

Data Structures and Algorithms

(ESO207)

Lecture 40

- Search data structure for integers : **Hashing**

Data structures for searching

in $O(1)$ time

Motivating Example

Input: a given set S of 1009 positive integers

Aim: Data structure for searching

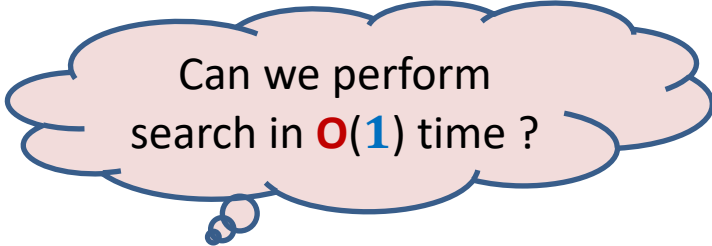
Example

```
{  
123, 579236, 1072664, 770832456778, 61784523, 100004503210023, 19,  
762354723763099, 579, 72664, 977083245677001238, 84, 100004503210023,  
...  
}
```

Data structure : Array storing S in sorted order

Searching : Binary search

$O(\log |S|)$ time



Can we perform
search in $O(1)$ time ?

Problem Description

U :

$S \subseteq U$,

$n = |S|$,

$n \ll m$

A search query:

Aim: A data structure for a given set S

that can facilitate search in $\mathbf{O}(1)$ time

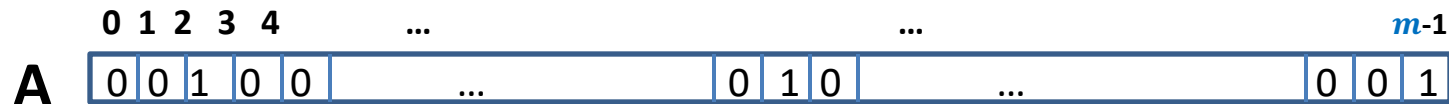
A trivial data structure for $O(1)$ search time

Build a 0-1 array **A** of size m such that

$A[i] = 1$ if $i \in S$.

$A[i] = 0$ if $i \notin S$.

Time complexity for searching an element in set S : $O(1)$.



This is a totally Impractical data structure because $n \ll m$!

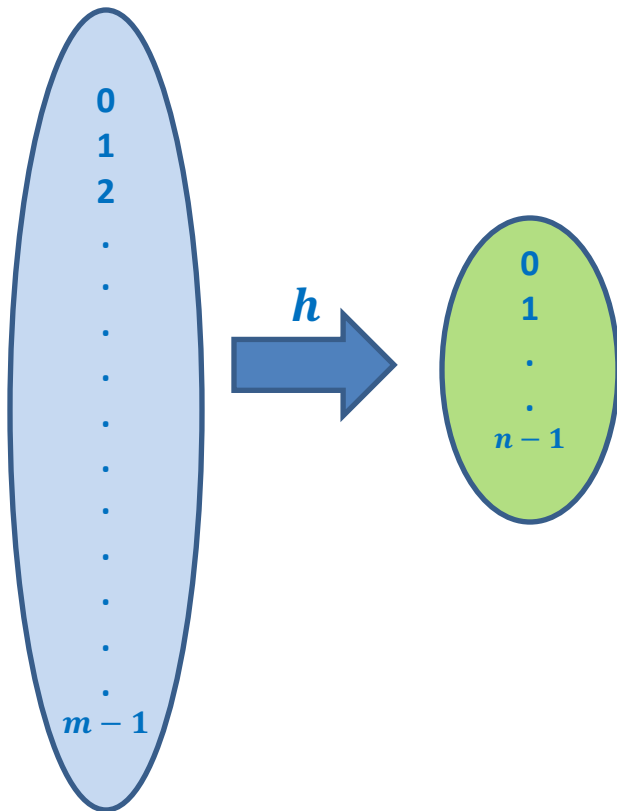
Example: n = few thousands, m = few trillions.

Question:

Can we have a data structure of $O(n)$ size that can answer a search query in $O(1)$ time ?

