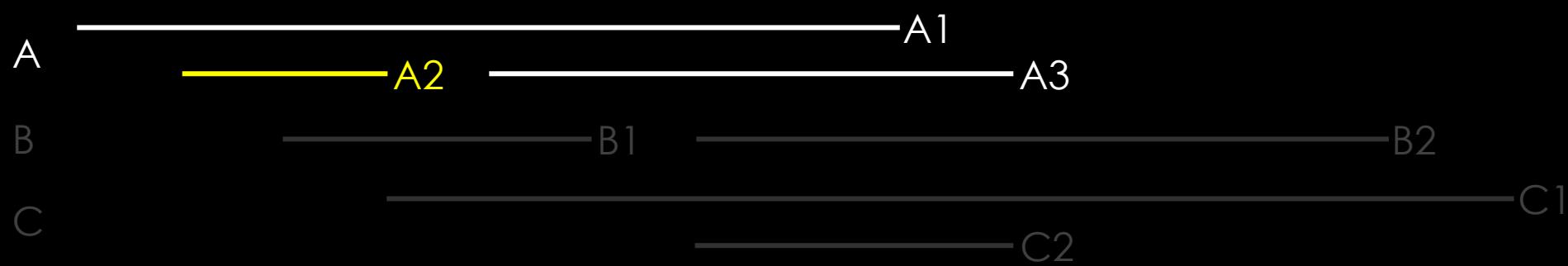
C1

index(A1, 1, 9)

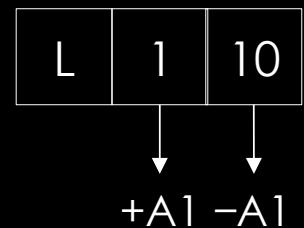


+A1 -A1

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

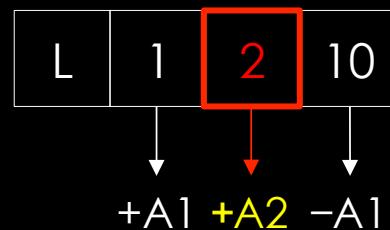
A3

B2

C

C2

index(A2, 2, 4)



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

A

A2

A1

B

B1

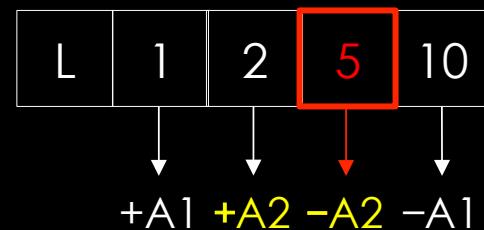
A3

B2

C

C2

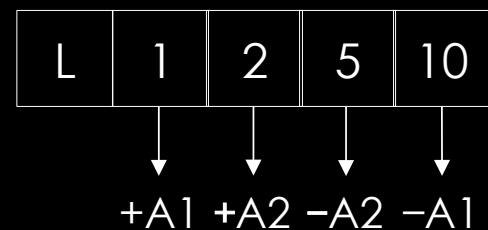
index(A2, 2, 4)



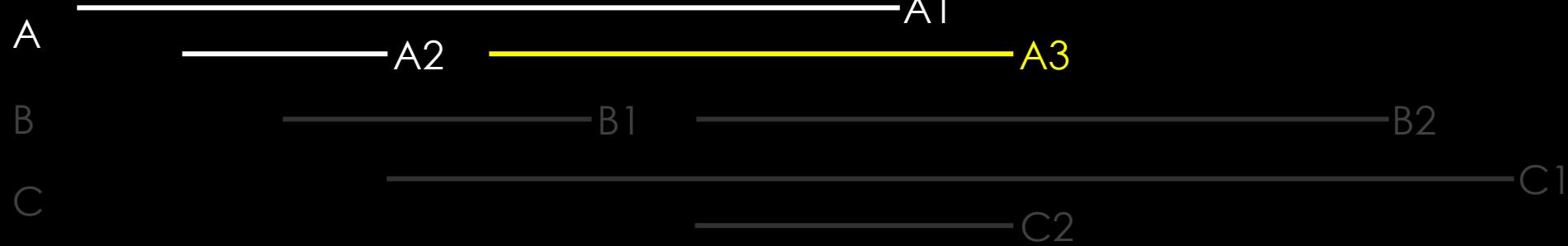
1 2 3 4 5 6 7 8 9 10 11 12 13 14



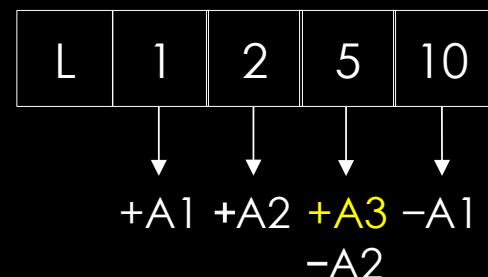
index(A3, 5, 10)



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

— A1

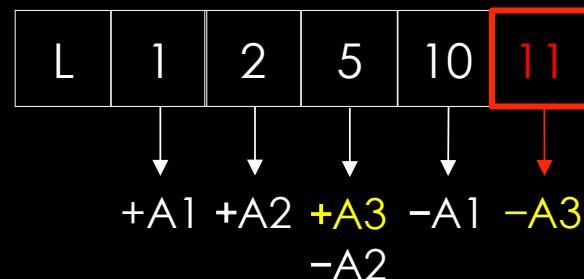
— A2 — A3

— B1 — B2

C

— C1
— C2

index(A3, 5, 10)

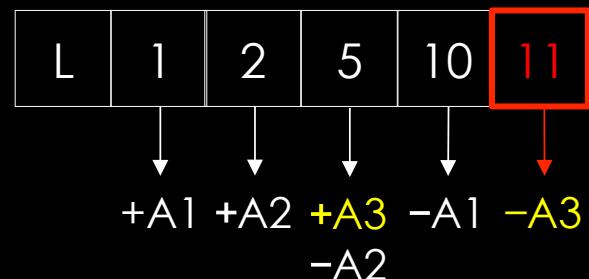


1 2 3 4 5 6 7 8 9 10 11 12 13 14



index(A3, 5, 10)

