

# CS186 Discussion #5

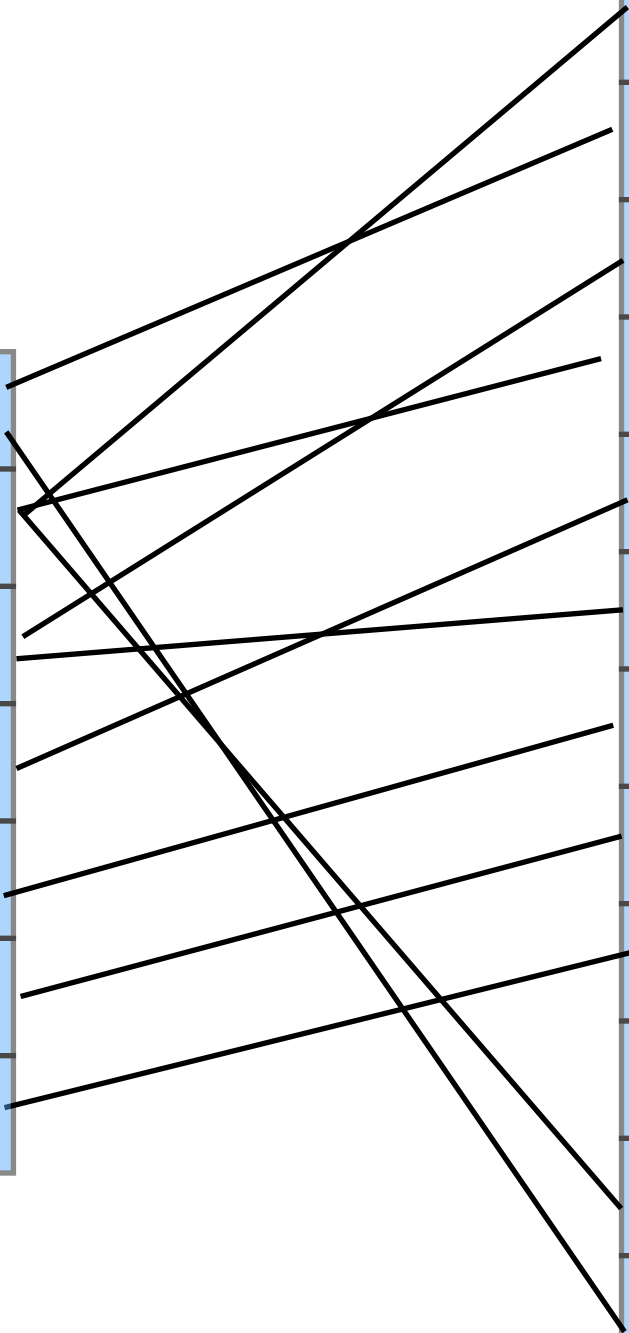
(Indexes, B+ Trees)

Fill out the Index definitions  
on your worksheet.

Bob	CS
Joe	Biology
Jim	English
Sue	CS
Jill	Math
Tim	English
Alice	Physics
Rob	History
Sue	Art
Brad	Chem
Jen	CS
Ted	Biology

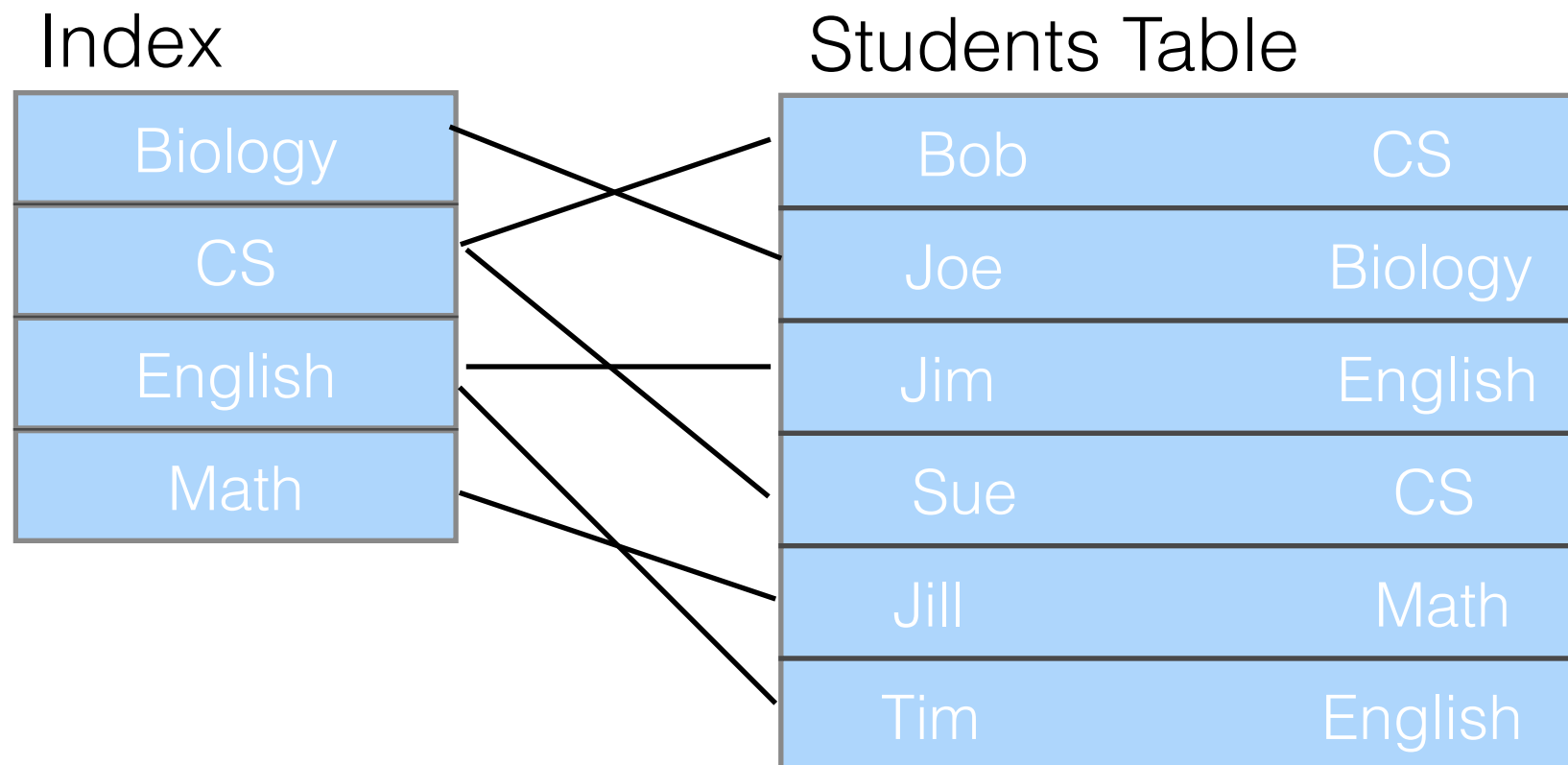
Biology
CS
English
Math
Physics
History
Art

Bob	CS
Joe	Biology
Jim	English
Sue	CS
Jill	Math
Tim	English
Alice	Physics
Rob	History
Sue	Art
Brad	Chem
Jen	CS
Ted	Biology



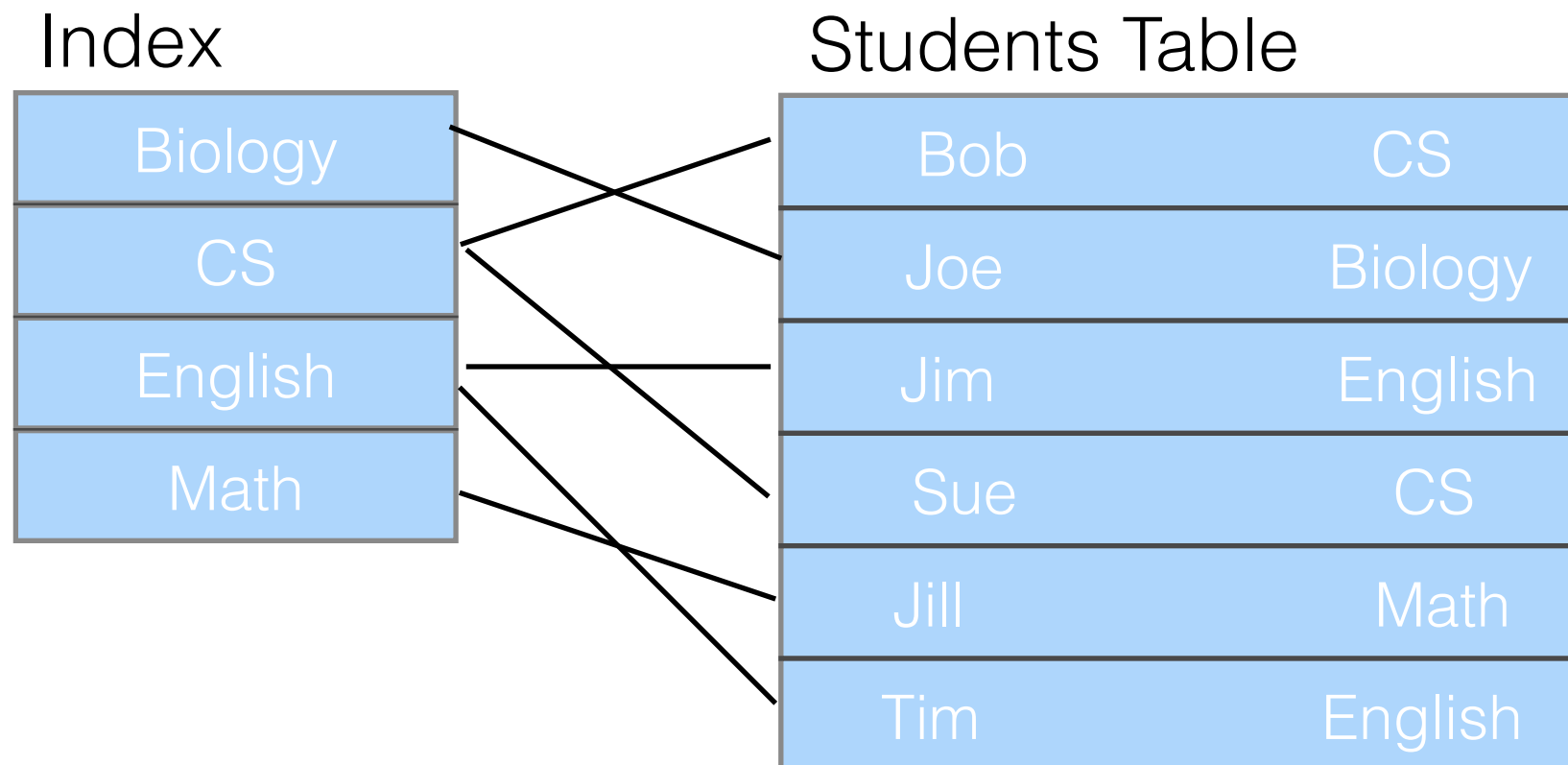
# Indexes

- Disk-based data structure for fast lookup by value (search key)
  - Ex: Find students in the CS department



# Hash Index

- Good for equality search



# Hash Index

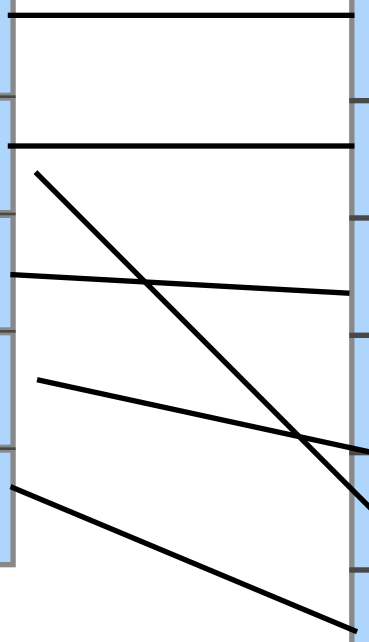
- Good for equality search

Index on GPA

3.4
2.0
2.2
3.5
3.0

Students Table

Bob	3.4
Joe	2.0
Jim	2.2
Sue	3.5
Jill	2.0
Tim	3.0



# Hash Index

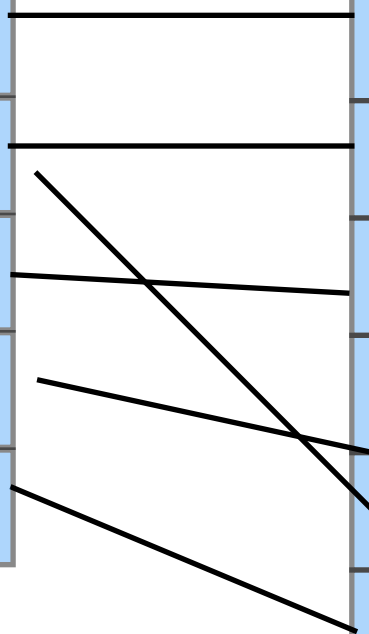
- Good for equality search
- Can't do range queries

Index on GPA

3.4
2.0
2.2
3.5
3.0

Students Table

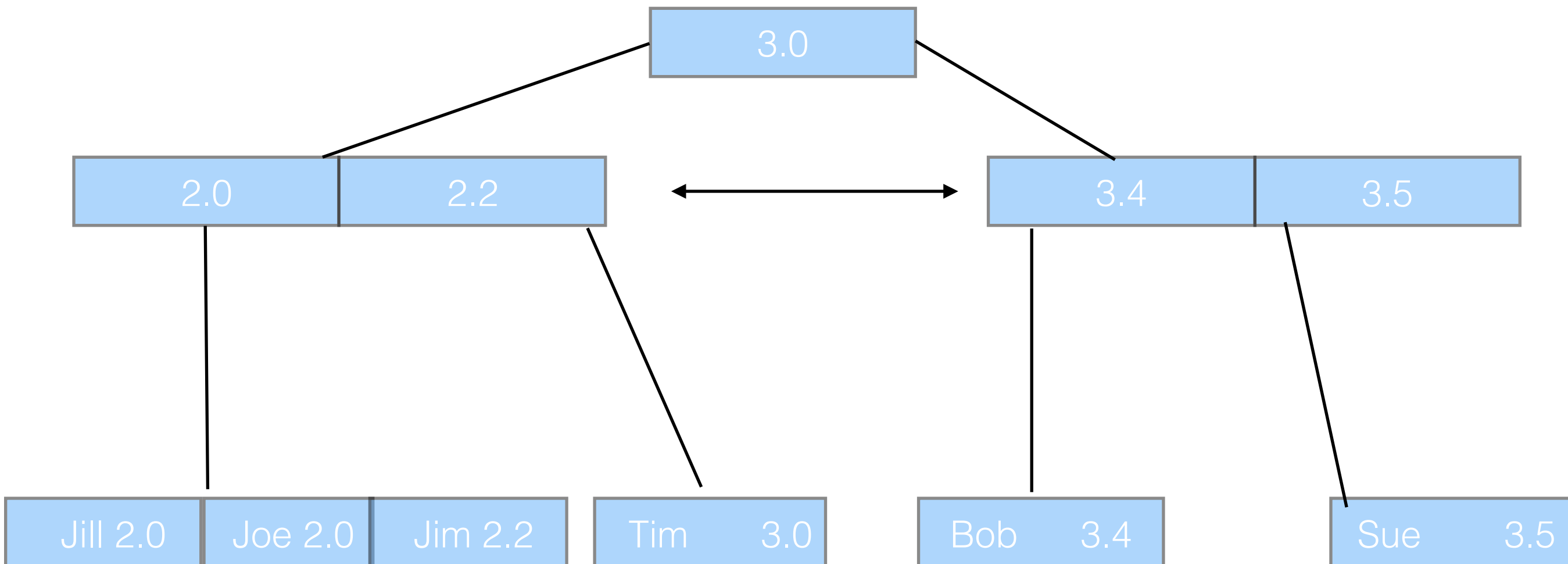
Bob	3.4
Joe	2.0
Jim	2.2
Sue	3.5
Jill	2.0
Tim	3.0





