# Efficient Query Execution over Large Databases through Semantic Caching of Bitmap Indices

## Sarah McClain and David Chiu (advisor)
### Department of Mathematics and Computer Science, University of Puget Sound

## Background

As popular applications become increasingly data-intensive the need for novel techniques to query large-scale data stores becomes more prevalent. Because computers' hard drives are slow, the more data is stored, the longer it takes to access useful information. For this reason, it is imperative to use efficient data structures such as bitmap indices and caches to provide fast data access.

### What is a bitmap index?

A bitmap index is a way of consolidating a large data set by representing certain columns using a single bit (0=false, or 1=true) for each row. The data is then represented by expanding the columns, or *attributes*, into what are called *bins* and represent the value of the row, or *tuple*, by a 0 or 1 depending on whether it falls in a bin.

| Tuple | Name | Sex | Age |
|-------|------|-----|-----|
| t1 | Kelsie | F | 9 |
| t2 | Megan | F | 26 |
| t3 | Nalin | M | 21 |
| t4 | Zac | M | 80 |
| ... | ... | ... | ... |

Sample Database

| Tuple | Name | | | | Sex | | Age | | |
|-------|------|---|---|---|-----|---|-----|---|---|
| | Kelsie | Megan | Nalin | Zac | M | F | 0-18 | 19-50 | 51-99 |
| t1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| t2 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| t3 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 |
| t4 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 |
| ... | ... | ... | ... | ... | | | | | |

Corresponding Bitmap Index

### What is a semantic cache?

Due to limited space on the (CPU), most data is stored in secondary storage, such as the disk, which is very slow. A semantic cache is a form of architecture that caches the result of selection queries after being fetched from the disk. Semantic caching involves splitting a query into two pieces, a *probe query* and a *remainder query*. The probe query addresses cached regions that contribute to the new query. The remainder query refers to data not contained in the cache (but relevant to the query) that much be fetched from secondary memory. Semantic caching allows the computer to reduce access to the disk by reusing partial matches in the cache.

## Mechanism

**Select \* from users where age >= 15 AND age < 75**

The bins for age correspond to columns 20 through 26 for this bitmap. The SQL query is looking to retrieve columns 21 through 25 of the bitmap and execute an "OR" operation to get all the users between ages 15 and 75. The cache contains already columns [21,23].

| Column # | 20 | 21 | 22 | 23 | 24 | 25 | 26 | ... |
|----------|----|----|----|----|----|----|----|-----|
| ... | | | | Age (years) | | | | ... |
| ... | <15 | 15-29 | 30-44 | 45-59 | 60-74 | 75-99 | >99 | ... |
| ... | 1 | 1 | 0 | 1 | 1 | 1 | 1 | ... |
| ... | 0 | 1 | 0 | 1 | 0 | 1 | 0 | ... |
| ... | 0 | 0 | 1 | 0 | 0 | 0 | 0 | ... |
| ... | 1 | 1 | 1 | 1 | 0 | 1 | 1 | ... |

**1** Probe Query
**Select \* from users age >= 15 AND age < 60**

Remainder Query **3**
**Select \* from users where age >= 60 AND age < 75**



### Cache

cacheLookup(R,21,25)

R = getData(21,23)
cacheLookup(R,24,25)

Return $end_Q = 23$

## Query File

### What is an "OR" operator?

An "OR" operator (denoted by the "|" symbol) is a Boolean operator that returns true (1=true, 0=false) if either (or both) operands is true

OR operation

### What is an "AND" operator?

An "AND" operator (denoted by the "&" symbol) is a Boolean operator that only returns true if both operands are true.

AND operation

### What is a point query?

A point query is searching for a single column in the bitmap. We represent point queries as such:

[0,0]

### What is a range query?

A range query is retrieving data for multiple columns and computing an "or" operation between them.