Answer: Hashing

Hash function



Hash function:

h is a mapping from U to $\{0, 1, \dots, n-1\}$ with the following characteristics.

- **Space** required for h
- $h(i)$ computable in \cdot

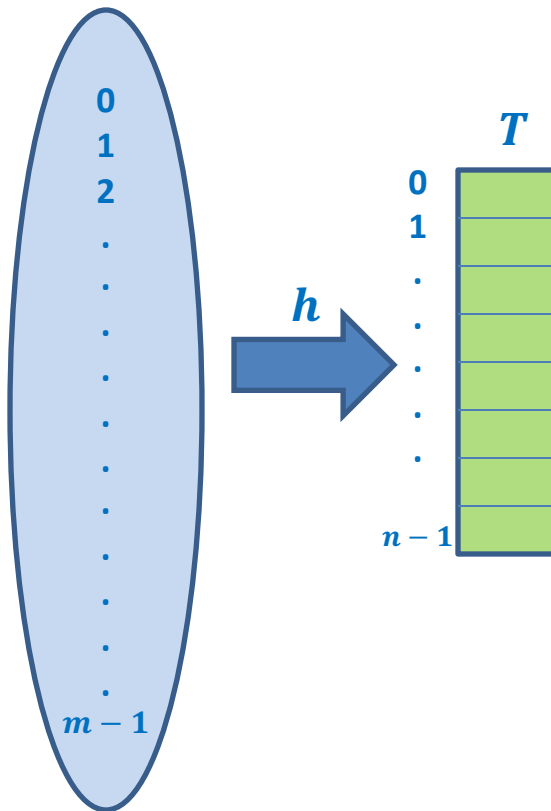
Example:

Hash value:

For a given hash function h , and $i \in U$.

$h(i)$ is called hash value of i

Hash function, hash value



Hash function:

h is a mapping from U to $\{0, 1, \dots, n-1\}$ with the following characteristics.

- **Space** required for h : a few **words**.
- $h(i)$ computable in **$O(1)$** time in **word RAM**.

Example: $h(i) = i \bmod n$

Hash value:

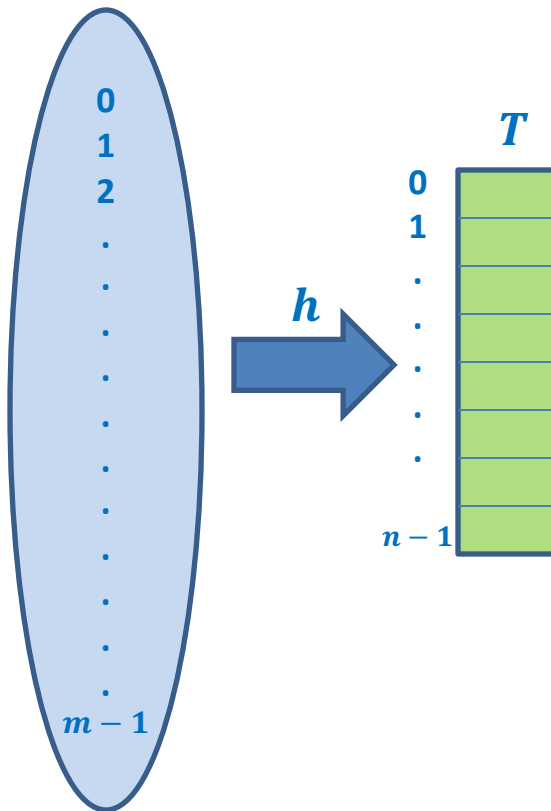
For a given hash function h , and $i \in U$.

$h(i)$ is called hash value of i

Hash Table:

An array $T[0 \dots n-1]$

Hash function, hash value, hash table



Hash function:

h is a mapping from U to $\{0, 1, \dots, n-1\}$ with the following characteristics.

- **Space** required for h : a few **words**.
- $h(i)$ computable in **$O(1)$** time in **word RAM**.