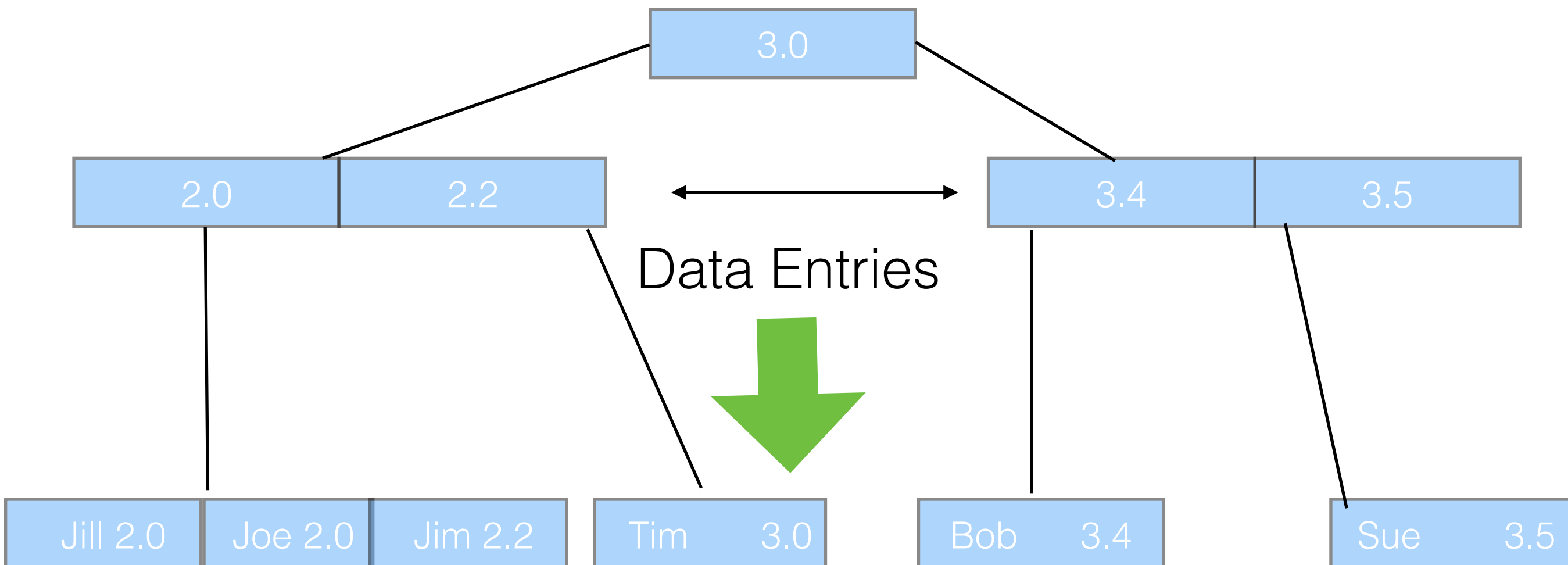
# Tree Index

- Fast equality search, insertion, deletion
- Can do range search



# Tree Index

- Fast equality search, insertion, deletion
- Can do range search

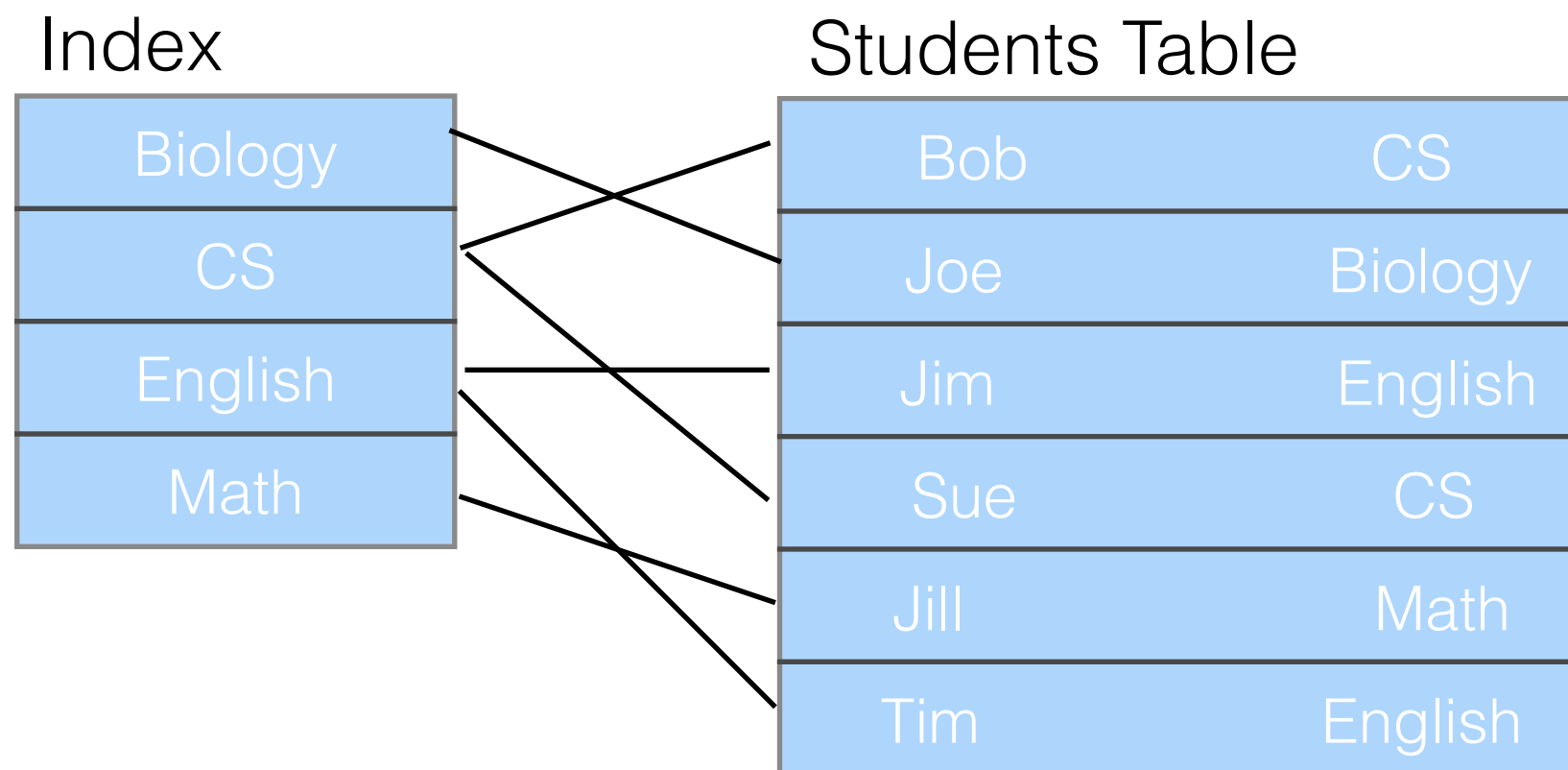


# 3 Ways to Store Entries in Index

- Alt 1: Actual data stored at index
  - Can have at most one index per table
- Alt 2: Store key and record ID of the matching record
- Alt 3: Store key and list of record IDs