— Too Many Keys —

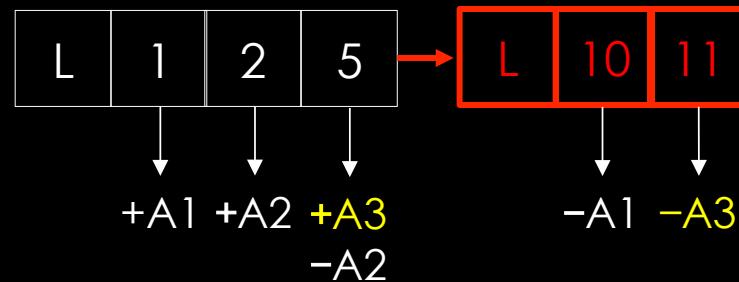


1 2 3 4 5 6 7 8 9 10 11 12 13 14



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

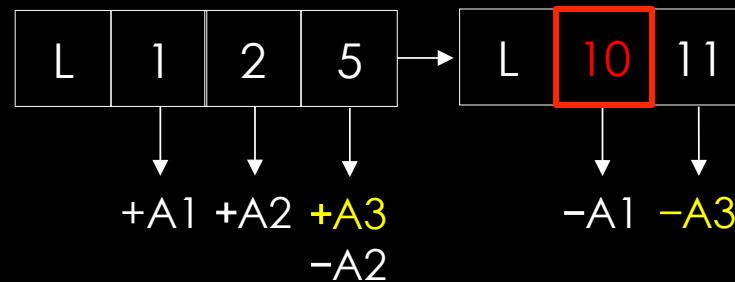
C

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

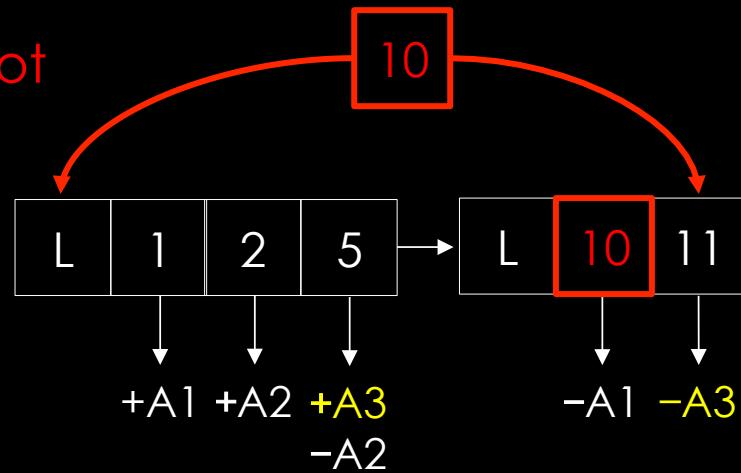
B2

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

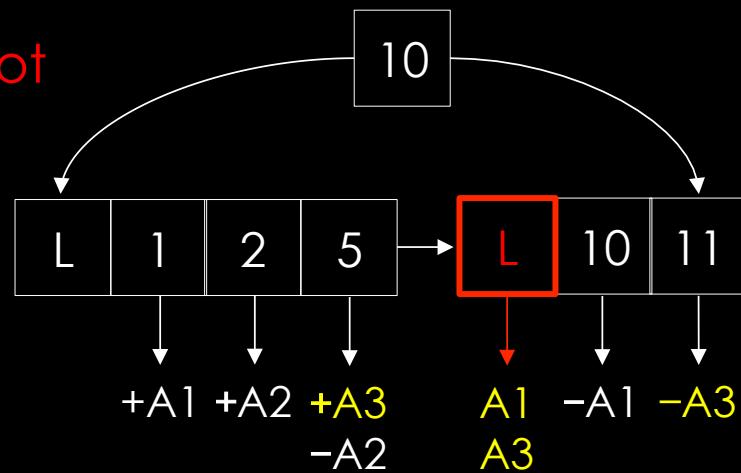
C

C2

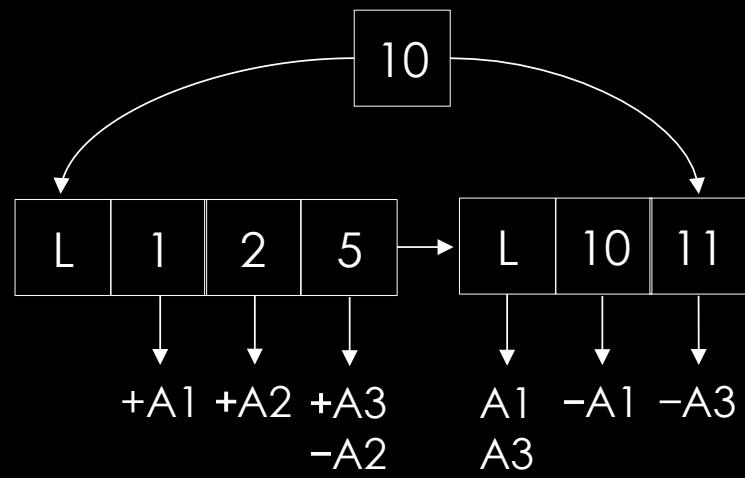
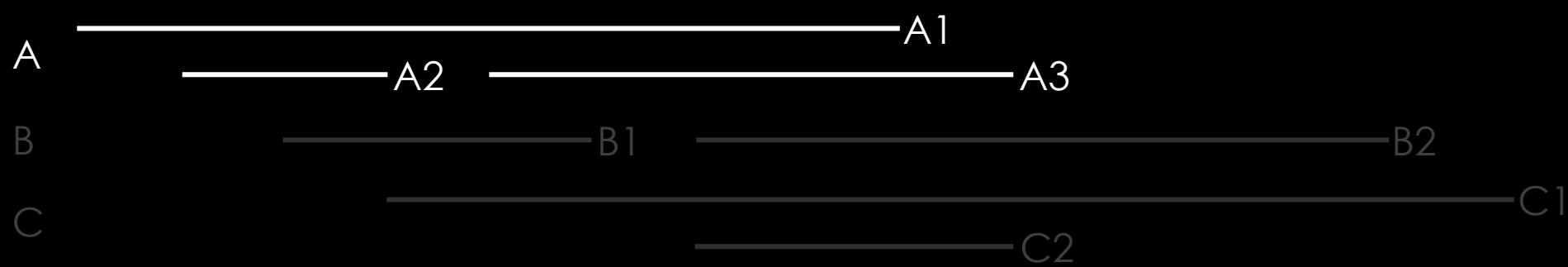
C1