Example: $h(i) = i \bmod n$

Hash value:

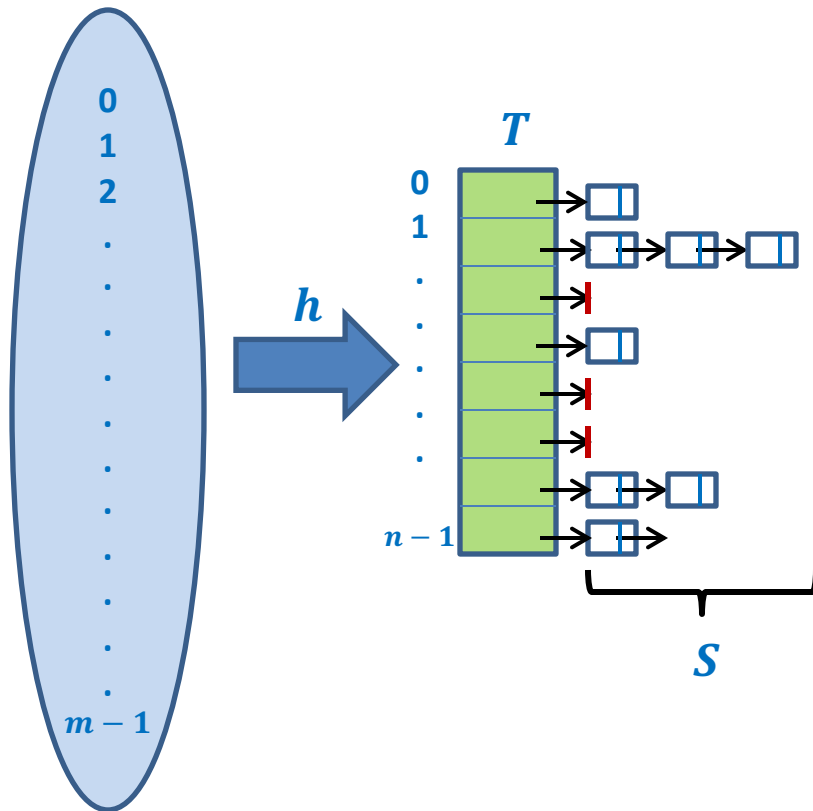
For a given hash function h , and $i \in U$.

$h(i)$ is called hash value of i

Hash Table:

An array $T[0 \dots n-1]$

Hash function, hash value, hash table



Hash function:

h is a mapping from U to $\{0, 1, \dots, n-1\}$ with the following characteristics.

- **Space** required for h : a few **words**.
- $h(i)$ computable in $\mathbf{O}(1)$ time in **word RAM**.

Example: $h(i) = i \bmod n$

Hash value:

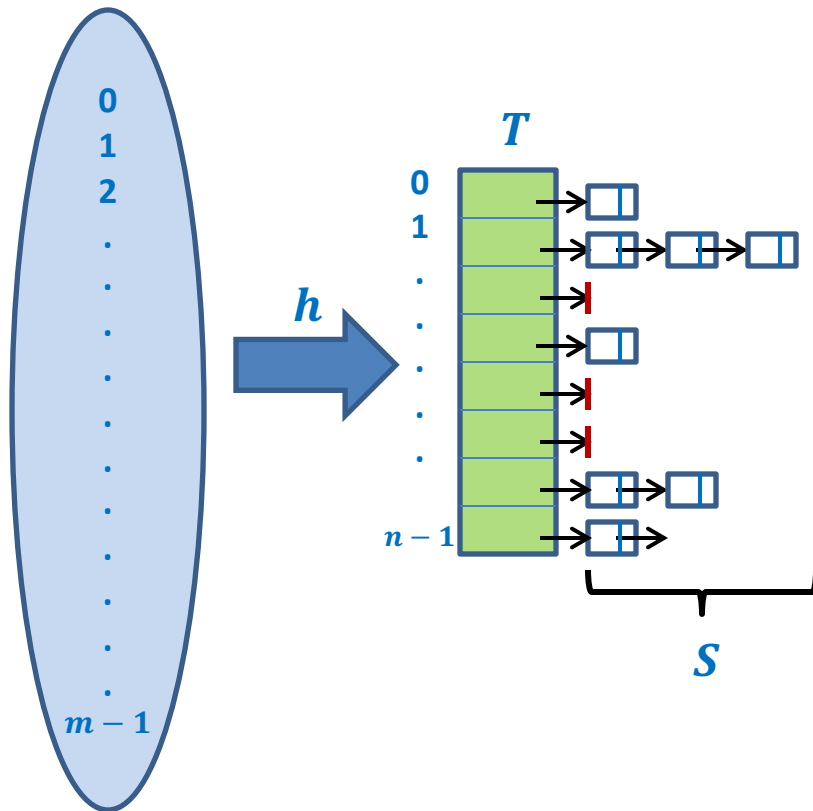
For a given hash function h , and $i \in U$.