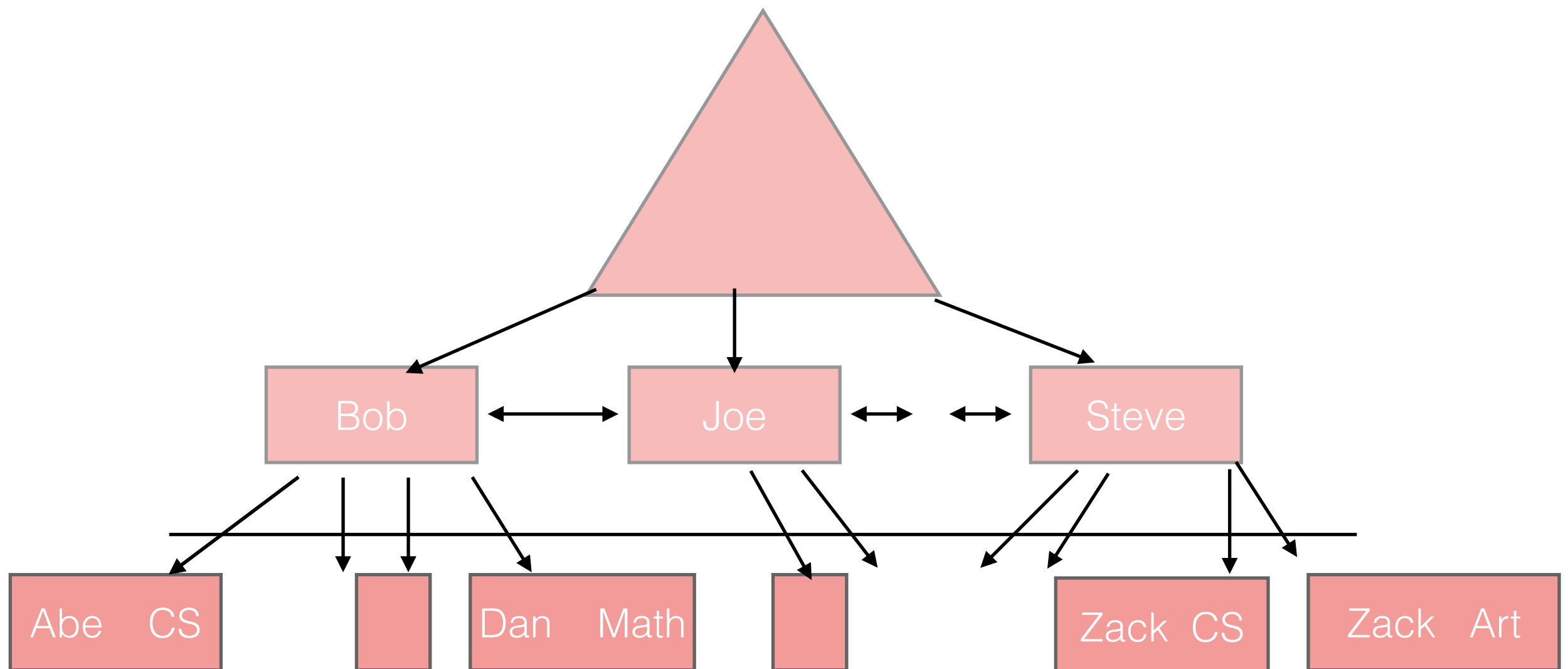
# Alt 1

Actual data stored at the index



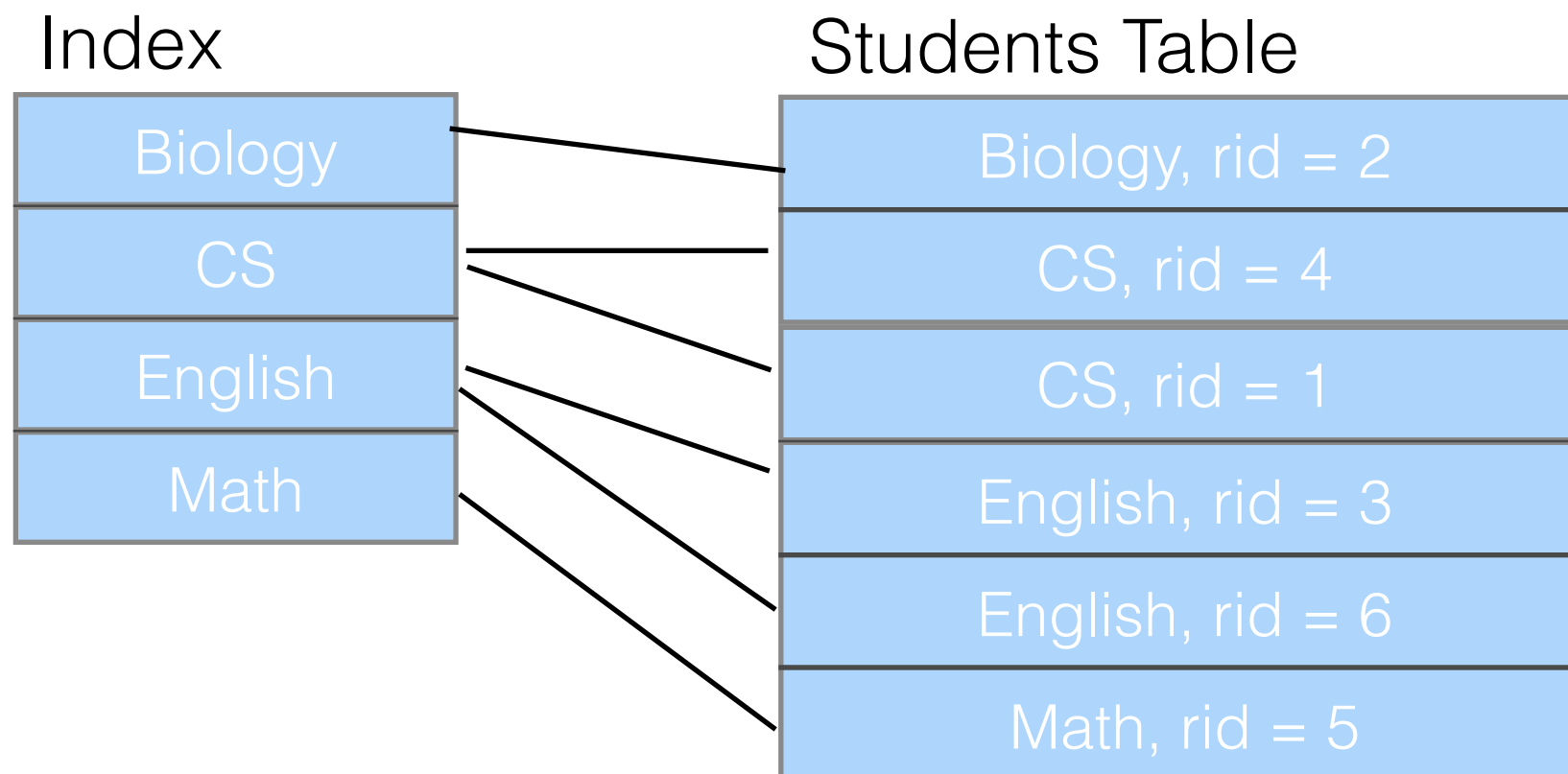
# Alt 1

Actual data stored at the index



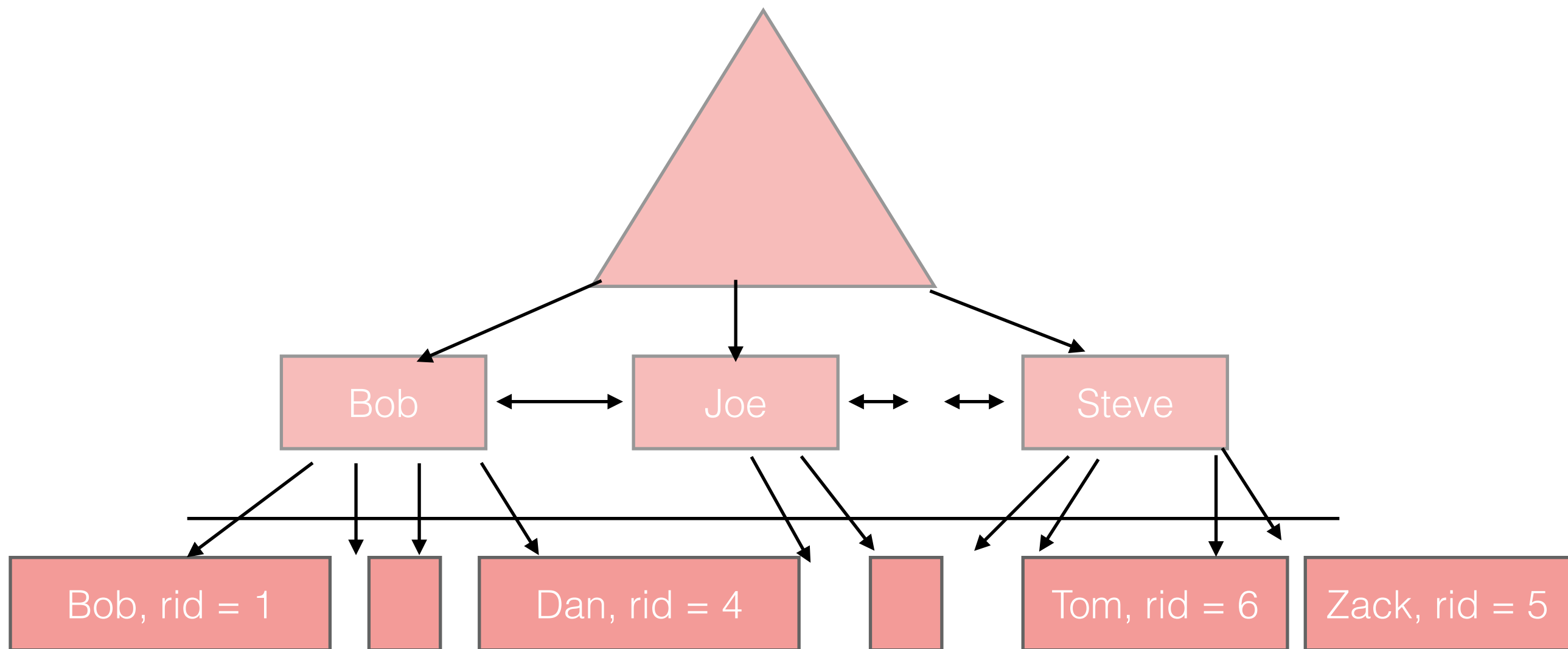
# Alt 2

<key, record ID>



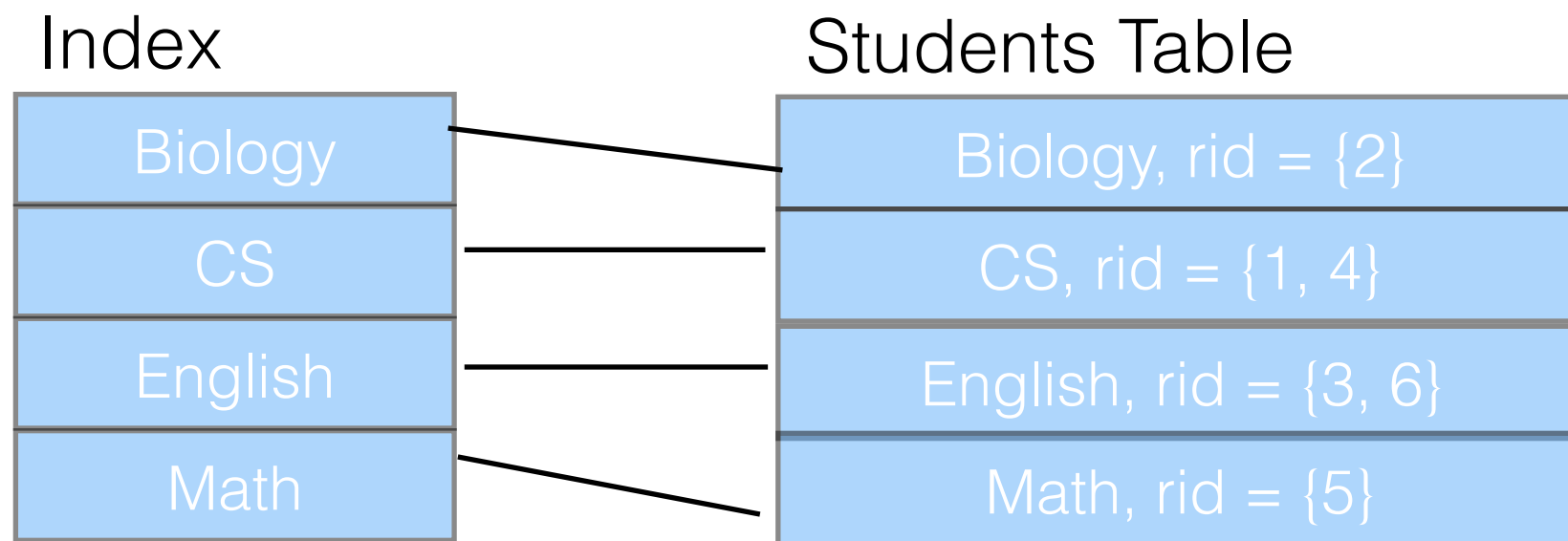
# Alt 2

<key, record ID>



# Alt 2

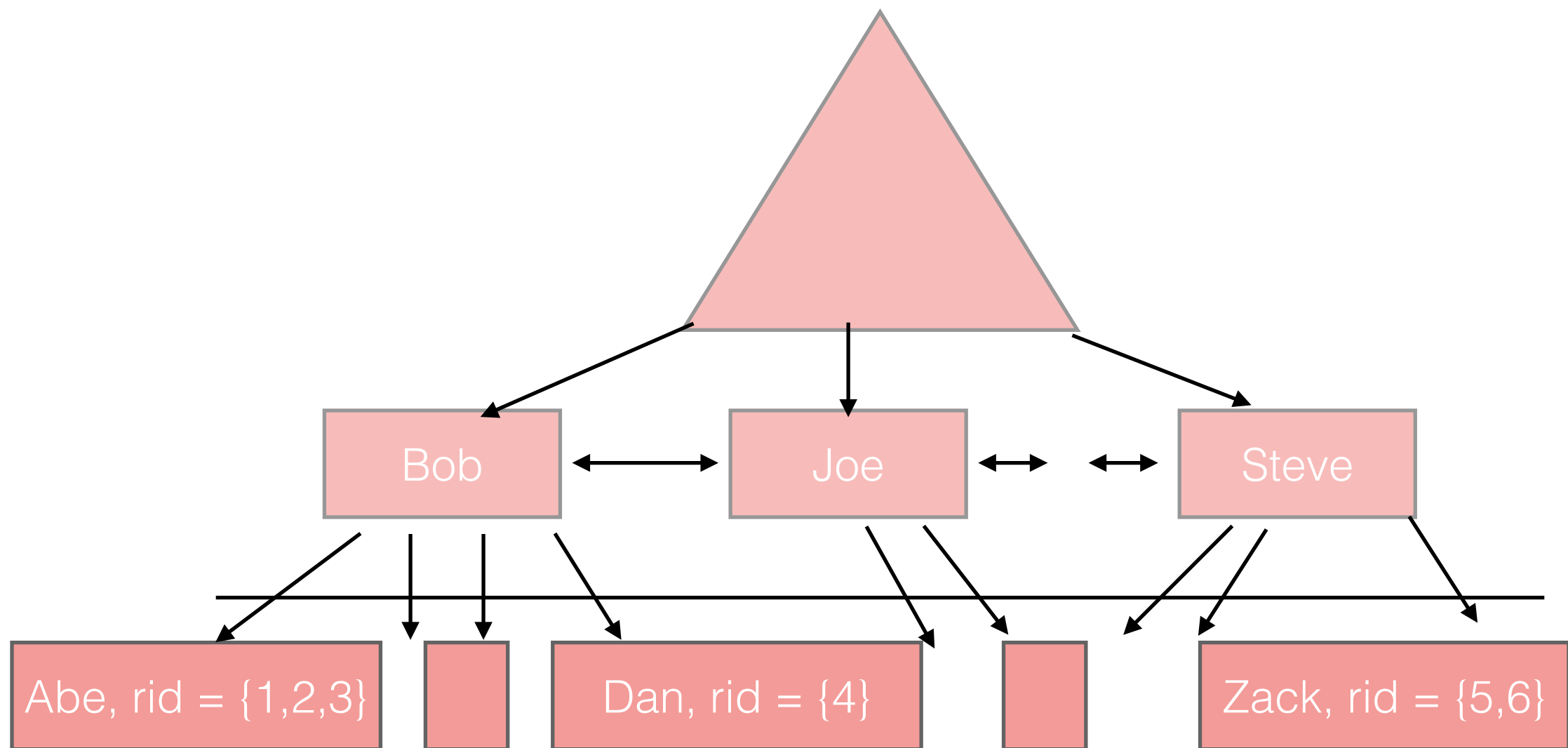
<key, record ID>





# Alt 3

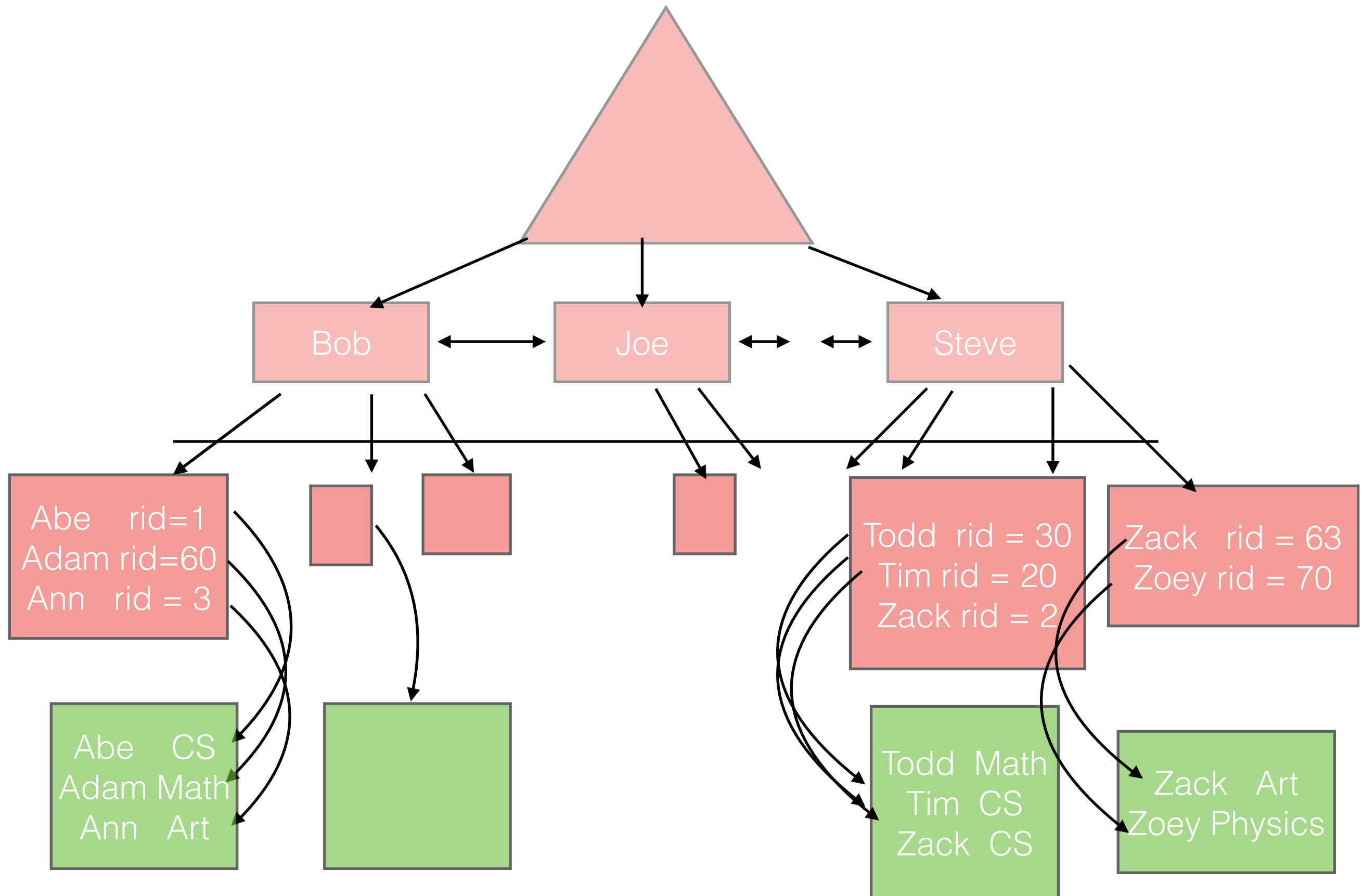
<key, list of matching record IDs>



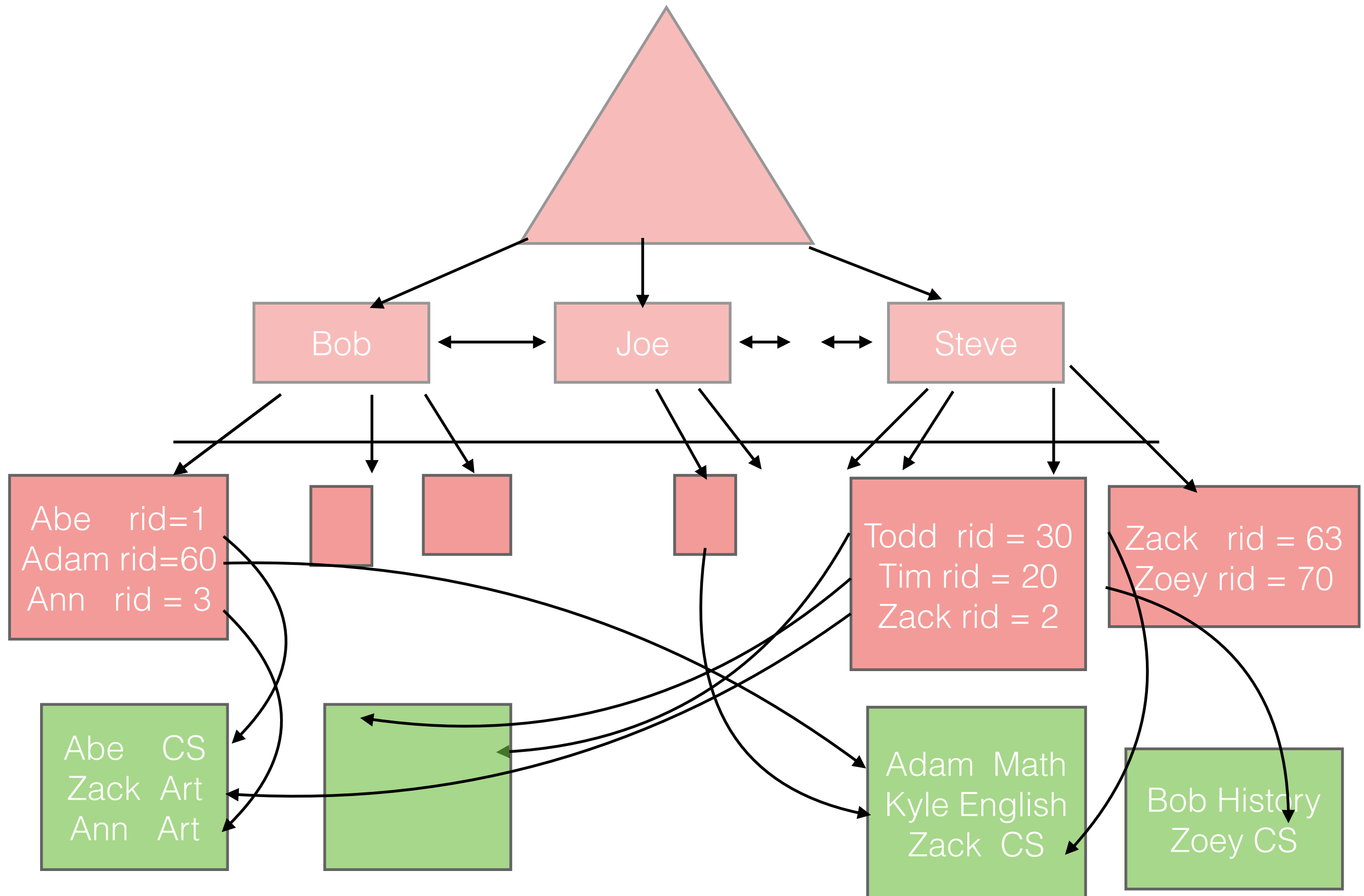
# Clustered vs Unclustered

- Clustered - index data entries are stored in (approximate) order by value of search keys in data records
- Can be clustered on at most one key
- Alternative 1 is always clustered

# Clustered



# Unclustered



Operation	Heap File	Sorted File	Clustered File
Scan all records	B	B	1.5B
Equality Search	0.5B	$\log_2(B)$	$(\log_2 1.5B + 1)$
Range Search	B	$\log_2(B) + \# \text{ pages matched}$	$(\log_2 1.5B + \# \text{ matches})$
Insert	2	$\log_2(B) + (B/2) * 2$	$(\log_2 1.5B + 2)$
Delete	$0.5B + 1$	$\log_2(B) + (B/2) * 2$	$(\log_2 1.5B + 2)$

# Worksheet: Indexing

What are important factors in determining whether or not you should add an index to add to a table?

What are important factors in determining whether or not you should add an index to add to a table?