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

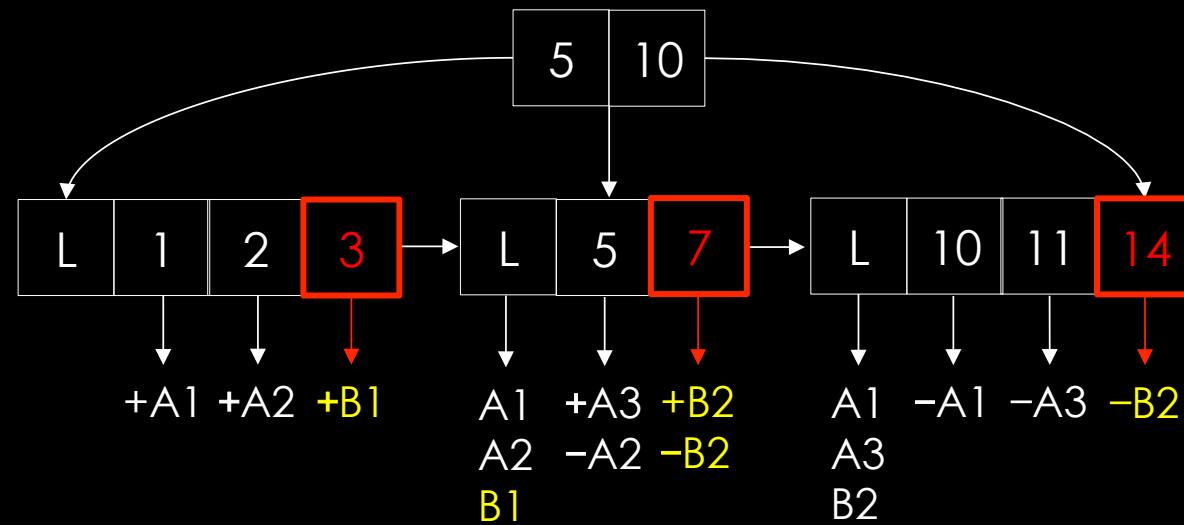
B1

B2

C

C1

C2



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

A

A2

A1

B

B1

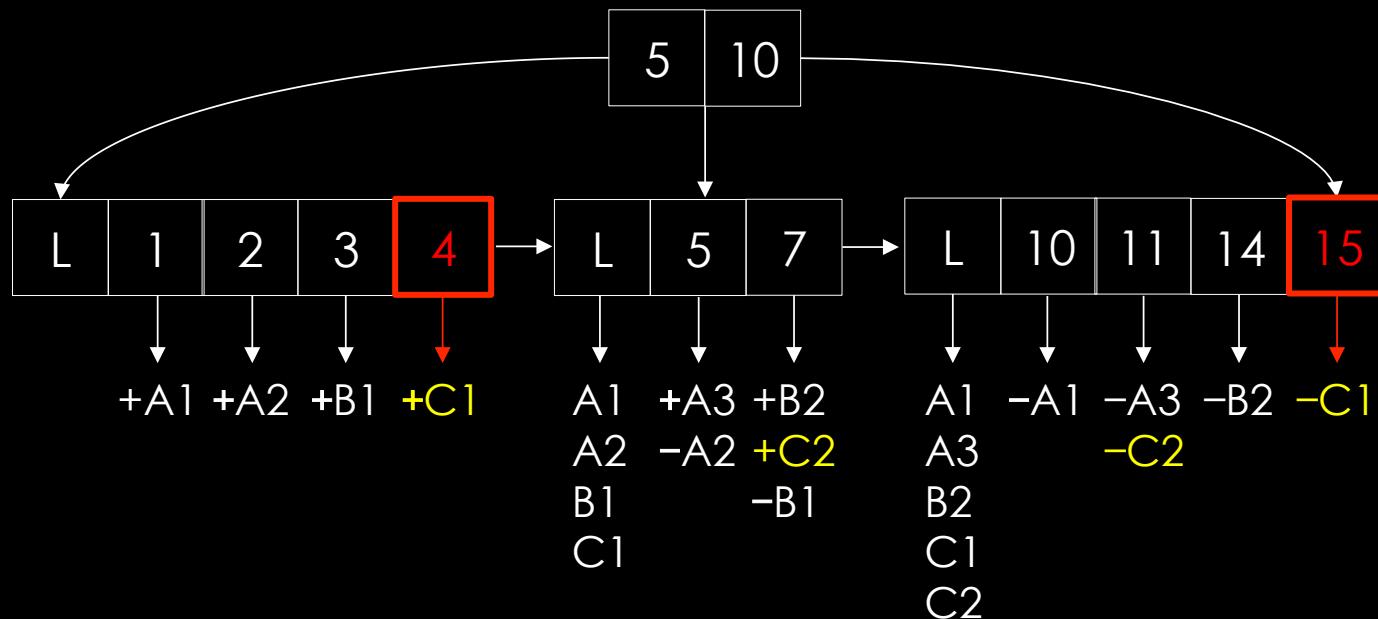
A3

B2

C

C2

C1

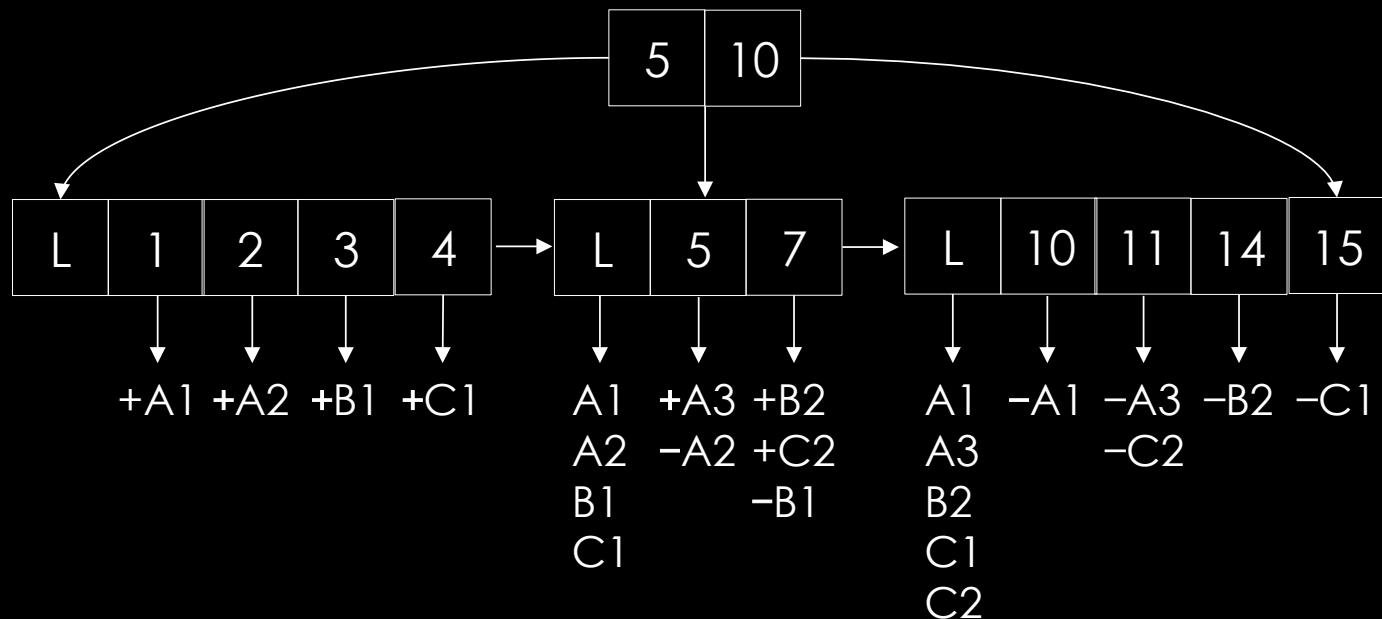


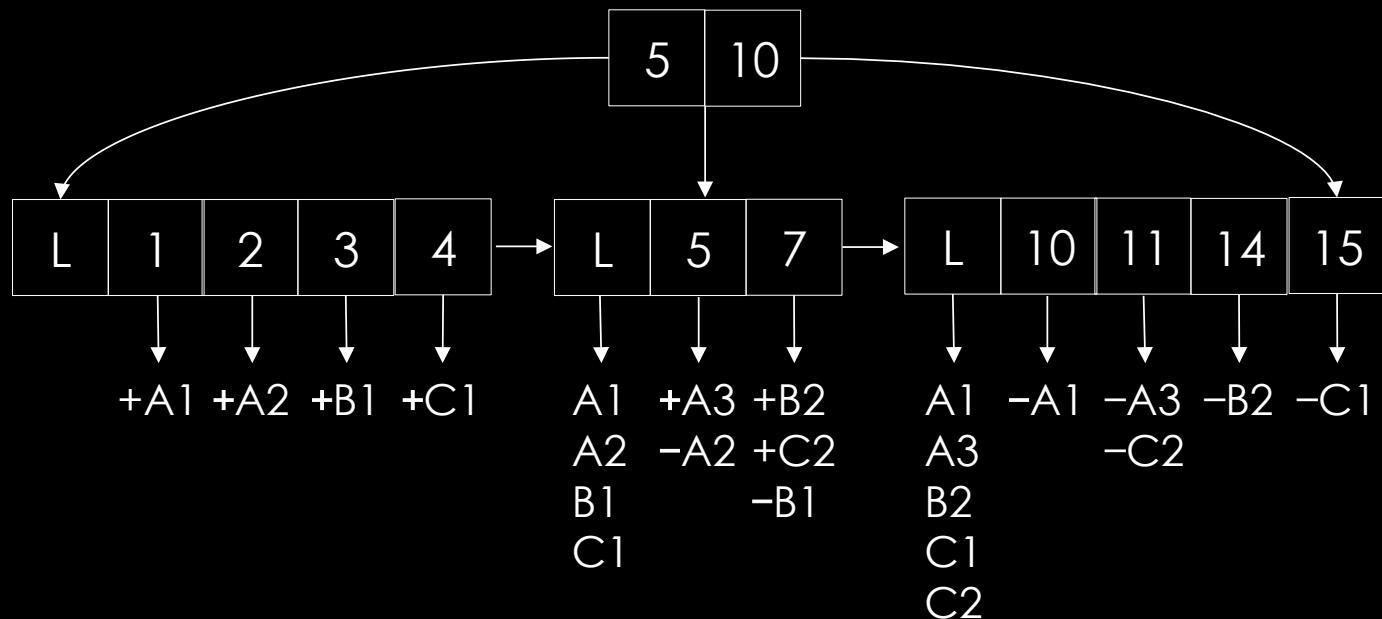
1 2 3 4 5 6 7 8 9 10 11 12 13 14



```
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

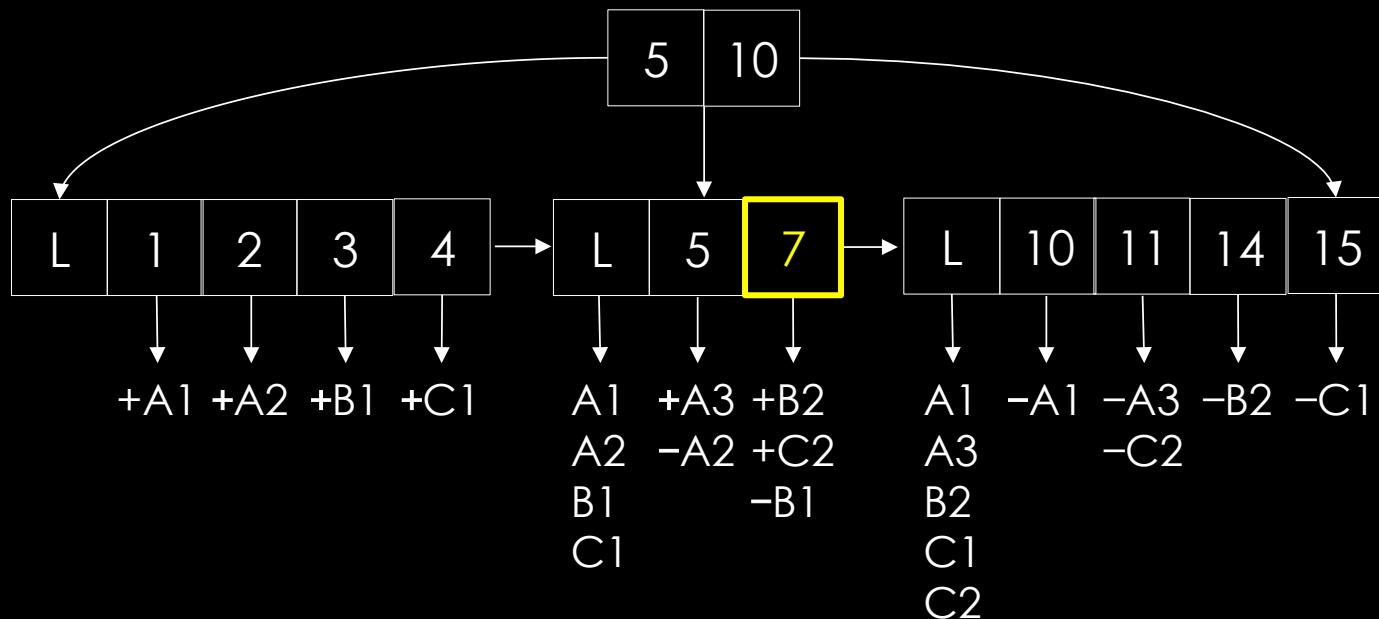
A1
A2
A3

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

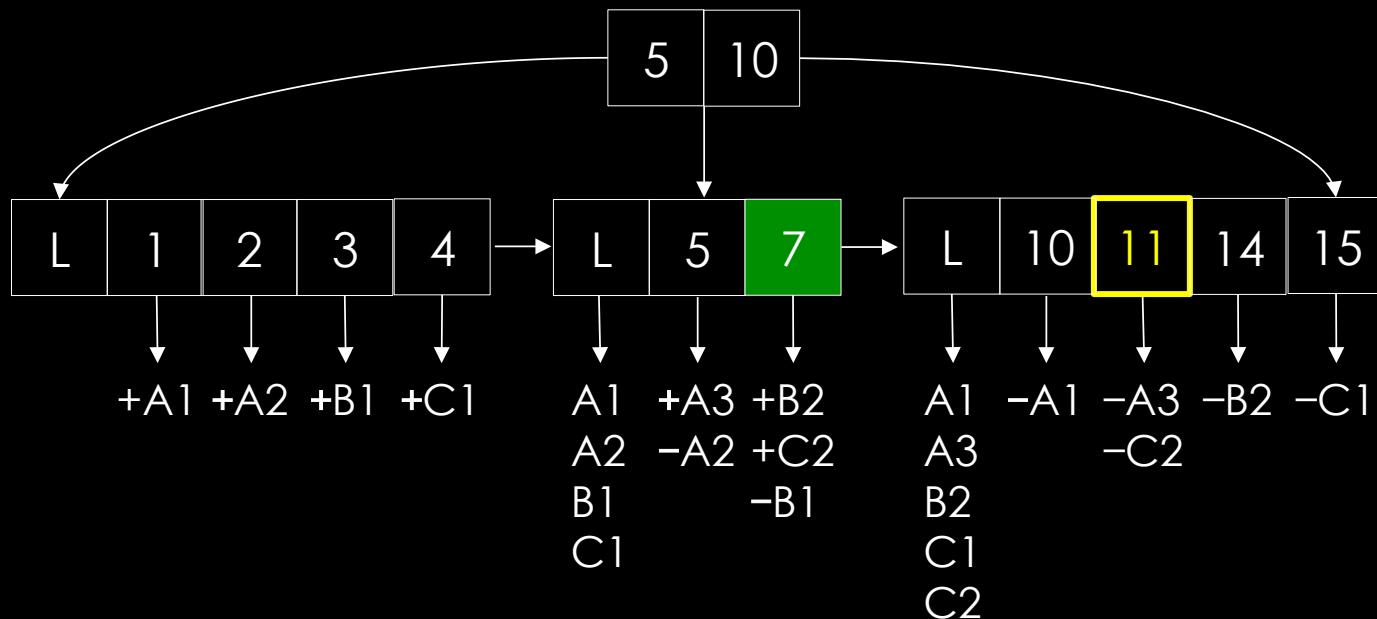
A1
A2
A3

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

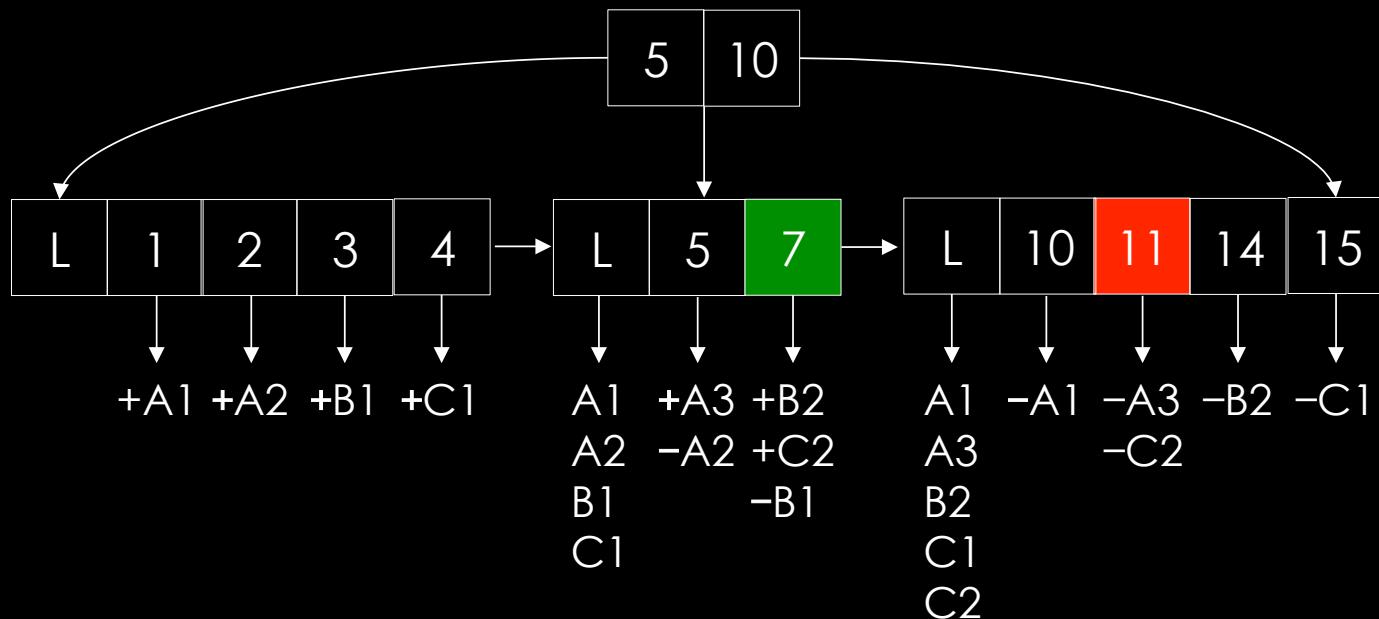
A1
A2
A3

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