$h(i)$ is called hash value of i

Hash Table:

An array $T[0 \dots n-1]$ of pointers storing S .

Hash function, hash value, hash table



Question:

How to use (h, T) for searching an element $i \in U$?

Answer:

$k \leftarrow h(i);$

Search element i in the list $T[k]$.

Time complexity for searching:

$\mathcal{O}(\text{length of the longest list in } T).$

Efficiency of Hashing depends upon hash function

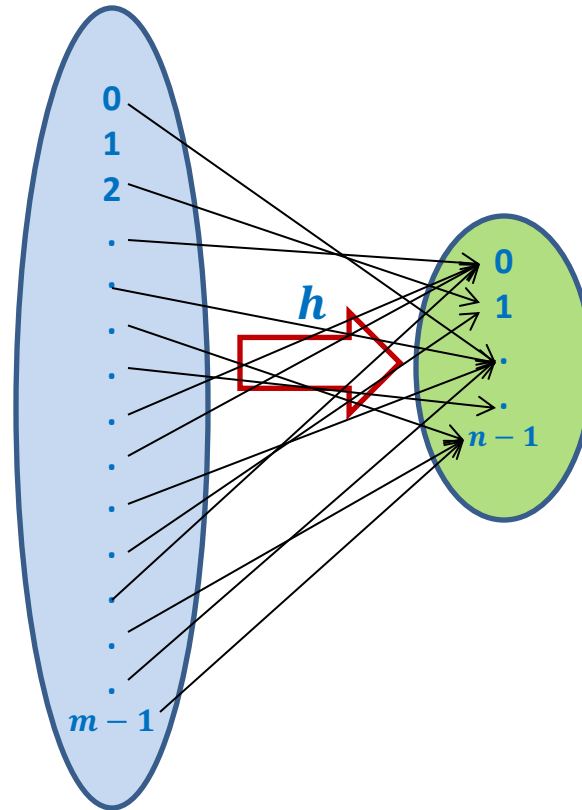
A hash function h is good if it can **evenly** distributes S .

Aim: To search for a good hash function for a given set S .



There **can not be** any hash function h which is good for every S .

Hash function, hash value, hash table



For every h , there exists a subset of $\left\lceil \frac{m}{n} \right\rceil$ elements from U which are hashed to same value under h .

So we can always construct a subset S for which all elements have same hash value

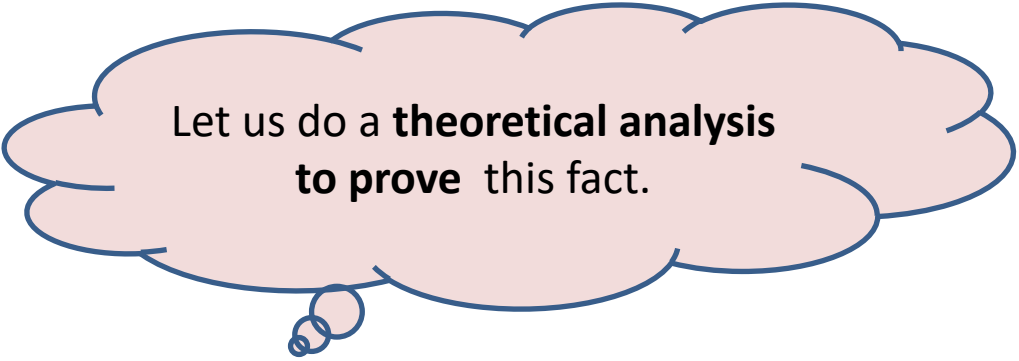
➔ All elements of this set S are present in a single list of the hash table T associated with h .