- Should know which field to cluster on (calculate based on typical queries run)
- Decide if you even want to cluster (high maintenance cost)



Clowns (Name text, SID int, Experience int,  
Age int)

RID: 022 - Casey | 134234 | 2 | 35

RID: 051 - Merk | 955551 | 2 | 35

RID: 101 - Test | 045671 | 11 | 15

Clowns (Name text, SID int, Experience int,  
Age int)

RID: 022 - Casey | 134234 | 2 | 35

RID: 051 - Merk | 955551 | 2 | 35

RID: 101 - Test | 045671 | 11 | 15

## Alternative 1

Casey | 134234 | 2 | 35

Merk | 955551 | 2 | 35

Test | 045671 | 11 | 15

Clowns (Name text, SID int, Experience int,  
Age int)

RID: 022 - Casey | 134234 | 2 | 35

RID: 051 - Merk | 955551 | 2 | 35

RID: 101 - Test | 045671 | 11 | 15

## Alternative 2

<key, RID>

<35, 022>

<35, 051>

<15, 101>

Clowns (Name text, SID int, Experience int,  
Age int)

RID: 022 - Casey | 134234 | 2 | 35

RID: 051 - Merk | 955551 | 2 | 35

RID: 101 - Test | 045671 | 11 | 15

## Alternative 3

<key, [RID list]>

<35, [022, 051]>

<15, [101]>

True or False? Given the table Students(sid char(20), gpa float, age integer), a clustered tree based index on gpa will increase the performance of the following query:

```
SELECT * FROM Students where age > 20 AND  
gpa > 3.5;
```

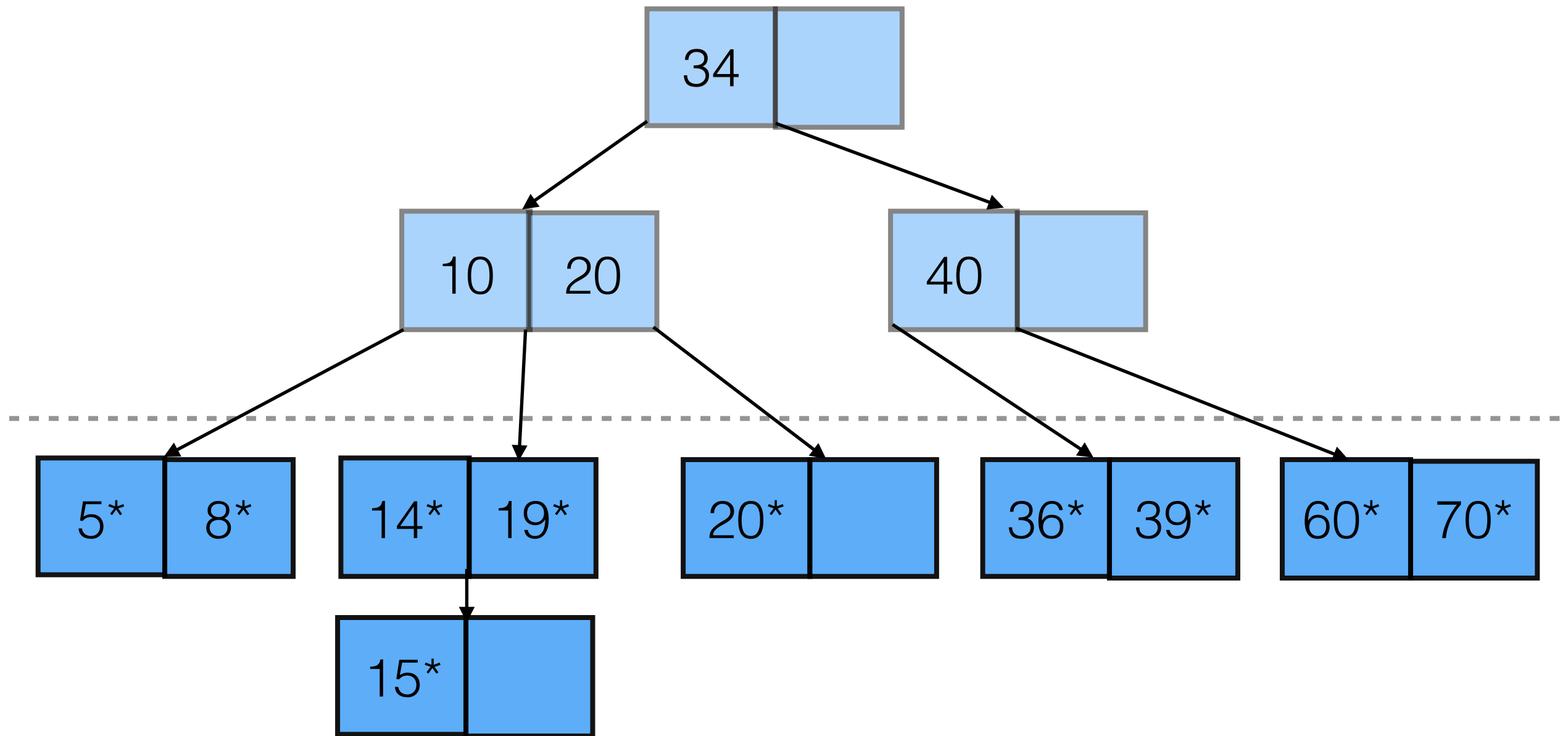
True or False? Given the table Students(sid char(20), gpa float, age integer), a clustered tree based index on gpa will increase the performance of the following query:  
`SELECT * FROM Students where age > 20 AND gpa > 3.5;`

TRUE

# ISAM

- Simple, static structure
- Created by:
  - Sorting records by index search key (e.g. “gpa”)
  - Building a tree on top of those records

# ISAM



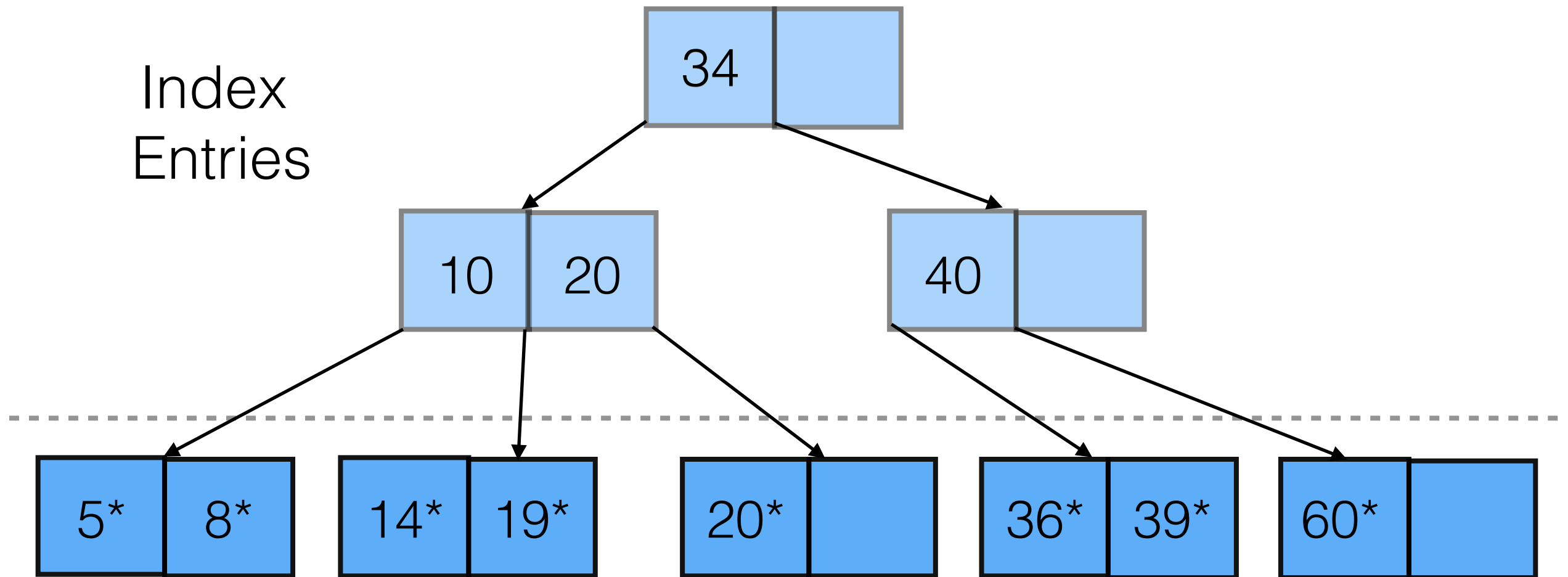


# ISAM - Insert X

- Traverse index pages to find correct leaf L
- If L has space:
  - Insert X in that page
- Else:
  - If an overflow page has space, insert X in that page
  - Else, create a new overflow page and insert X