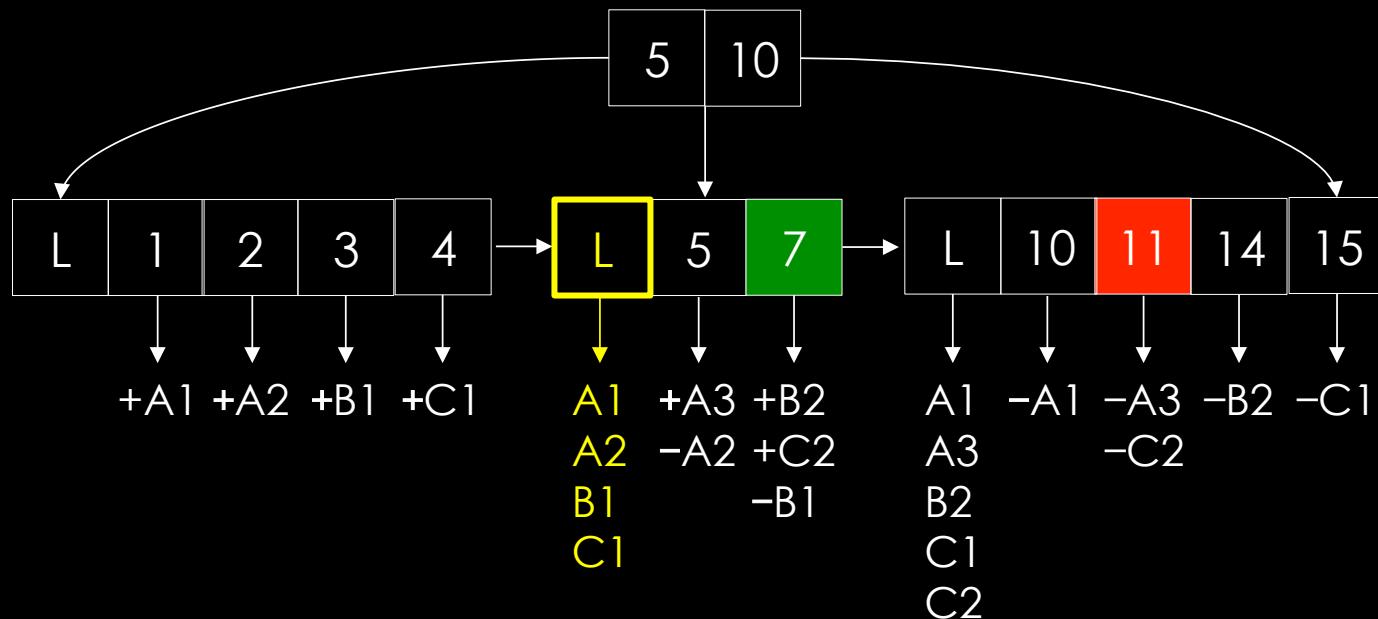
A1
A2
A3

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

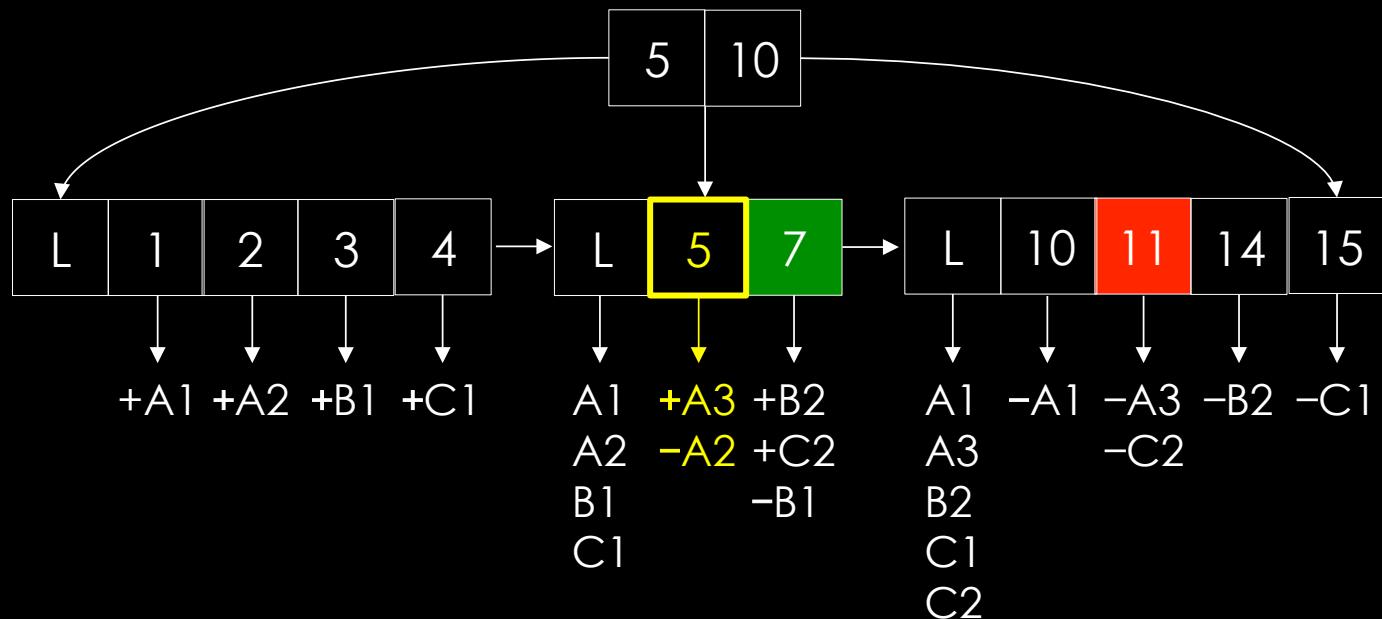
C

C2

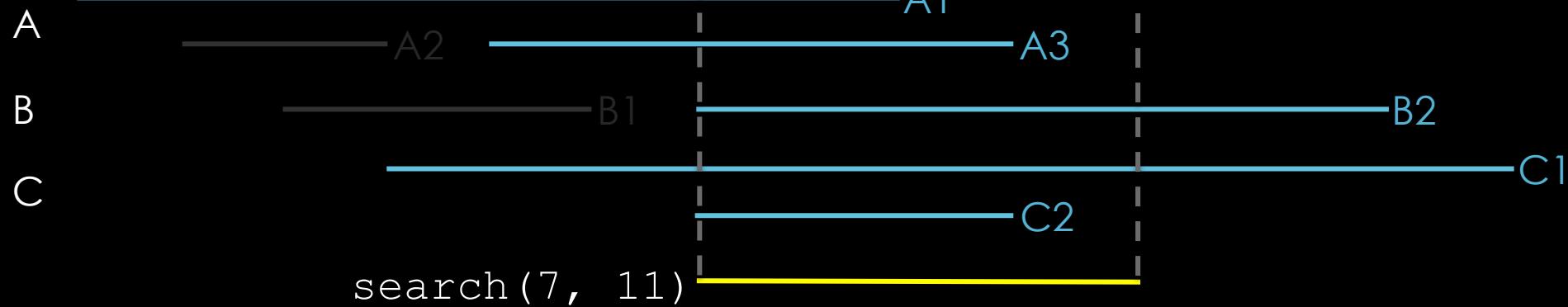
C1

search(7, 11)

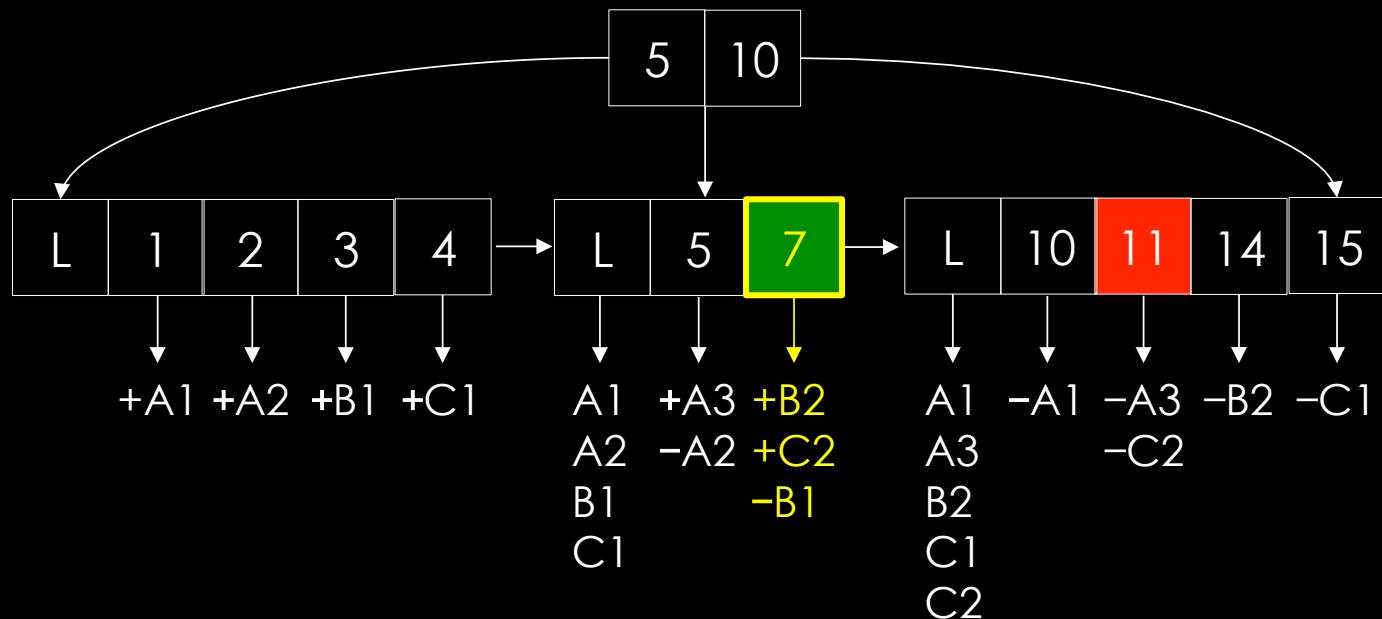
$$I = [A1 \ A2 \ B1 \ C1] + [A3] - [A2] = [A1 \ A3 \ B1 \ C1]$$



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



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



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

A

A2

A1

B

B1

A3

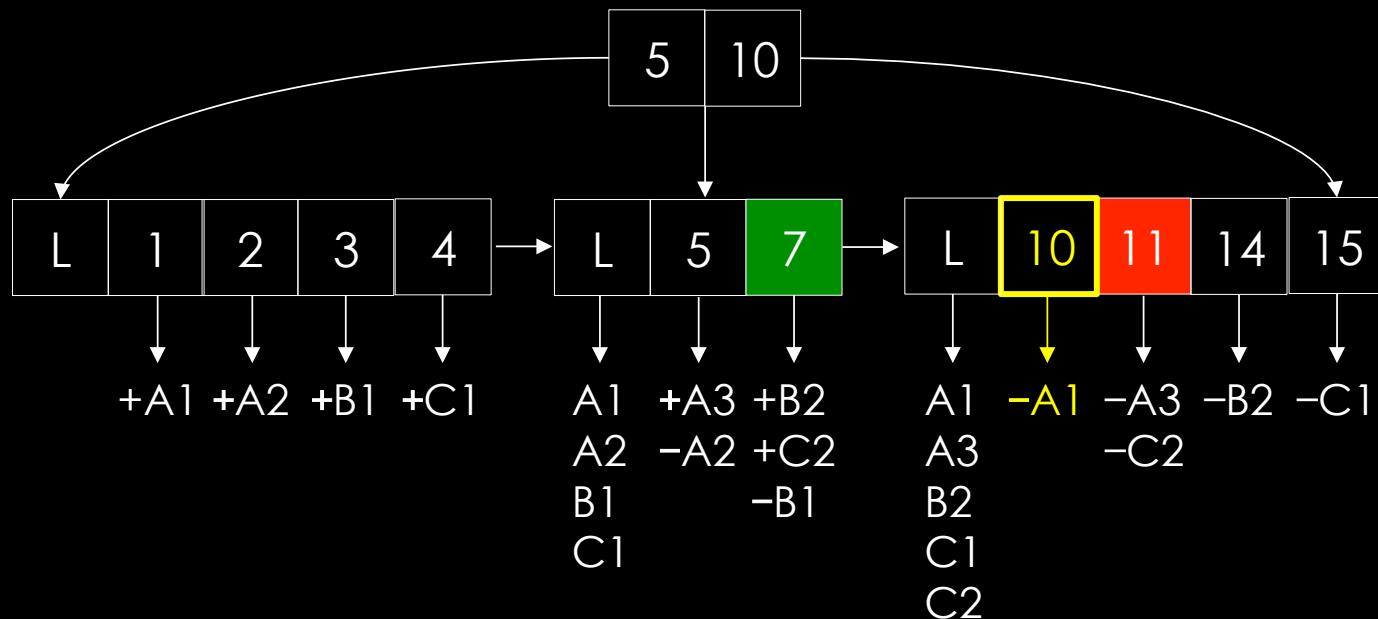
B2

C

C2

search(7, 11)

I = [A1 A3 B2 C1 C3]



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

A

A2

A1

B

B1

A3

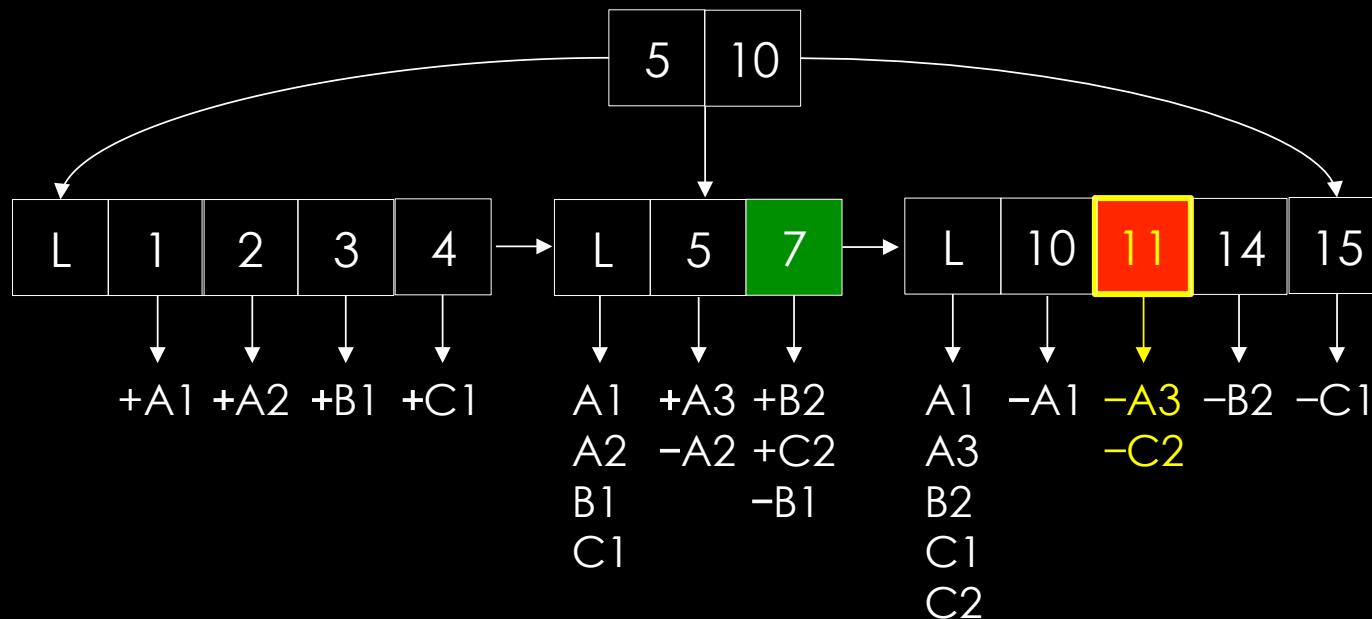
B2

C

C2

search(7, 11)

I = [A1 A3 B2 C1 C3]