[0,2]

## Algorithm

```
Algorithm 1 cacheLookup(R, startQ, endQ)
1:  Input
2:     R            Set of result vectors
3:     startQ       the starting column of full query
4:     endQ         the ending column of full query
5:  Output
6:     R*           Result vectors for probe query [startQ, end'Q]
7:     end'Q        ending column of probe query
8:  i ← hashCode(startQ)
9:  best ← list[i].head
10: for all [startQ, end'Q] in list[i] do
11:     if endQ = end'Q then
12:         R ← getData([startQ, endQ])
13:         R* ← R ∪ R*
14:         return endQ
15:     if best.difference > |calculateDifference(end'Q, endQ)| then
16:         best ← [startQ, end'Q]
17:         best.difference = |calculateDifference(end'Q, endQ)|
18: if end'Q ≥ endQ then
19:     R ← getData([startQ, end'Q])
20:     R* ← R ∪ R*
21:     return end'Q
22: return cacheLookup(R, end'Q + 1, endQ)
```
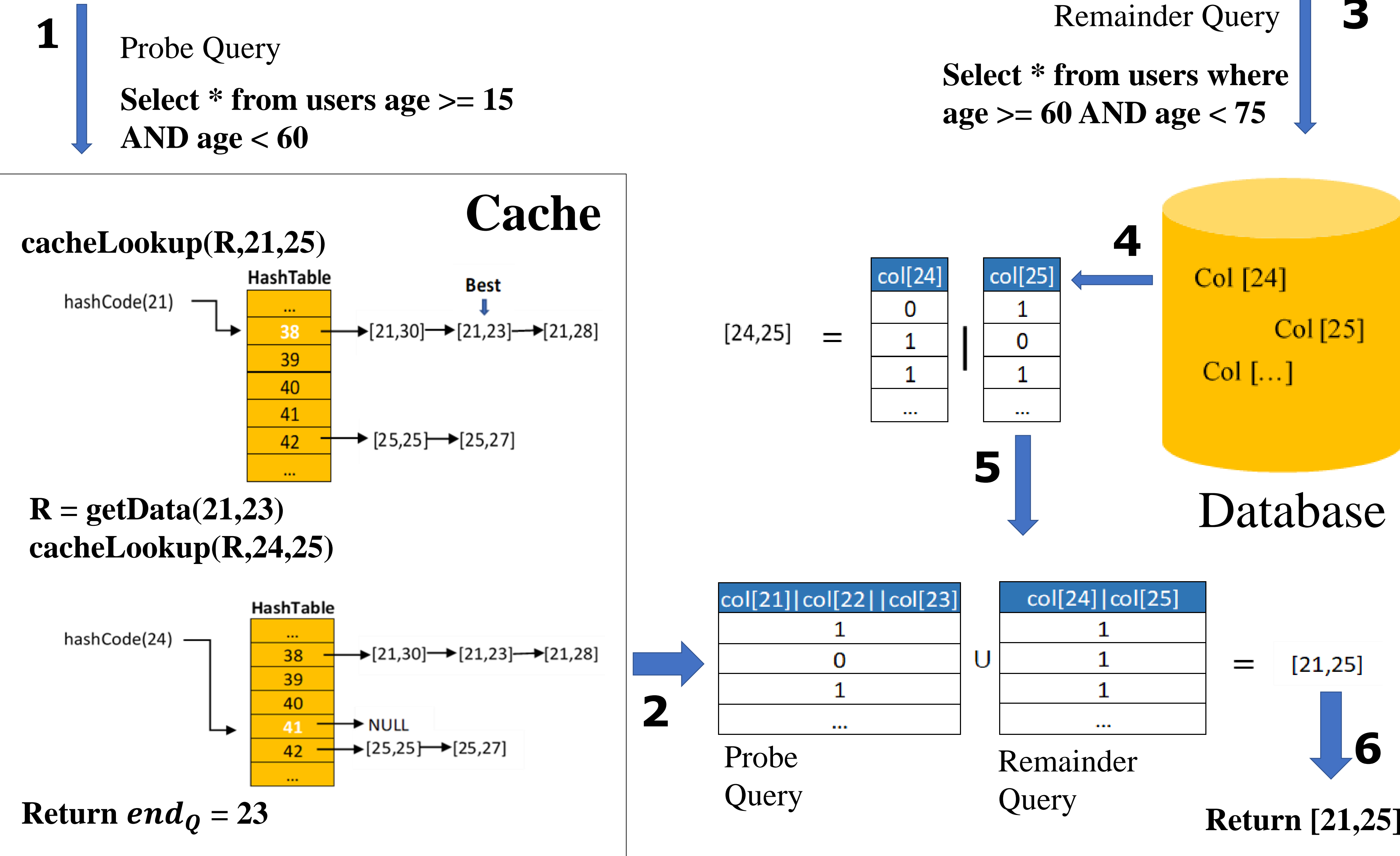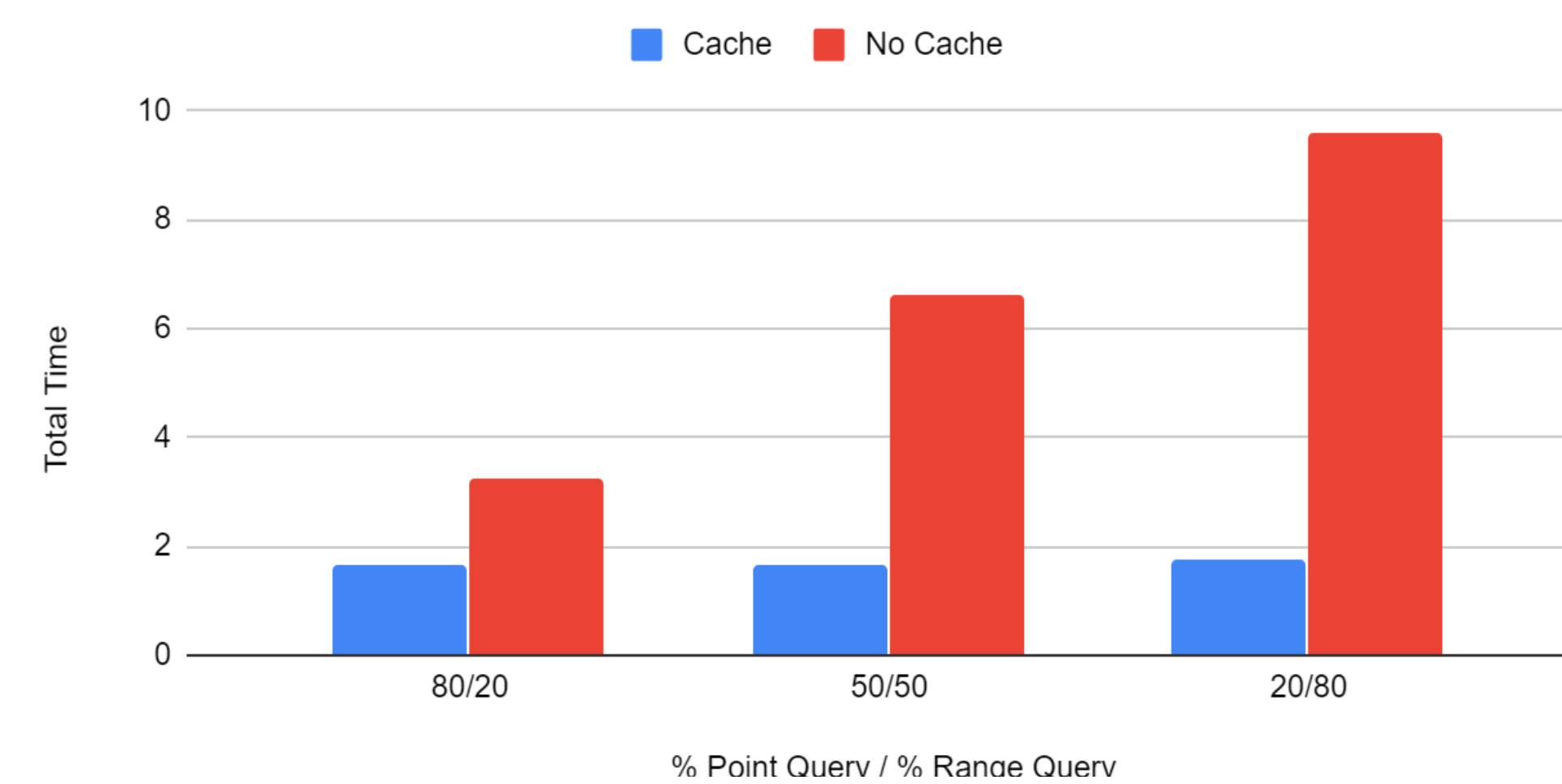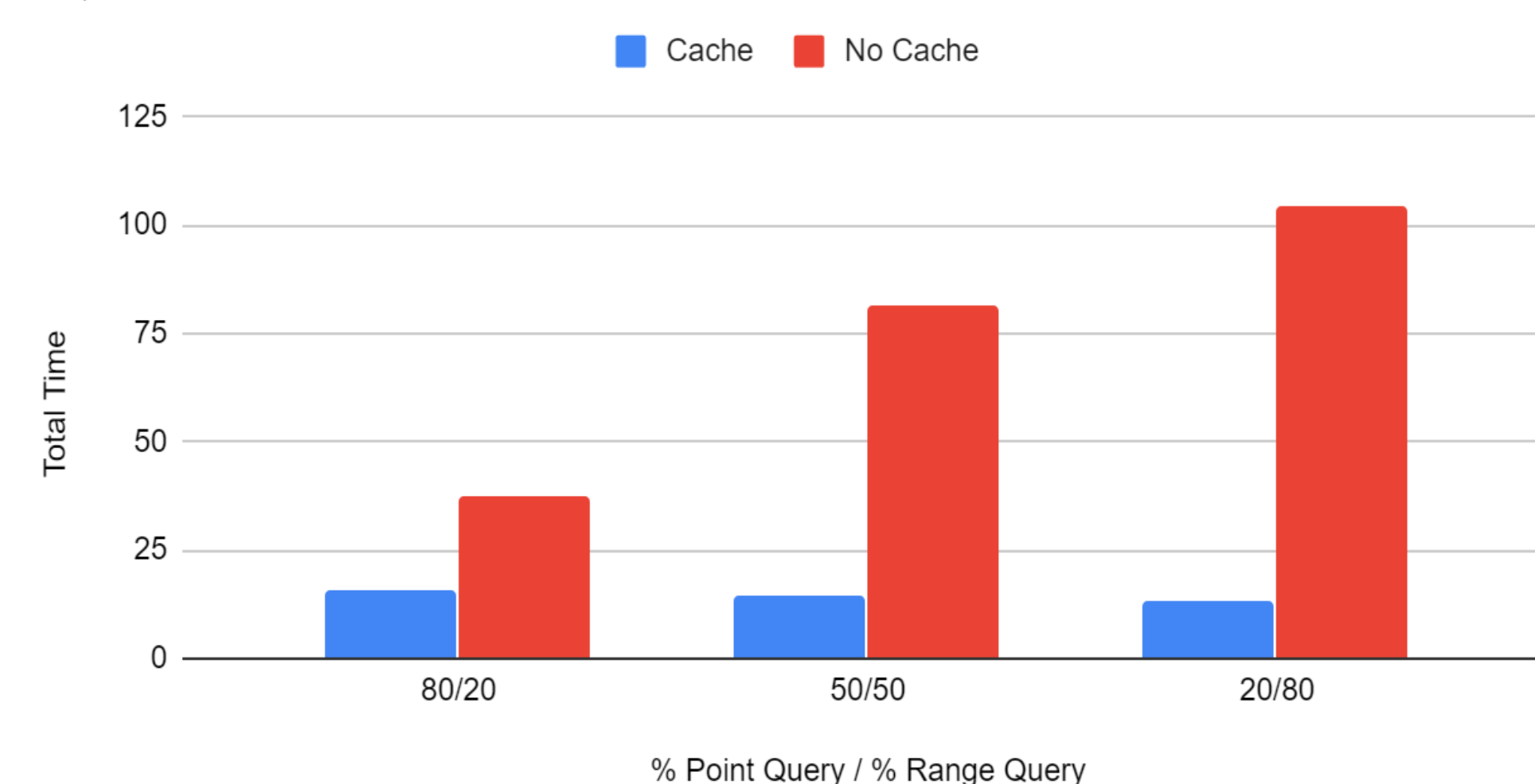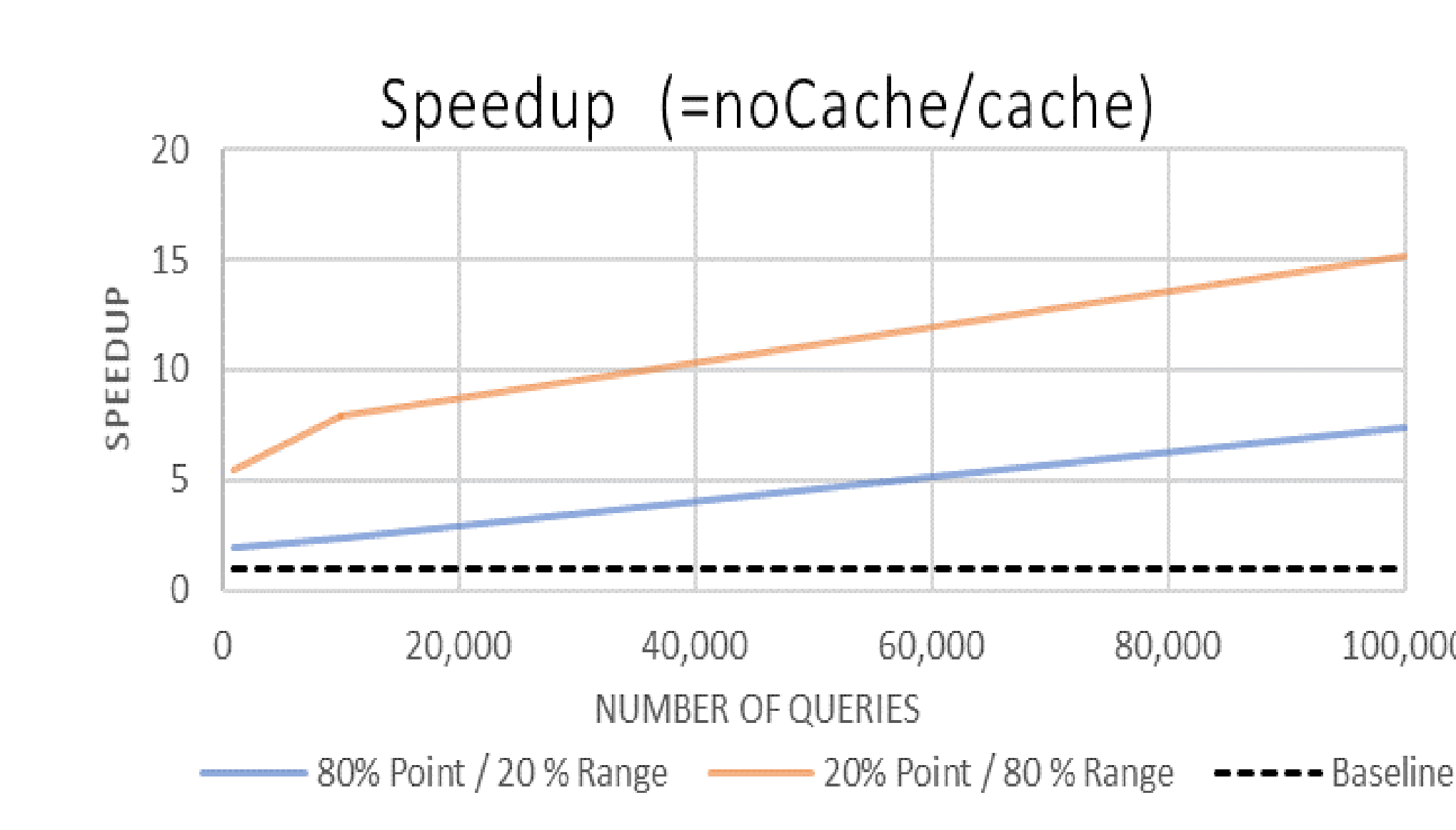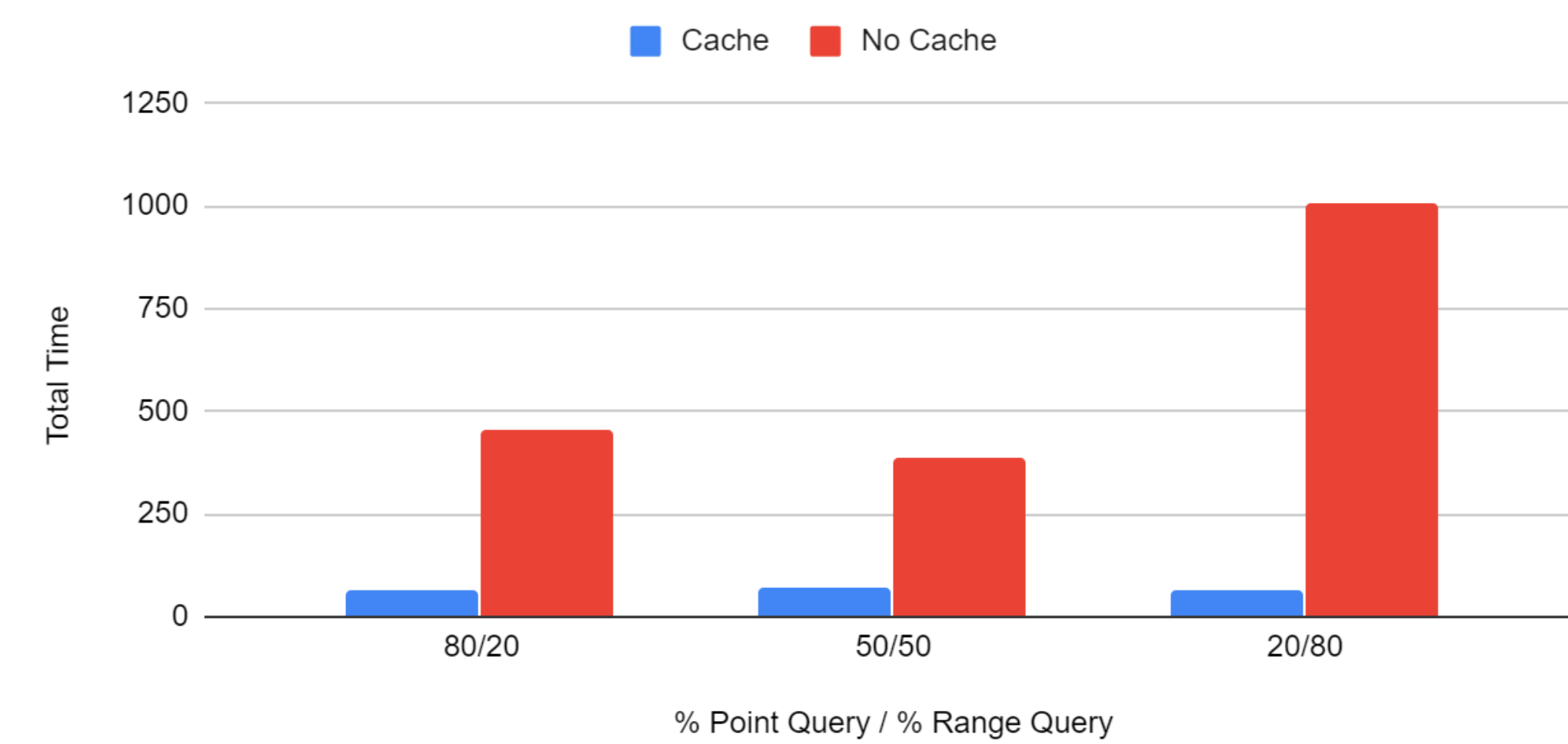
## Plots

**Total Time for Executing Query File With Cache and Without Cache**
1,000 Queries



**Total Time for Executing Query File With Cache and Without Cache**
100,000 Queries



**Total Time for Executing Query File With Cache and Without Cache**
10,000 Queries





Speedup (=noCache/cache)

## Results

The semantic cache tremendously decreased the amount of time needed to query the database. As shown in the figures, as the number of queries grows, it becomes more efficient to leverage this caching architecture. The semantic cache also has a greater speedup as the number of range queries increases. There is still more to be studied about semantic caching. More work on the effect of different lookup algorithms and caching compound ranges could increase efficiency even more.

## Acknowledgements