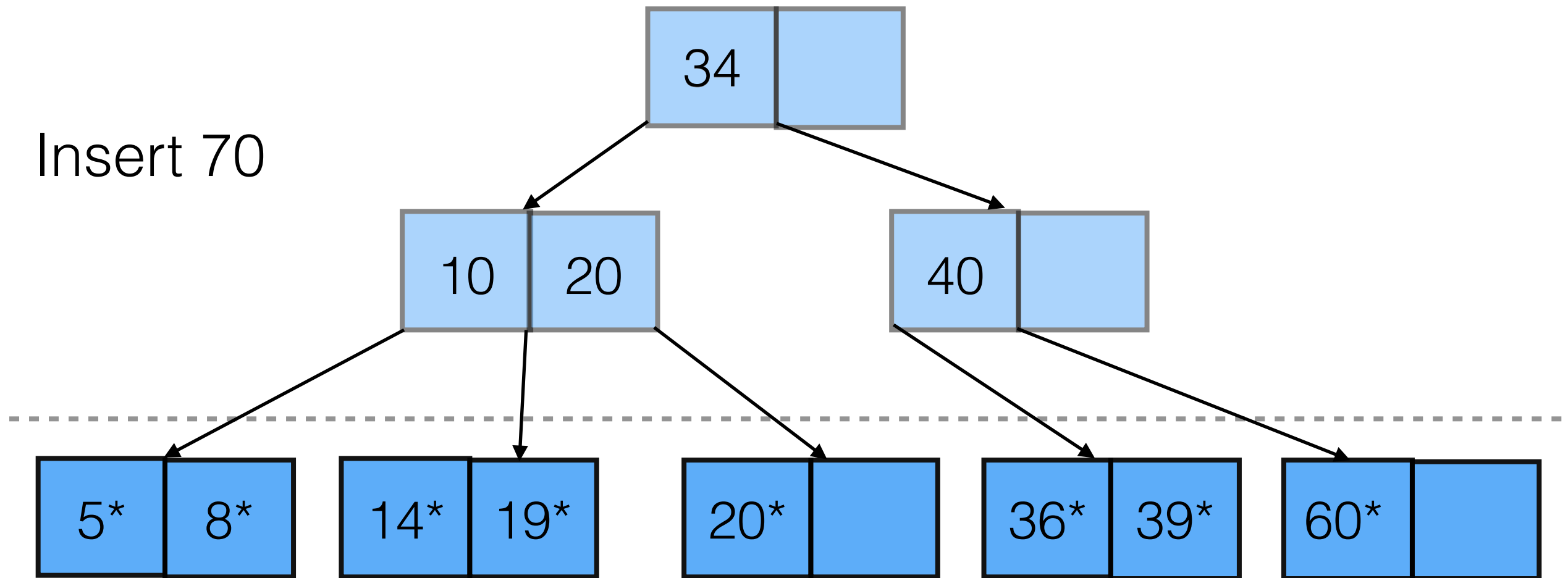
# ISAM

Index  
Entries



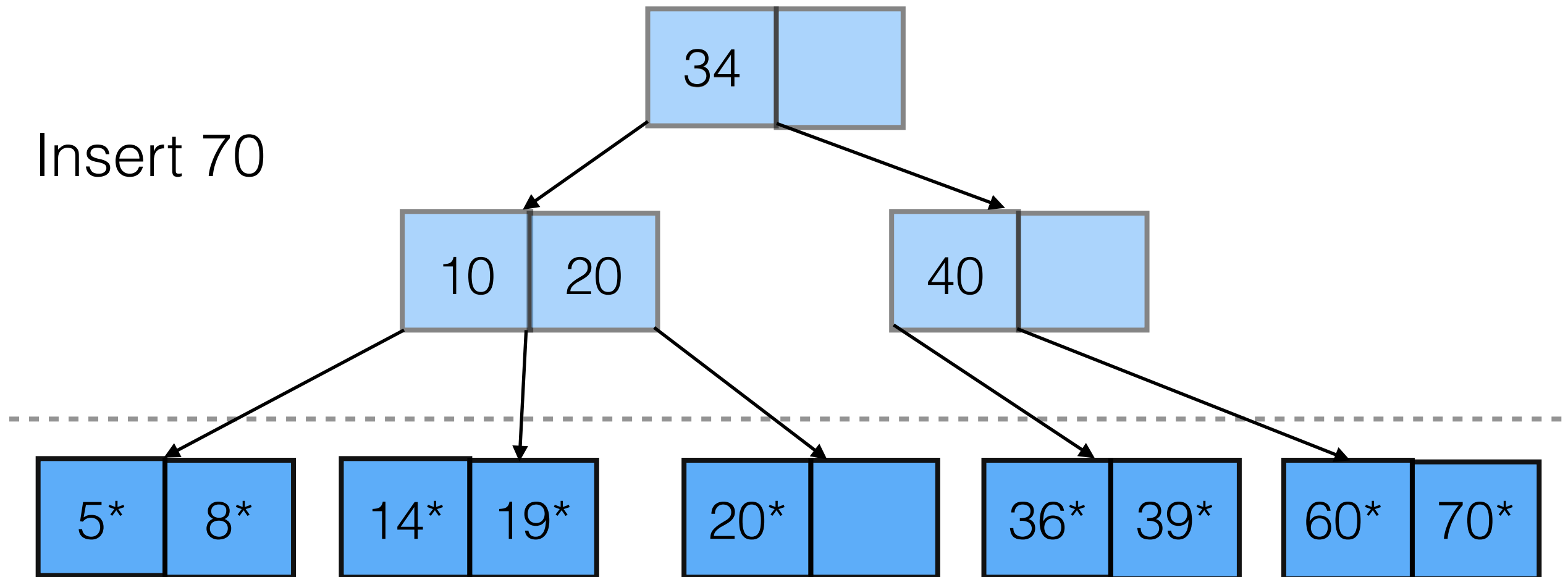
# ISAM

Insert 70



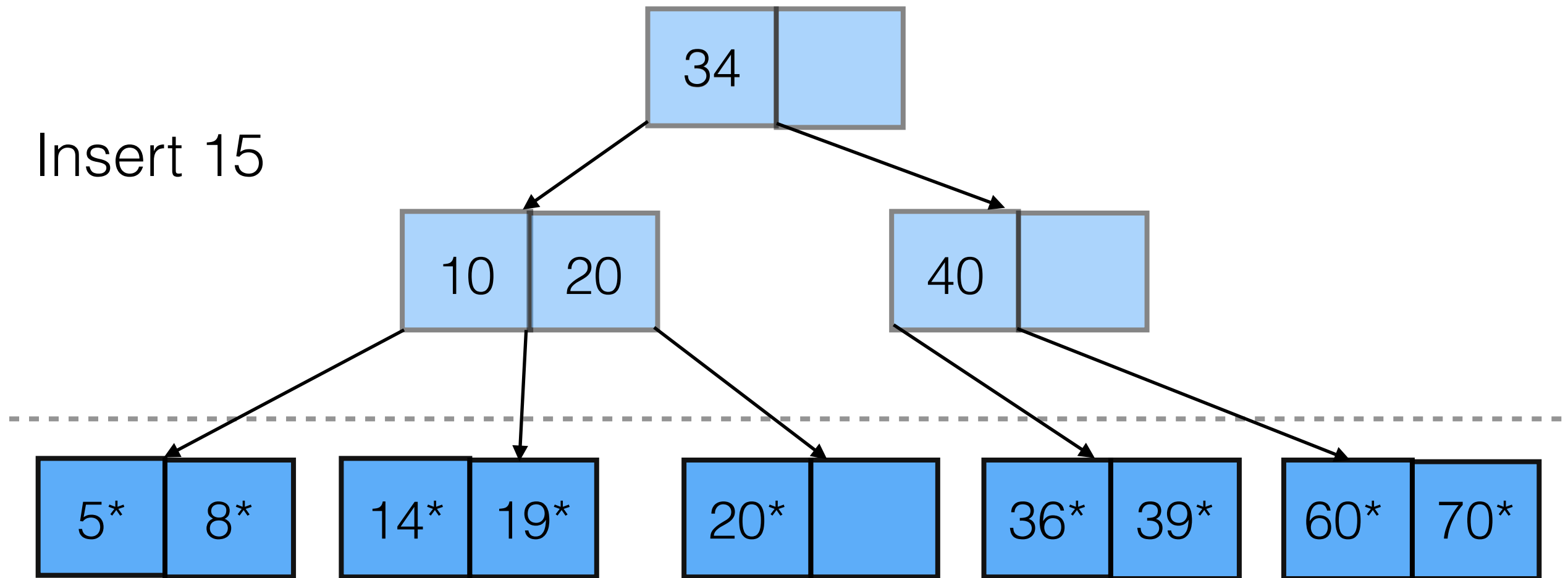
# ISAM

Insert 70



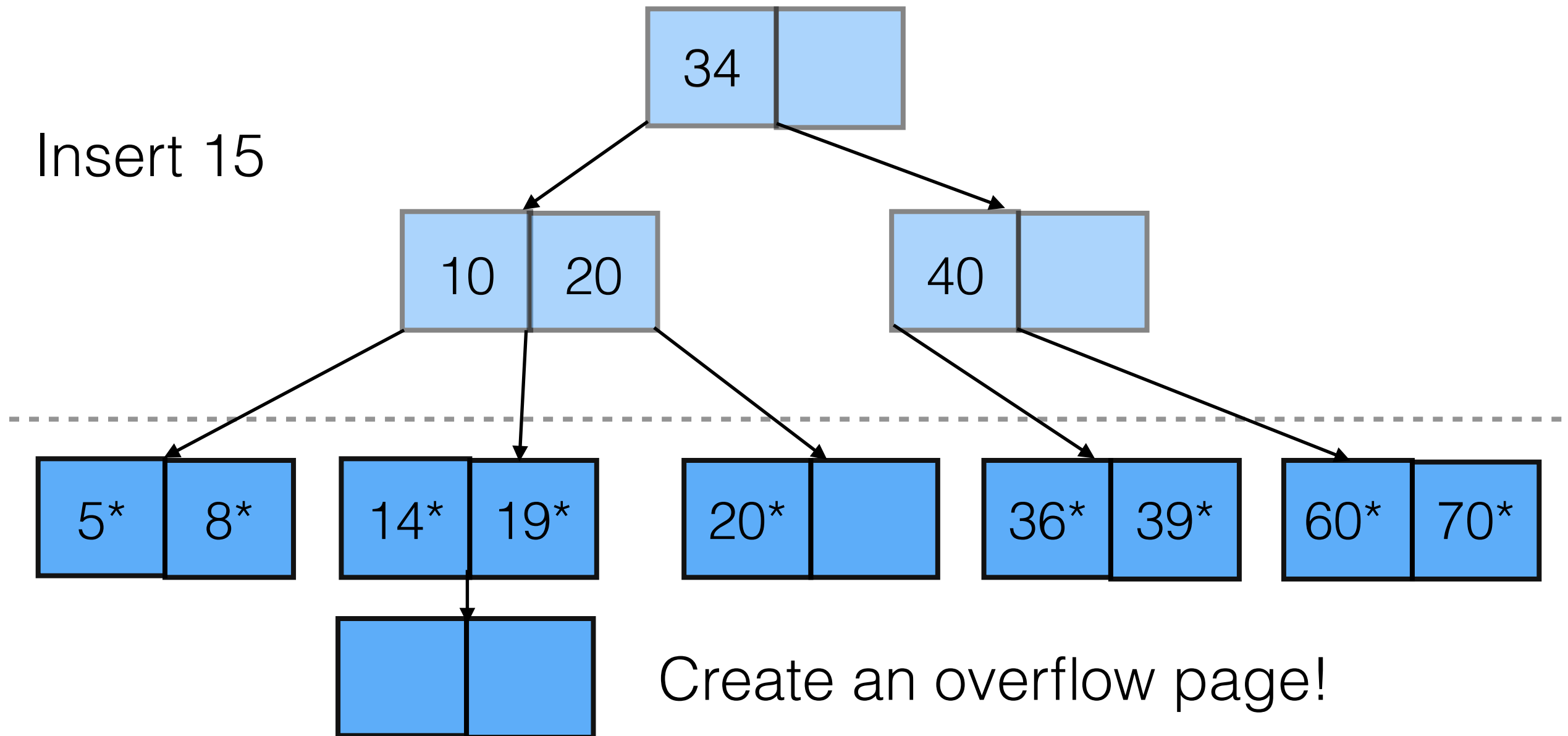
# ISAM

Insert 15



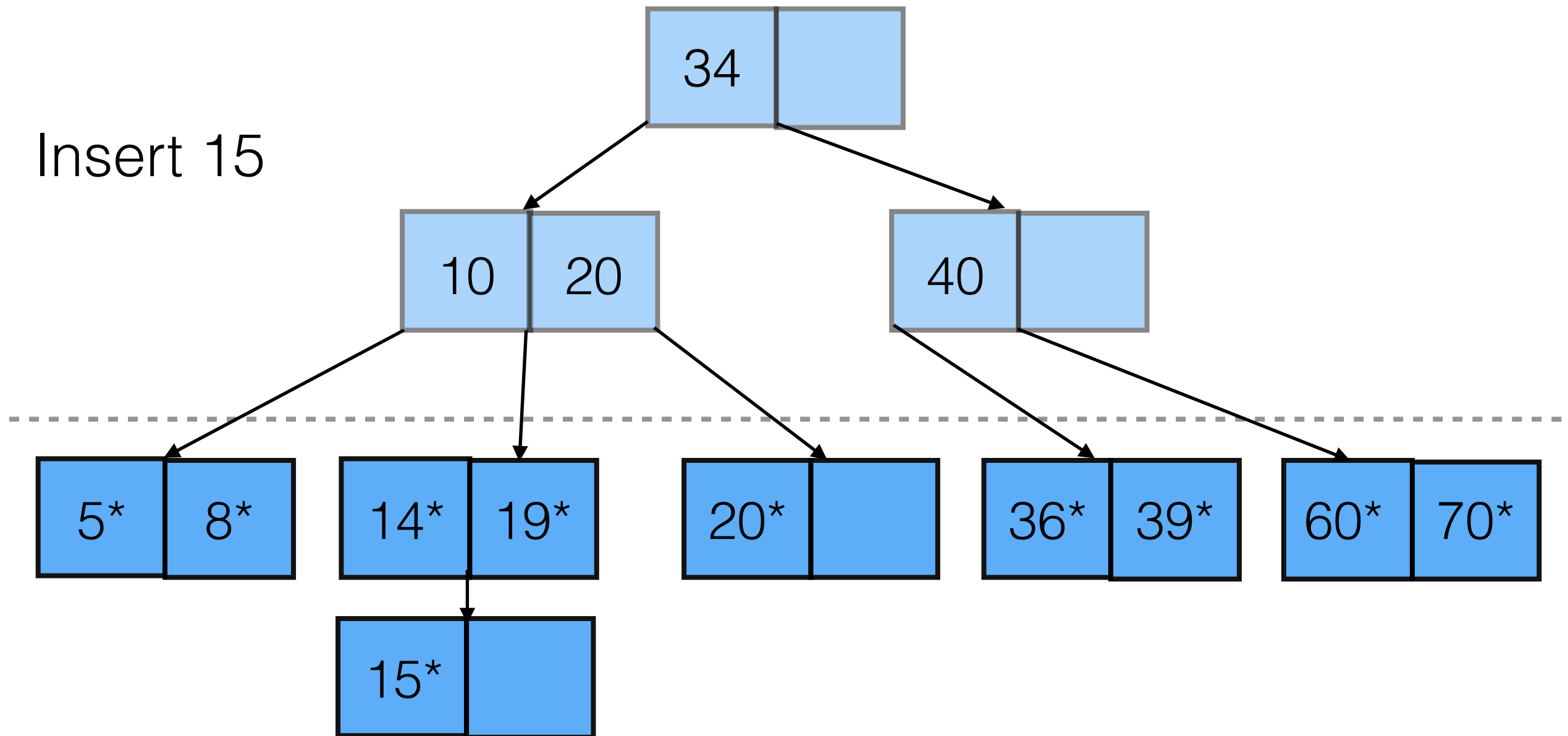
# ISAM

Insert 15



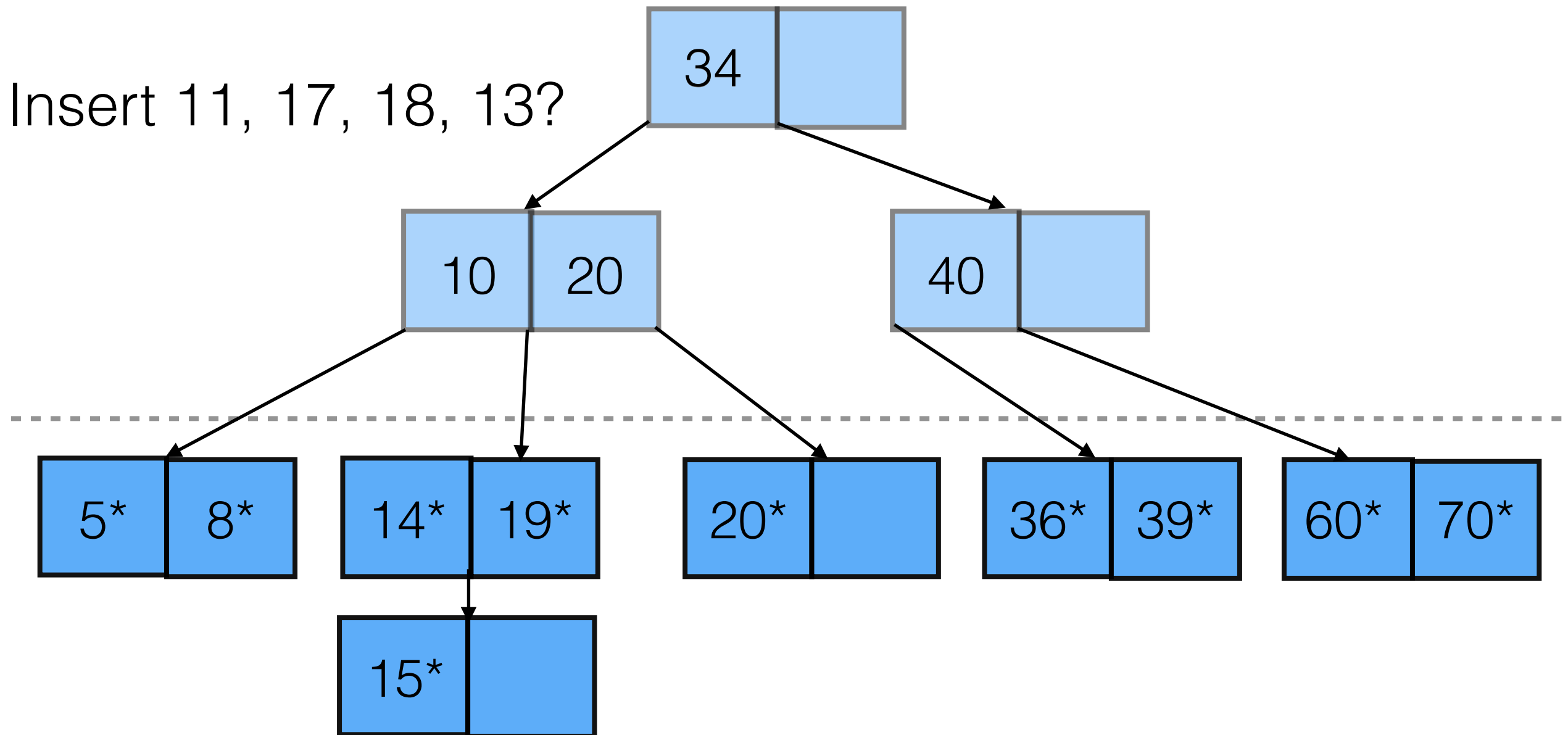
# ISAM

Insert 15



# ISAM

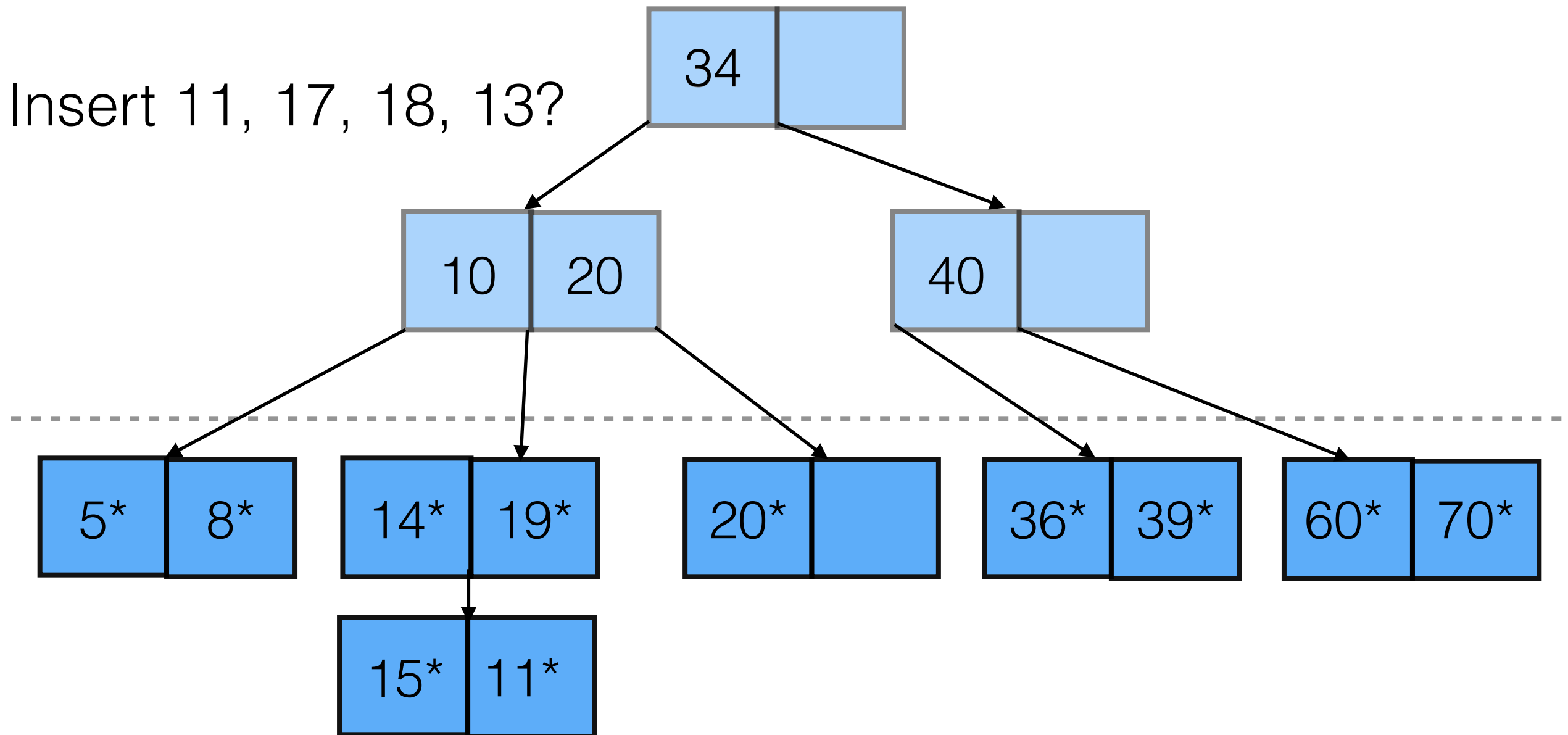
Insert 11, 17, 18, 13?





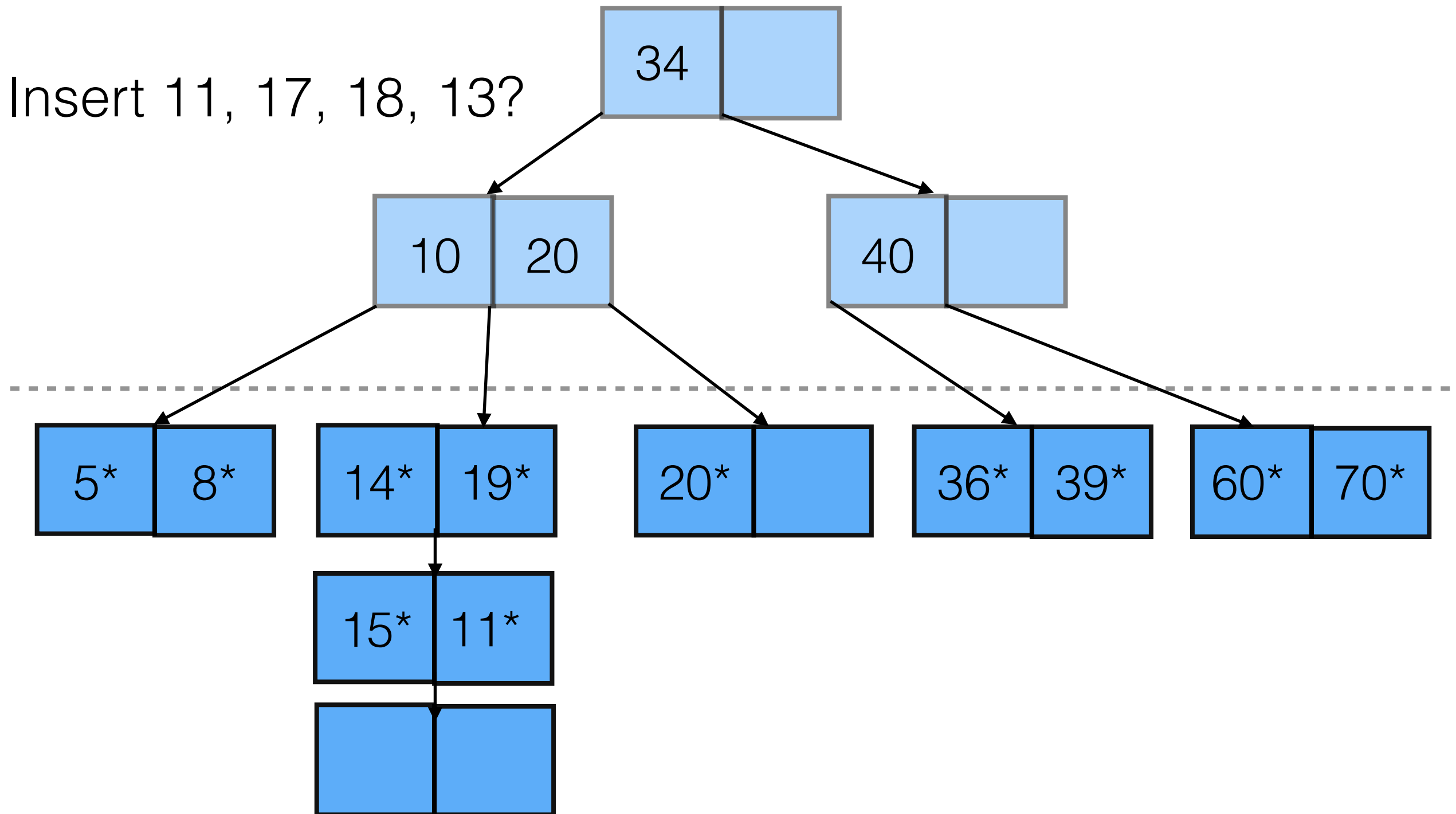
# ISAM

Insert 11, 17, 18, 13?