- B+ tree
 - Efficient on-disk database
 - Traversal between leaves
- Many files in one index
- Count intersections within the index

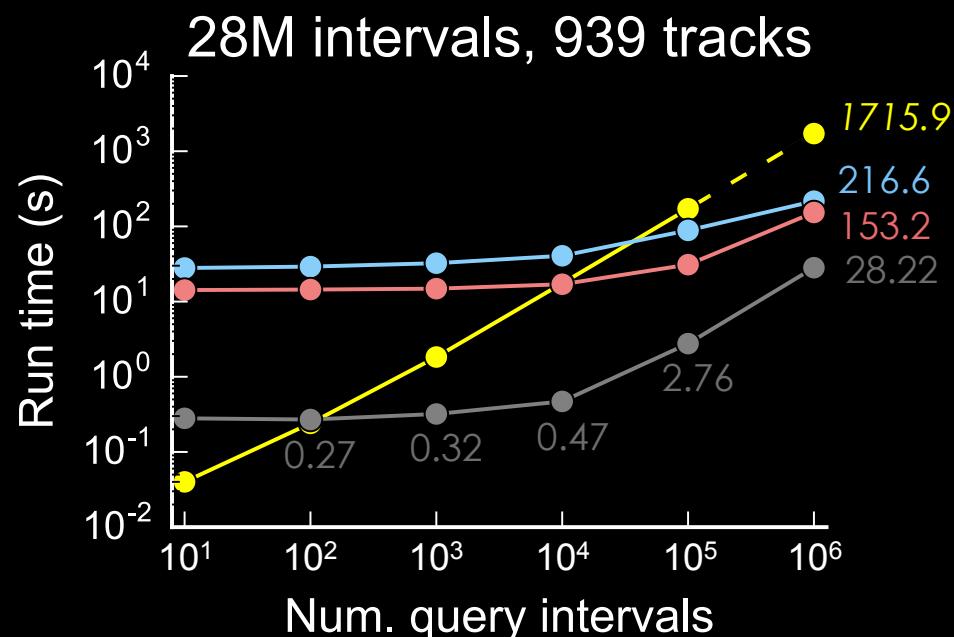
Count intersections for 10 - 1M 100bp intervals



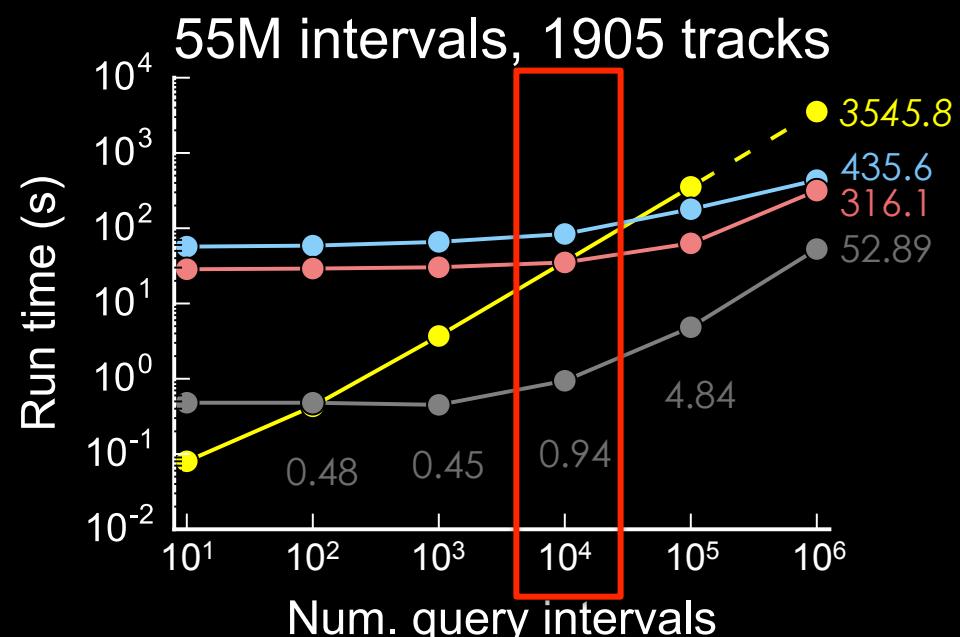
127 reference epigenomes
15 genomic states [CHROHMM Ernst 2012]
1905 tracks

2.6GHz Intel Xeon CPU, 20MB cache

SQLITE3 w/ UCSC binning
TABIX
BEDTOOLS sorted
GIGGLE



Speed up: 60.8X 5.4X 7.7X



67X 5.9X 8.2X



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

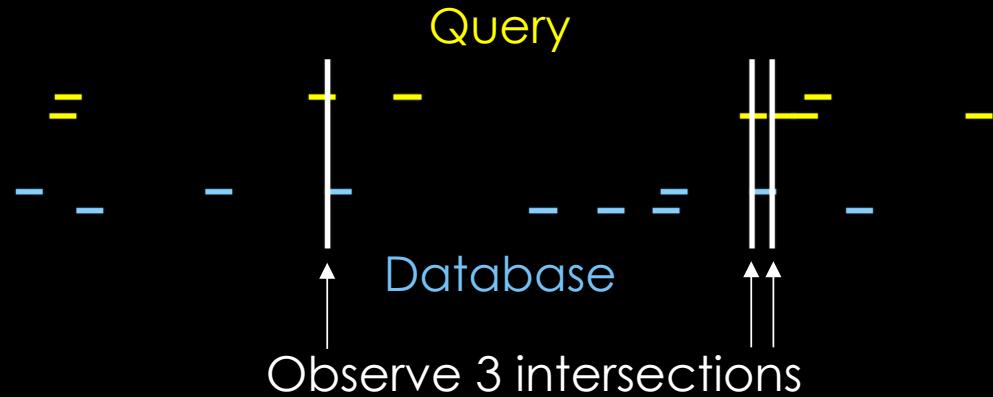
ChIP-seq of heart

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

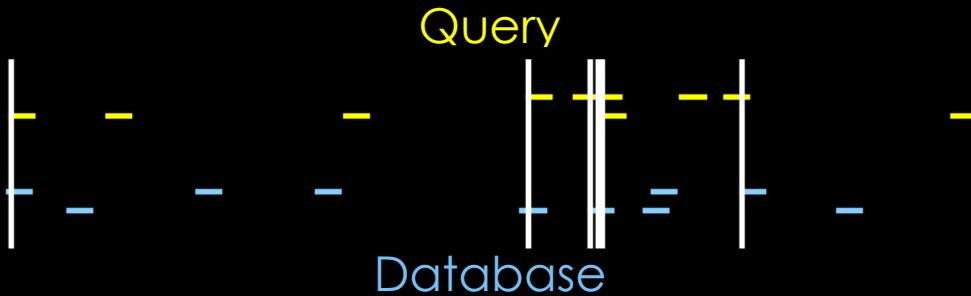
ChIP-seq of embryonic facial prominence