➔ $O(n)$ worst case search time.

Why does hashing work **so well** in Practice ?

$$h(i) = i \bmod n$$

Because the set **S** is usually a **uniformly random** subset of U .



Let us do a **theoretical analysis**
to prove this fact.

Why does hashing work so well in Practice ?

$$h(x) = x \bmod n$$

Let y_1, y_2, \dots, y_n denote n elements selected randomly uniformly from U to form S .

Question:

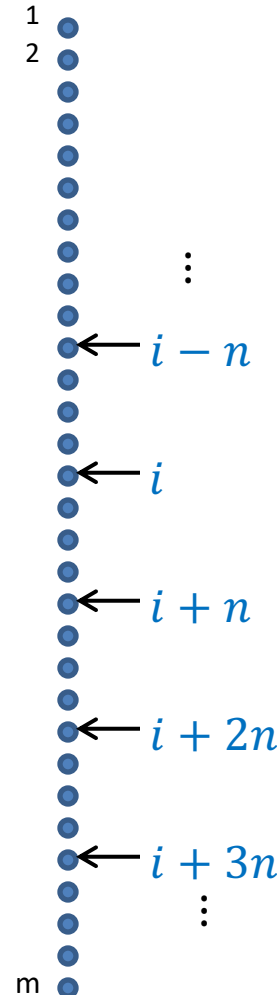
What is expected number of elements of S colliding with y_1 ?

Answer: Let y_1 takes value i .

$P(y_j \text{ collides with } y_1) = ??$

How many possible values can y_j take ? $m - 1$

How many possible values can collide with i ?



Why does hashing work **so well** in Practice ?

$$h(x) = x \bmod n$$

Let y_1, y_2, \dots, y_n denote n elements selected randomly uniformly from U to form S .

Question:

What is expected number of elements of S colliding with y_1 ?

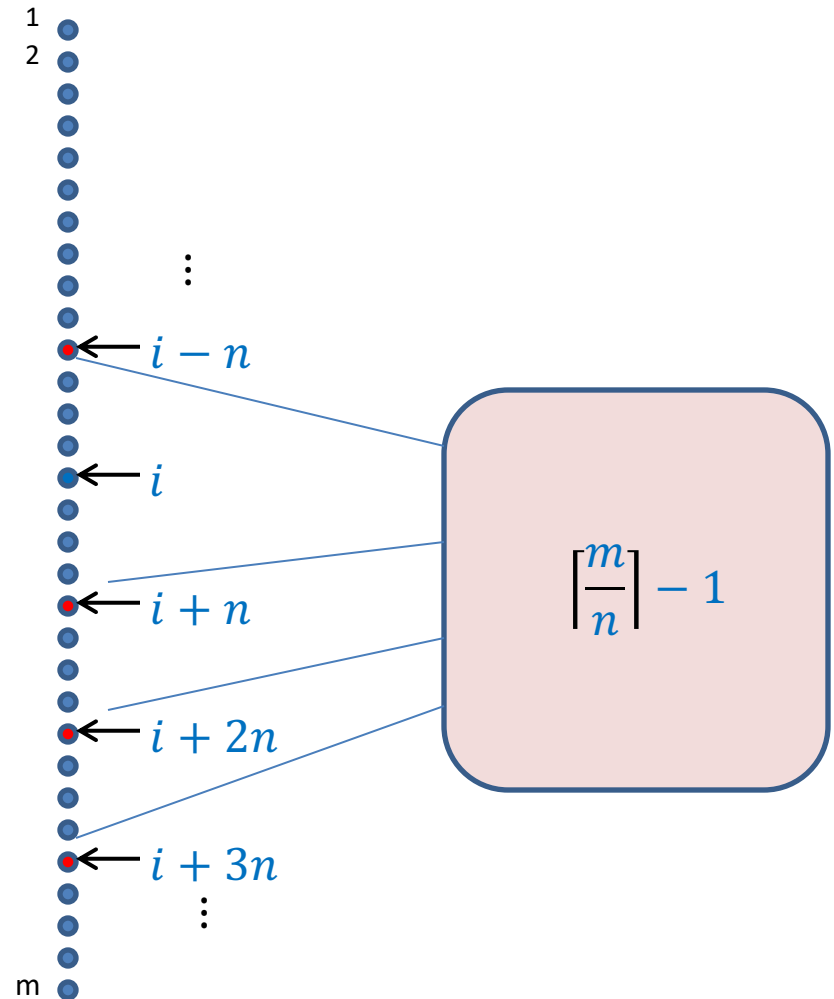
Answer: Let y_1 takes value i .