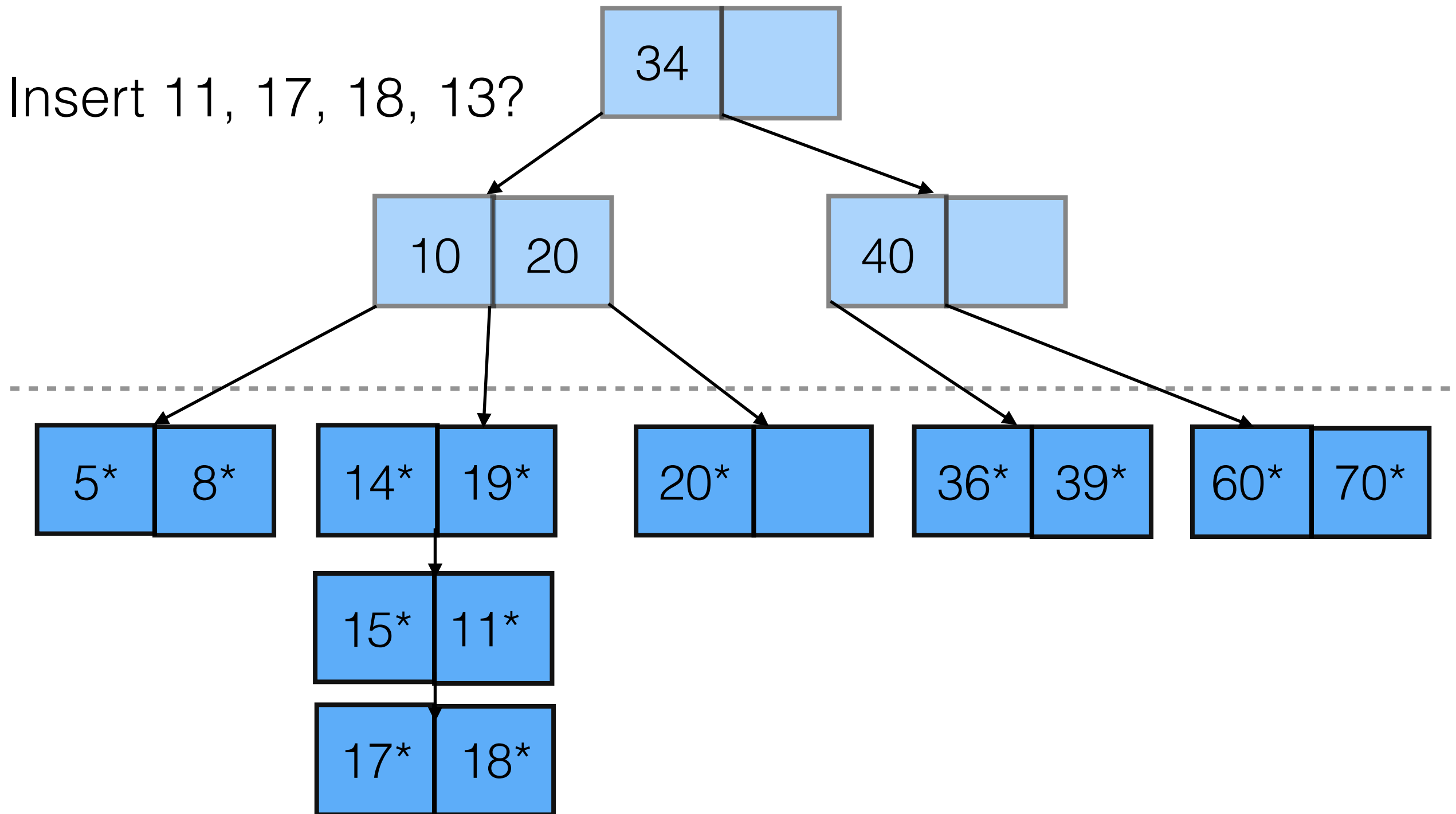
# ISAM

Insert 11, 17, 18, 13?



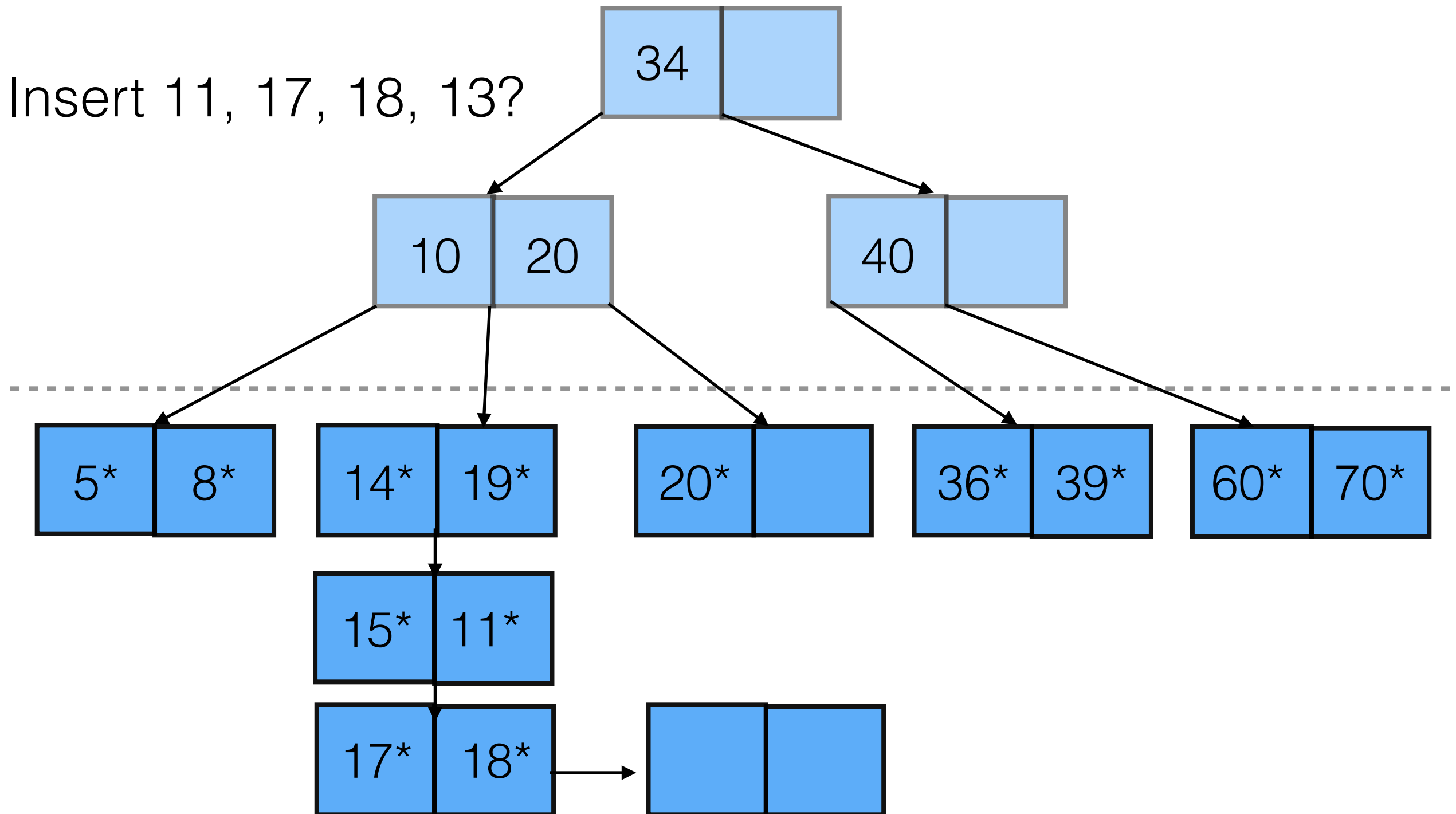
# ISAM

Insert 11, 17, 18, 13?



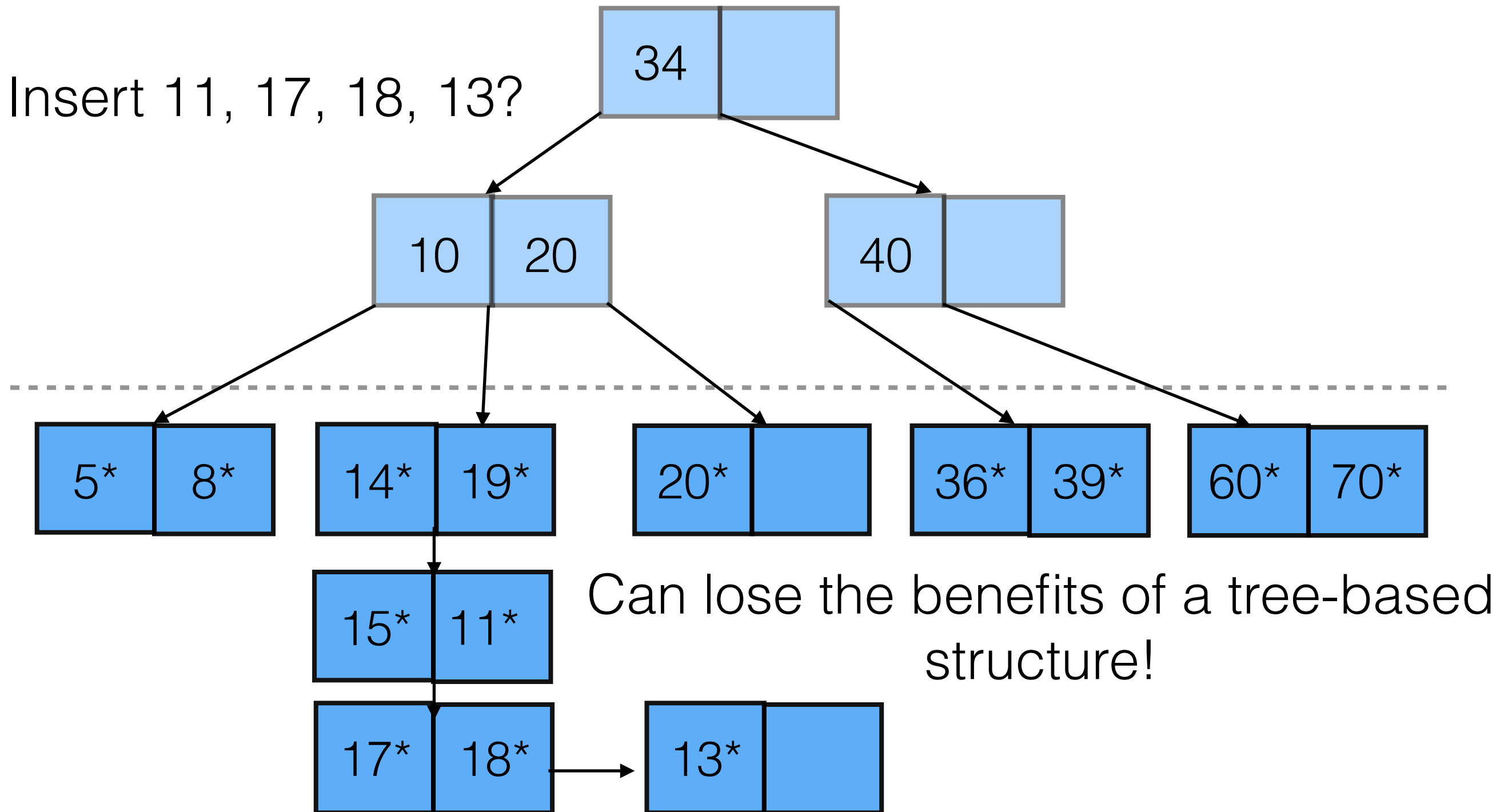
# ISAM

Insert 11, 17, 18, 13?



# ISAM

Insert 11, 17, 18, 13?



# B+ Trees

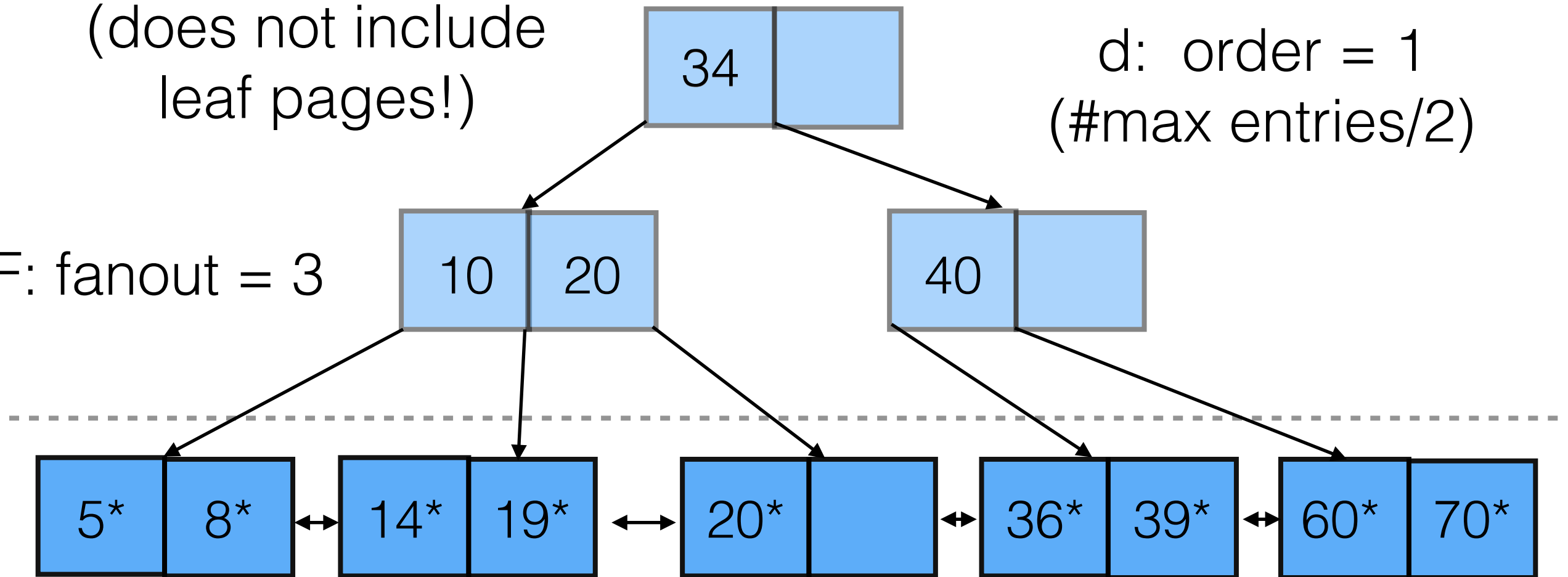
- Dynamic structure to keep tree height-balanced
- Adjusts under inserts and deletes
- Maintain minimum 50% occupancy for each page (except root)

# B+ Trees

H: height = 2  
(does not include  
leaf pages!)

d: order = 1  
(#max entries/2)