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

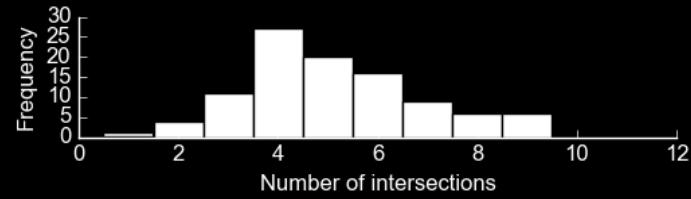
Monte Carlo simulations



Monte Carlo simulations

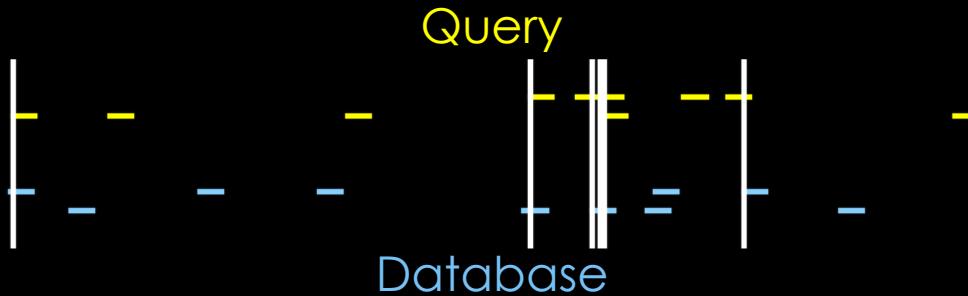


Observe 3 intersections

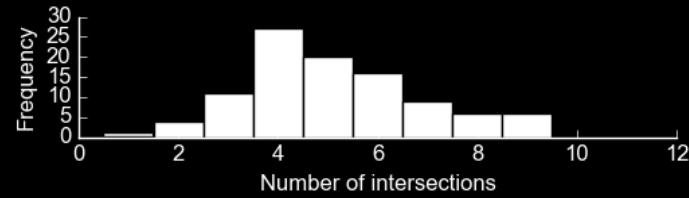


Expect ~4
0.75 "enrichment"

Monte Carlo simulations



Observe 3 intersections



Expect ~4
0.75 "enrichment"

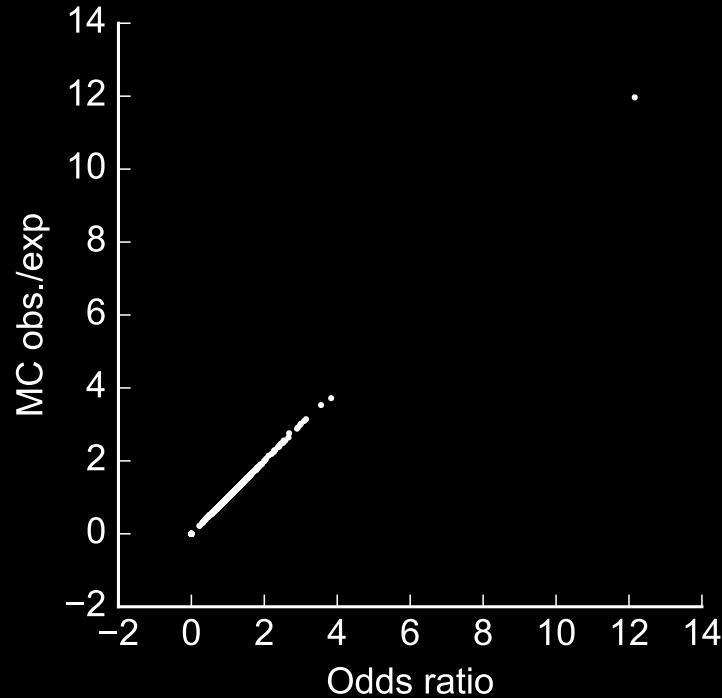
Contingency table

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

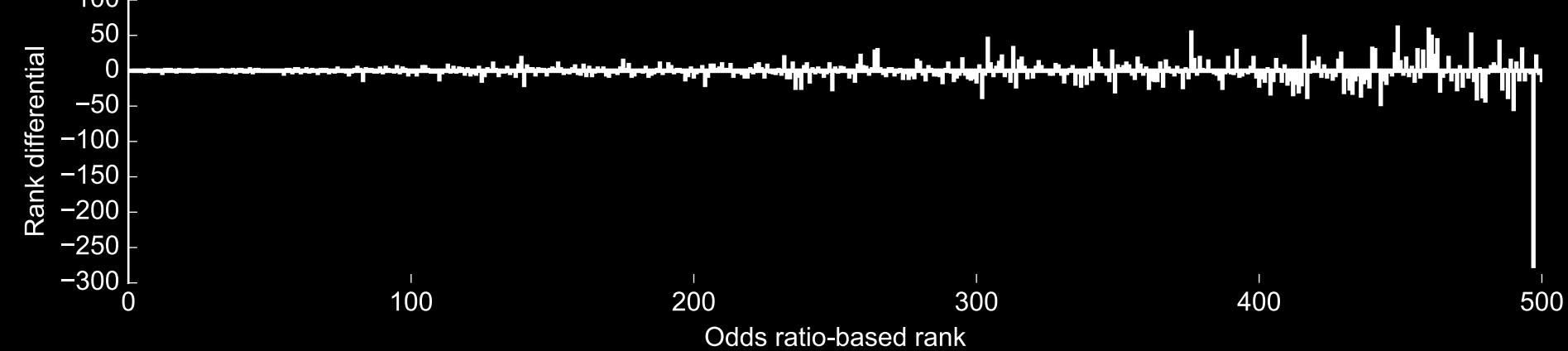


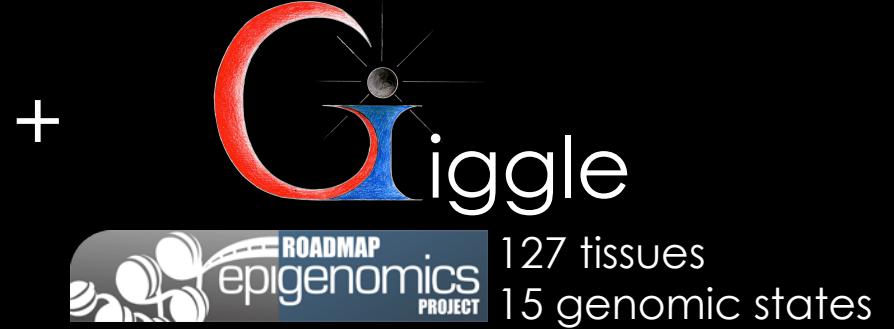
Brent Pedersen

Odds ratio vs. Monte Carlo(N=10000)



Difference in ranking by odds ration and Monte Carlo



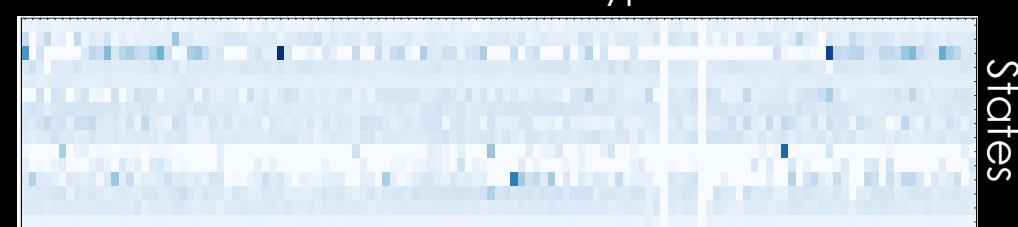


SFARI  40 families
20M variants

De Novos in:

Affected child

```
gqt query -i all.vcf.gz
-p "phenotype == 2"
-g "count(HET) == 1"
-p "phenotype == 1"
-g "HOM_REF"
```



Unaffected child

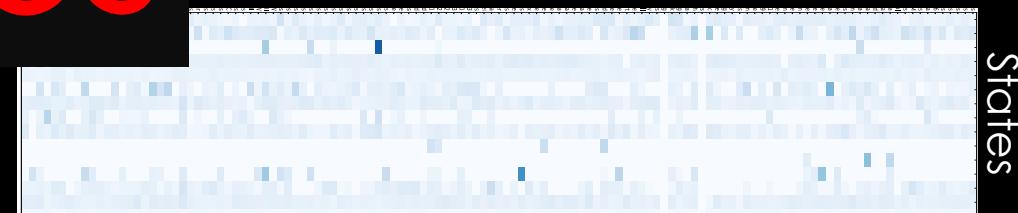
```
-p "phenotype == 1
& parental ... ^.."
-g "count(HET) == 1"
-p "phenotype != 1
|| parental"
-g "HOM_REF"
```

< 6s



Affected female

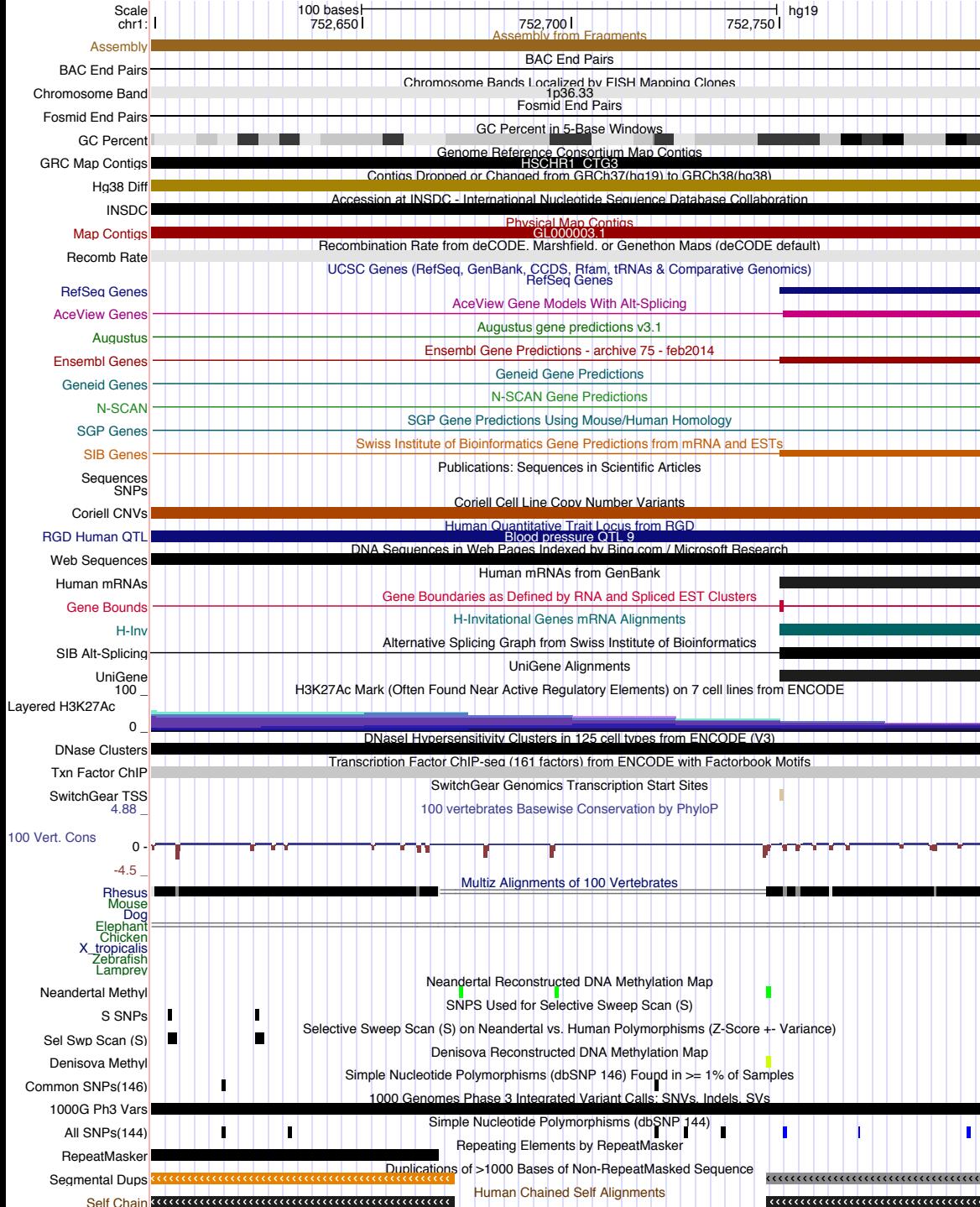
```
-p "phenotype == 2
& sex == 2"
-g "count(HET) == 1"
-p "phenotype == 1"
-g "HOM_REF"
```



Affected European