F: fanout = 3



N: leaf pages = 5

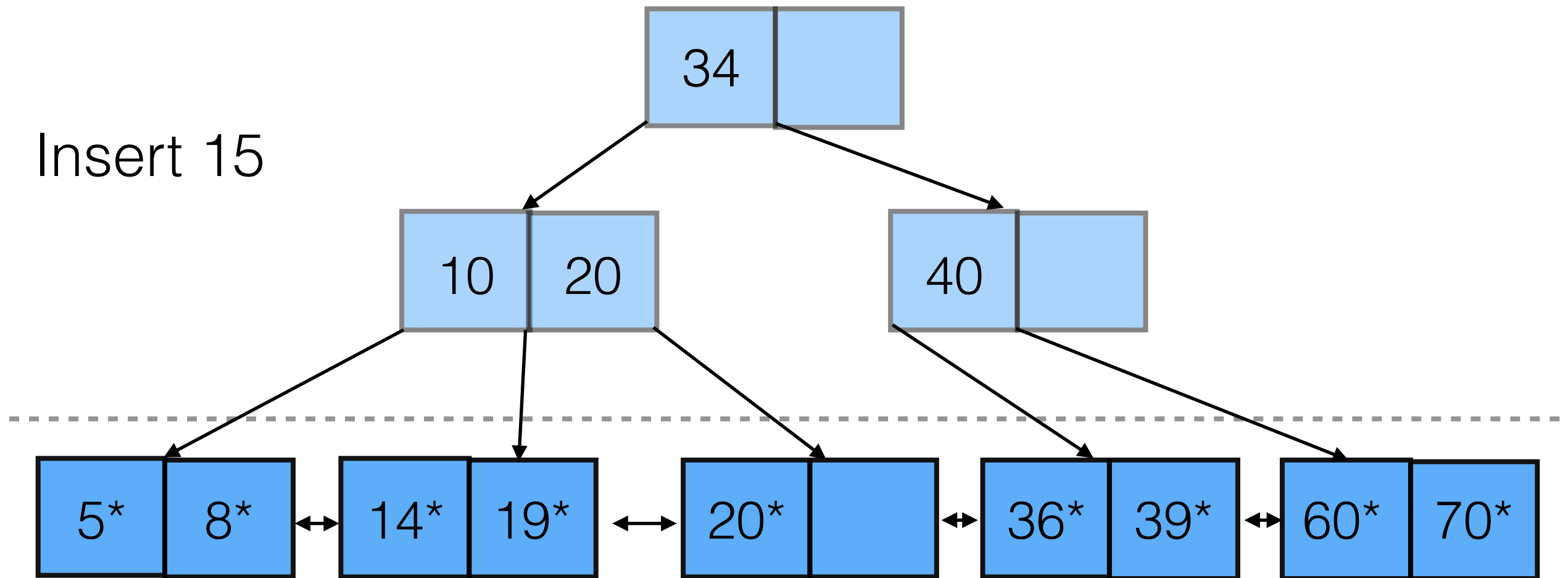
# B+ Trees - Insert X

- Find correct leaf L
- Put X in L
  - If not enough space in L:
    - Split L into L and L2
    - **Copy** up middle key to parent
    - If not enough space in parent:
      - Apply algorithm recursively, except **push** up middle key