```
-p "phenotype == 2
& population == 'EUR'"
-g "count(HET) == 1"
-p "phenotype == 1"
-g "HOM_REF"
```





Tonya Di Sera

Command line

```
giggle index -i "tracks/*gz" -o tracks_i  
giggle search -i tracks_i -r chrl:1-1000  
giggle search -9 tracks_i -q query.bed.gz
```

C API

```
struct giggle_index *gi =  
    giggle_load(index_dir_name, uint32_t_ll_giggle_set_data_handler);  
  
struct giggle_query_result *gqr = giggle_query(gi ,chr ,beg ,end, gqr);  
  
for(i = 0; i < gqr->num_files; i++) {  
    struct giggle_query_iter *gqi = giggle_get_query_itr(gqr, i);  
  
    while (giggle_query_next(gqi, &result) == 0)  
        printf("%s\n", result);  
  
    giggle_iter_destroy(&gqi);  
}
```

Python API COMING SOON



Brent Pedersen

TODO:

Index remote files

Faster indexing

Index updating (append)

Integration with Genotype Query Tools (GQT)



github.com/ryanlayer/giggle

Ryan M. Layer

ryan.layer@gmail.com

[@ryanlayer](https://twitter.com/ryanlayer)

Aaron Quinlan, Brent Pedersen, Tonya Di Sera