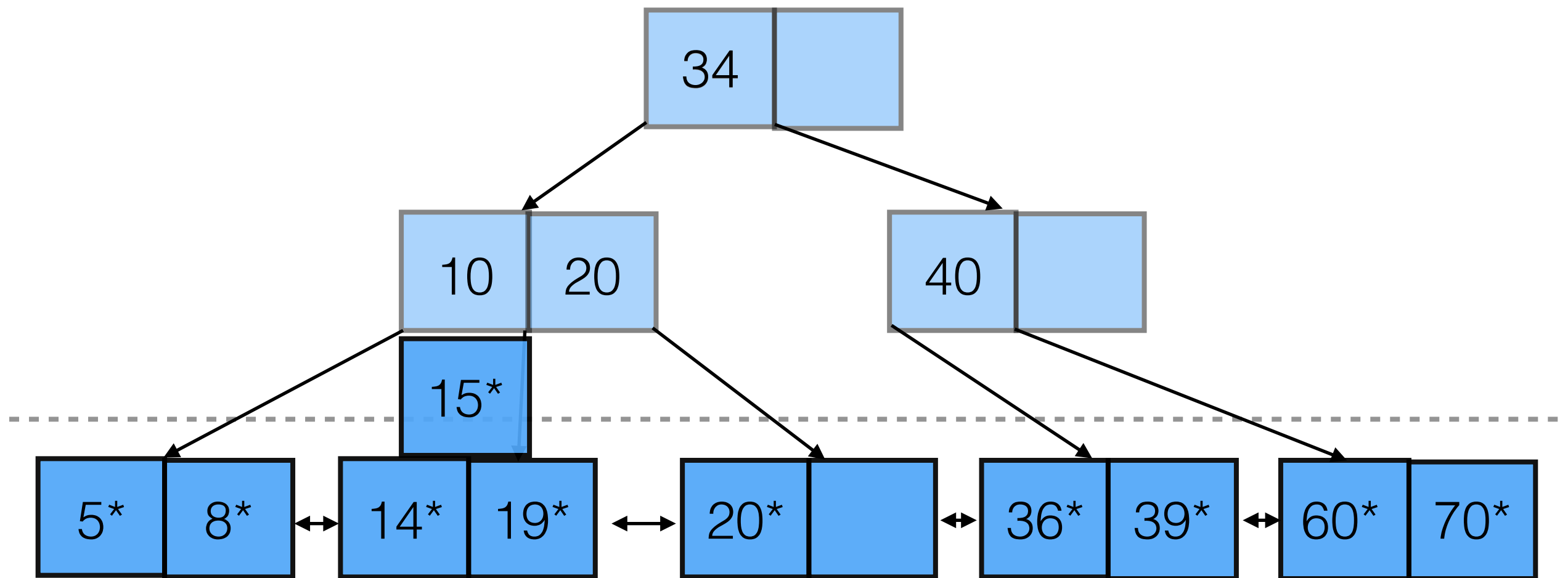


# B+ Trees

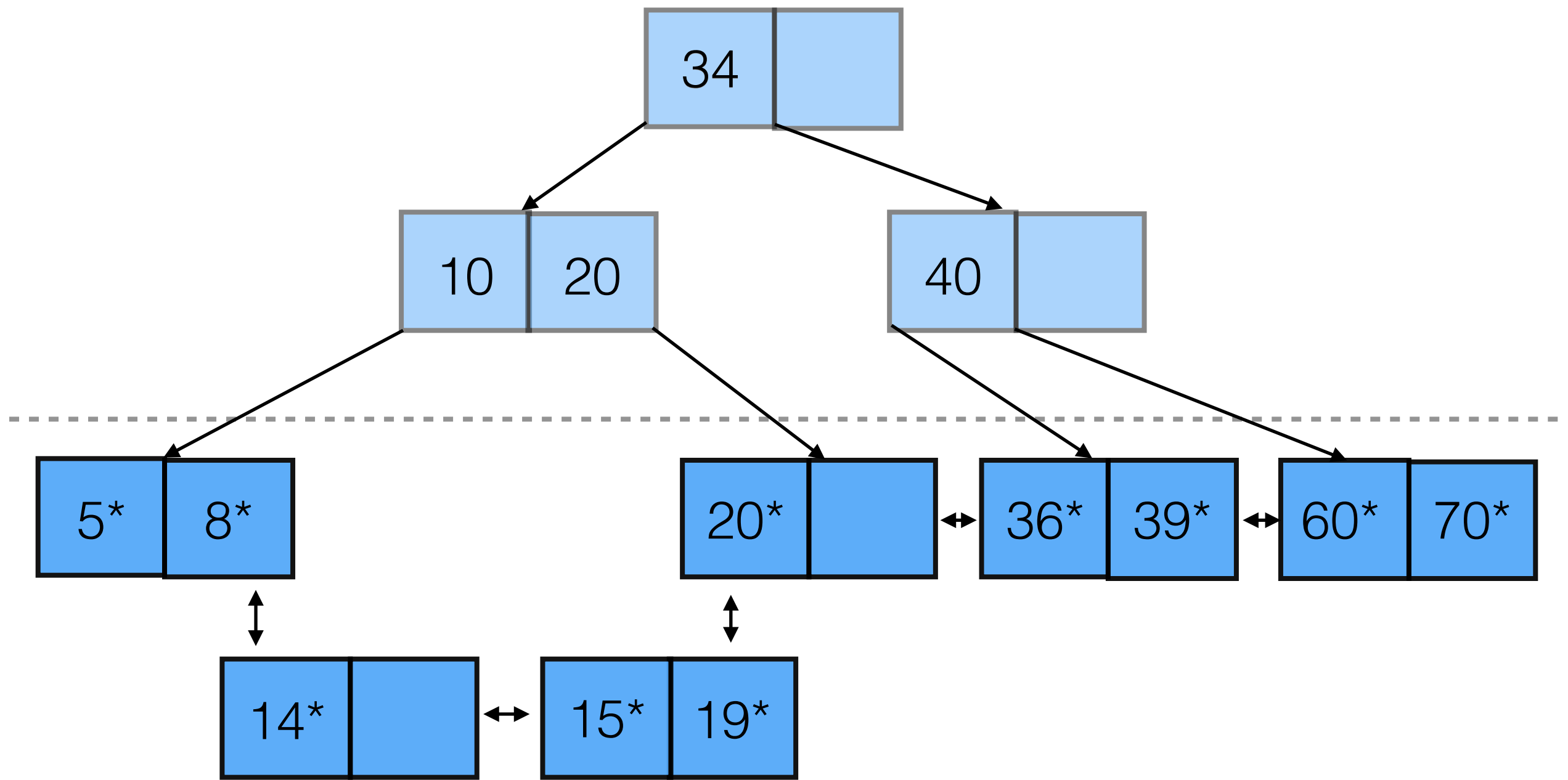
Insert 15



# B+ Trees

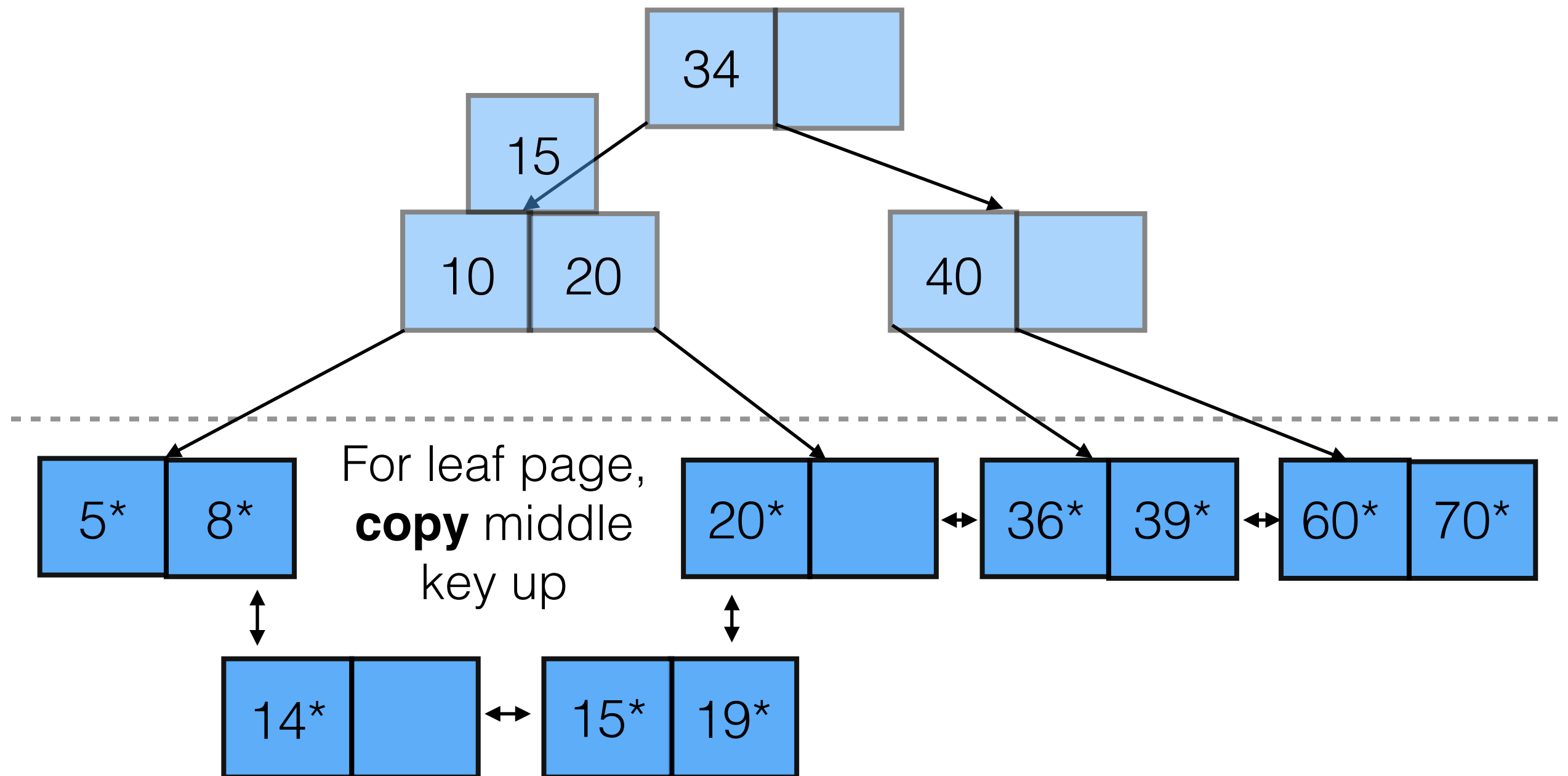


# B+ Trees



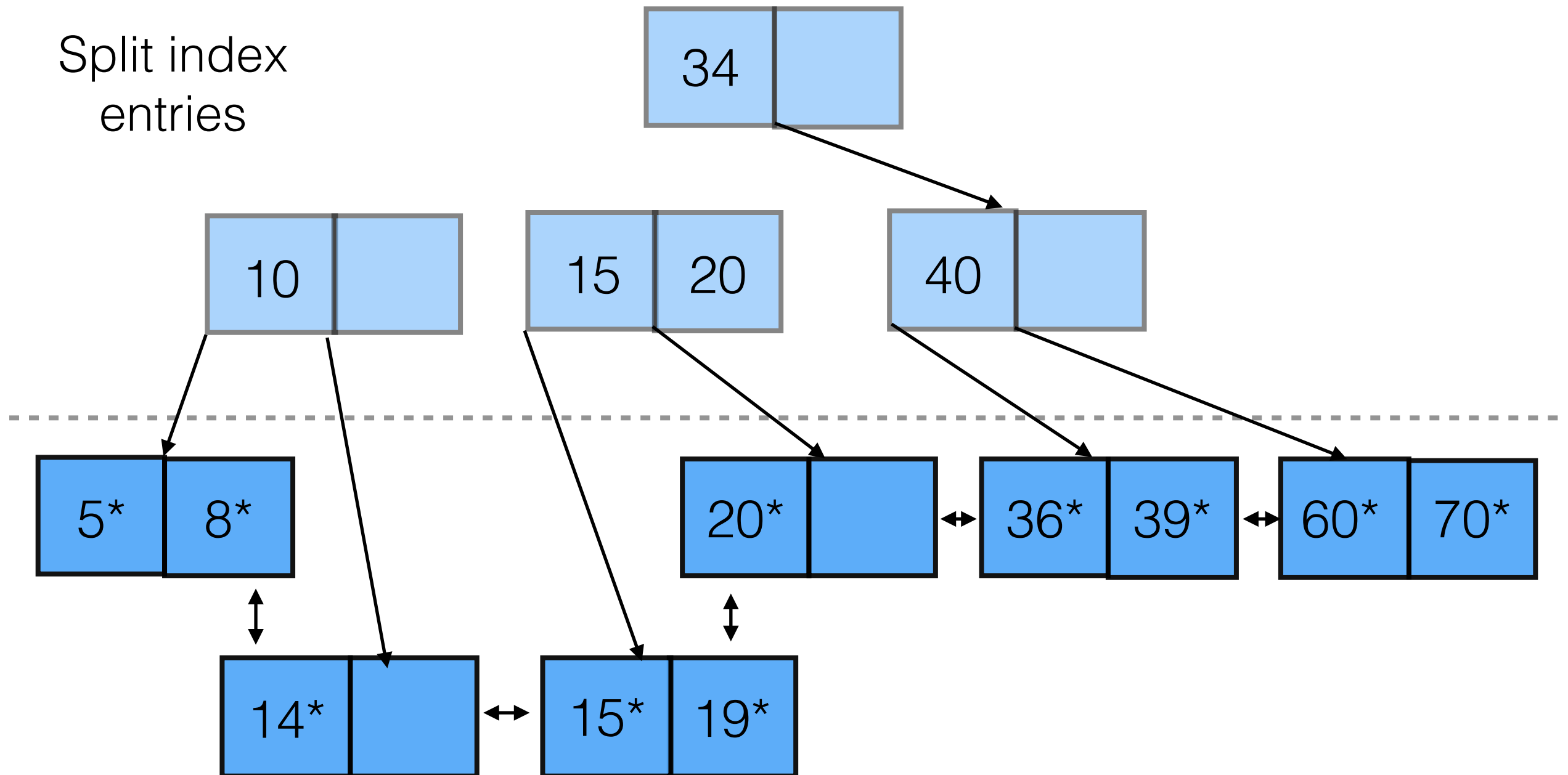
Split the leaf node

# B+ Trees



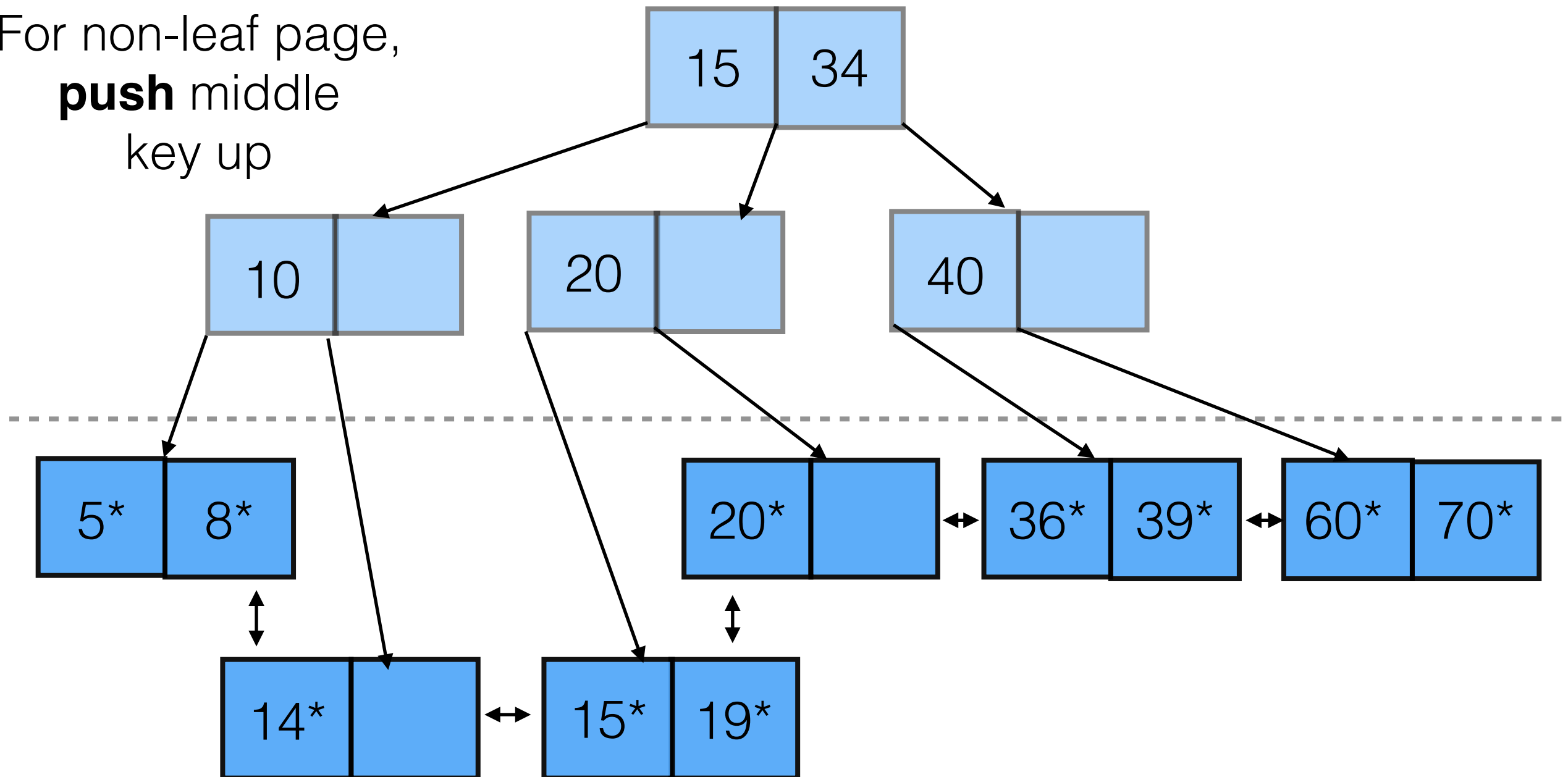
# B+ Trees

Split index  
entries



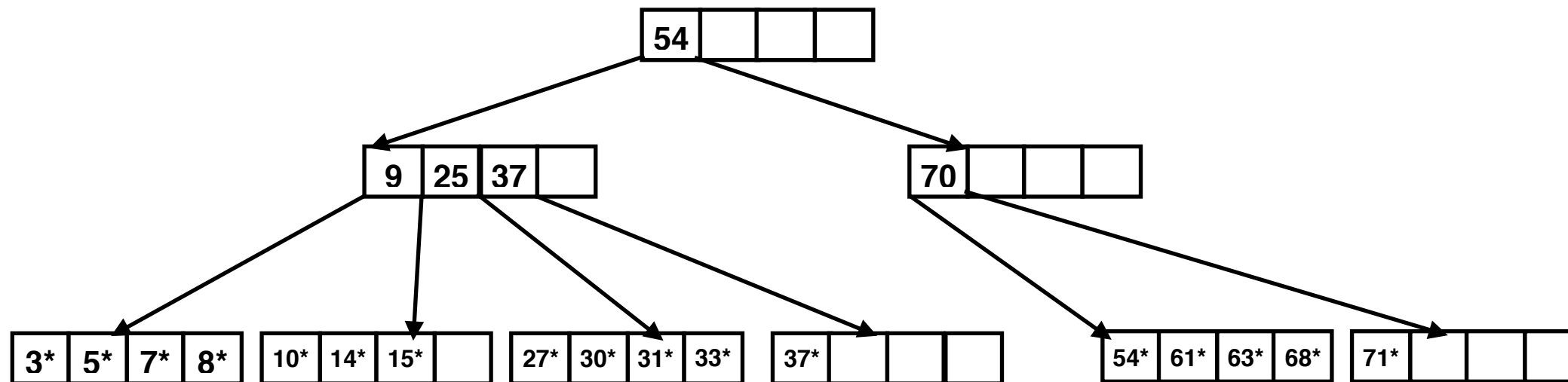
# B+ Trees

For non-leaf page,  
**push** middle  
key up



# Worksheet

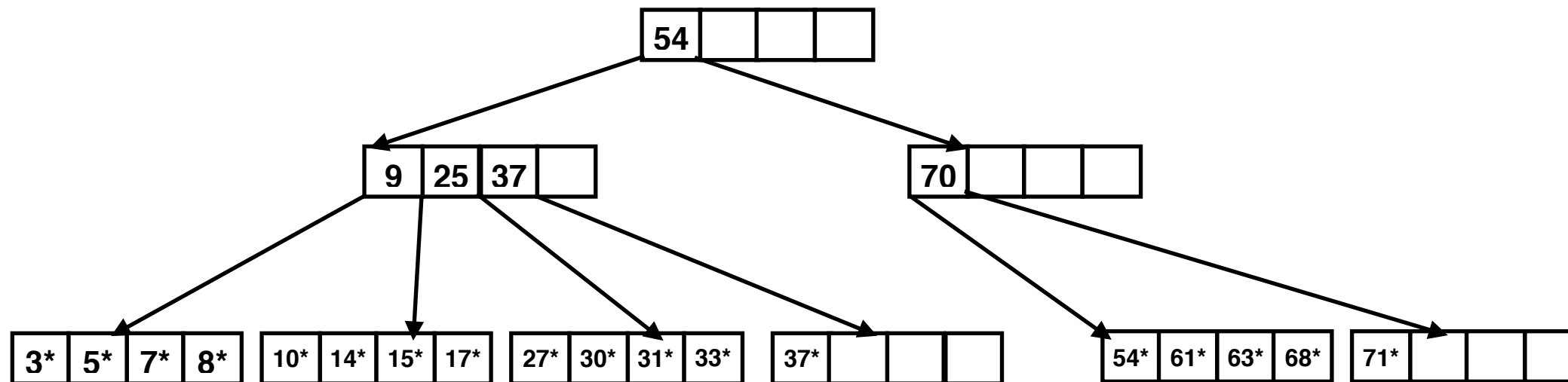
Consider the B+ Tree below and insert the following in order: 17, 18, 29.



Insert 17

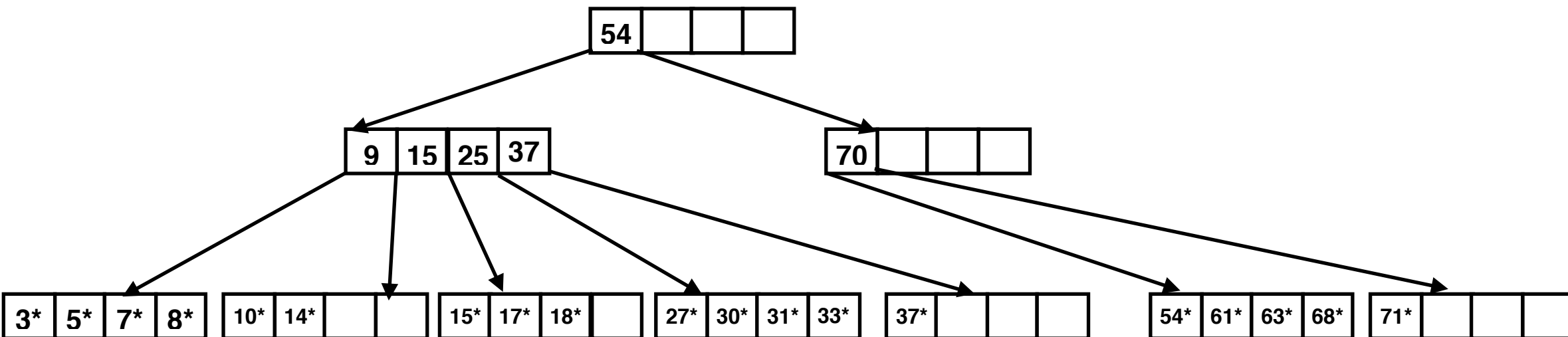


Consider the B+ Tree below and insert the following in order: 17, 18, 29.



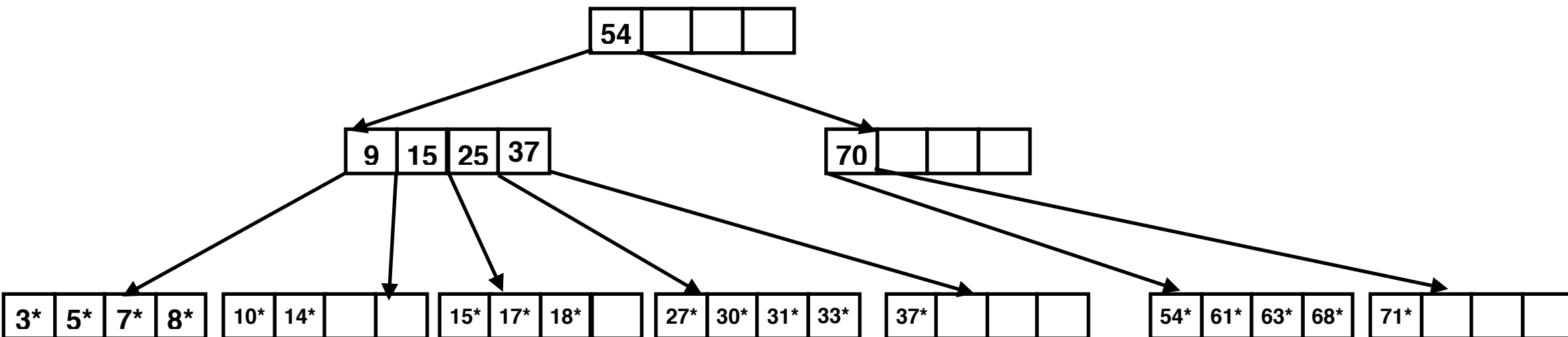
Insert 18

Consider the B+ Tree below and insert the following in order: 17, 18, 29.



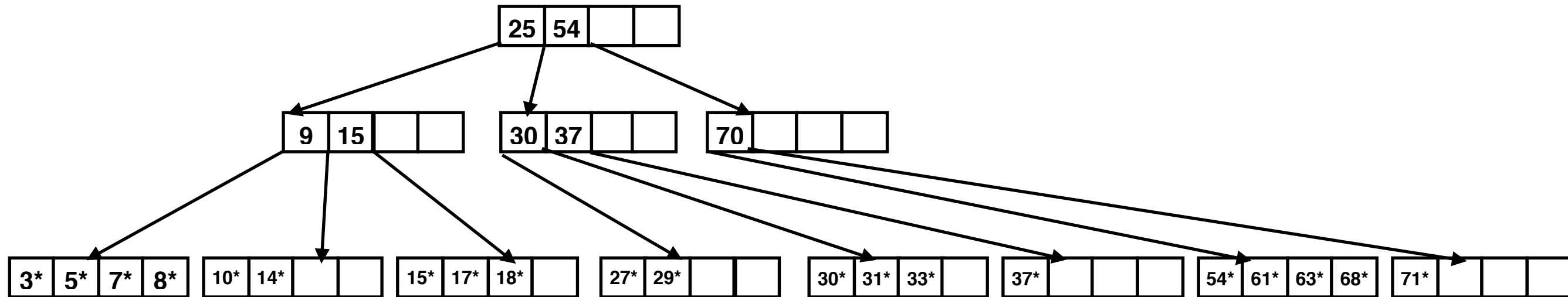
Insert 18

Consider the B+ Tree below and insert the following in order: 17, 18, 29.



Insert 29

Consider the B+ Tree below and insert the following in order: 17, 18, 29.



Insert 29

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many leaf pages are required?

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many leaf pages are required?

Records/page?

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many leaf pages are required?

$$\text{Records/page} = 16384 * .67 / 200 = \sim 54$$

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many leaf pages are required?

$$\begin{aligned}\text{Records/page} &= 16384 * .67 / 200 = \sim 54 \\ 10^9 / 54 &= \sim 18.5 * 10^6 \text{ pages.}\end{aligned}$$



We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

Assume each index entry takes 32 bytes. What is approximately the maximum fanout of the index?

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

Assume each index entry takes 32 bytes. What is approximately the maximum fanout of the index?

$$\begin{aligned} F &= \text{index entries/page} \\ &= (\text{bytes/page}) / (\text{bytes/index\_entry}) \\ &= (16384 * 0.67) / (32) = \sim 343 \end{aligned}$$

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

What is the height (# levels of non-leaf nodes) of the tree?

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

What is the height (# levels of non-leaf nodes) of the tree?

$$\begin{aligned} \text{Height} &= \log_F(\text{leaf pages}) \\ &= \log_{343}(18.5 * 10^6) = \sim 3. \end{aligned}$$

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many I/O operations are required to insert a new record

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many I/O operations are required to insert a new record

4 Reads (3 non-leaf reads, 1 leaf read) + 1 Write  
= 5 I/O's.

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many pages are required to store the non-leaf nodes?

We are using a B+ tree with alternative 1 to store one billion records. Each record is 200 bytes, each disk page has 16kB (16,384 Bytes) and will always be at most 67% full.

How many pages are required to store the non-leaf nodes?

$$1 + 343 + 343^2 = 117993$$



You have decided to create a clustered B+ Tree on the age field. The tree has a fanout of 200 and a height of 3. Assume that you are on average returning 50,000 users per query.

On average, how many I/O's are performed by such a query?

You have decided to create a clustered B+ Tree on the age field. The tree has a fanout of 200 and a height of 3. Assume that you are on average returning 50,000 users per query.

On average, how many I/O's are performed by such a query?

$$\begin{aligned} & \text{cost to traverse tree} + \text{matching pages} \\ &= \text{height} + (\text{size matching user entries} / \text{page size}) \\ &= 3 + (50,000 * 2 / (16 * .67)) \\ &= 9332 \end{aligned}$$

Assume your B+ tree is unclustered. In the worst case, how many I/O's do you need now? Assume that you are still returning 50,000 users per query on average, and that an index entry is 3 times smaller than a user entry.

Assume your B+ tree is unclustered. In the worst case, how many I/O's do you need now? Assume that you are still returning 50,000 users per query on average, and that an index entry is 3 times smaller than a user entry.

**3** I/Os to descend non-leaf index pages

$\text{ceil}(50,000 * 2/3 / (16 * .67)) = \mathbf{3110}$  I/Os to read leaf  
index pages

**50,000** I/Os to read unordered data pages.

$3 + 3110 + 50,000 = 53,113$  I/Os.