$P(y_j \text{ collides with } y_1) =$

$$\frac{\left\lceil \frac{m}{n} \right\rceil - 1}{m-1}$$

Expected number of elements of S colliding with $y_1 =$

$$\begin{aligned} &= \frac{\left\lceil \frac{m}{n} \right\rceil - 1}{m-1} (n-1) \\ &= O(1) \end{aligned}$$



Why does hashing work **so well** in Practice ?

Conclusion

1. $h(i) = i \bmod n$ works so well because
for a uniformly random subset of U ,
the **expected** number of collision at an index of T is $O(1)$.

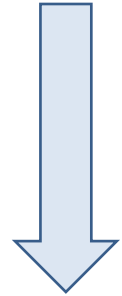
It is easy to fool this hash function such that it achieves $O(s)$ search time.
(do it as a simple exercise).

This makes us think:

“How can we achieve worst case $O(1)$ search time for a given set S .”

Hashing: theory

1953



1984

$U : \{0, 1, \dots, m - 1\}$

$S \subseteq U,$

$n = |S|,$

Theorem [FKS, 1984]:

A hash table and hash function can be computed in average $O(n)$ time for a given S s.t.

Space : $O(n)$

Query time: worst case $O(1)$

Ingredients :

- elementary knowledge of **prime numbers**.
- The algorithms use **simple randomization**.

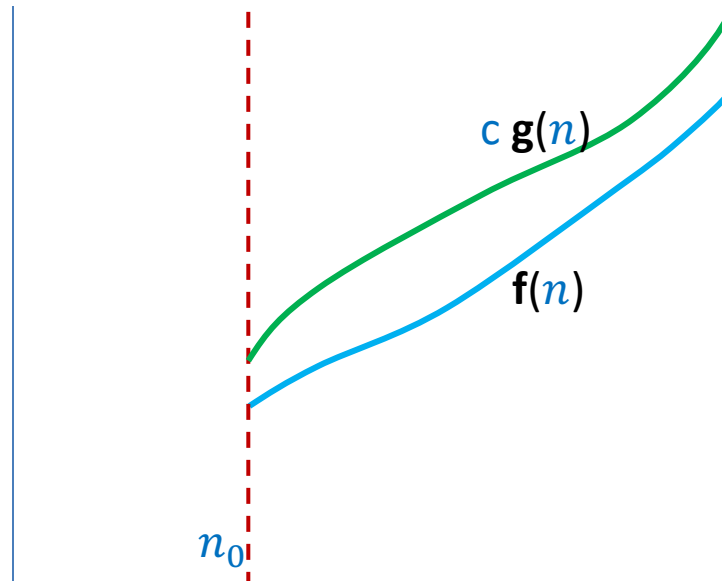
Order notation

Definition: Let $f(n)$ and $g(n)$ be any two increasing functions of n .

$f(n)$ is said to be

if there exist constants c and n_0 such that

$$f(n) \leq c g(n) \quad \text{for all } n > n_0$$



$$f(n) = O(g(n))$$

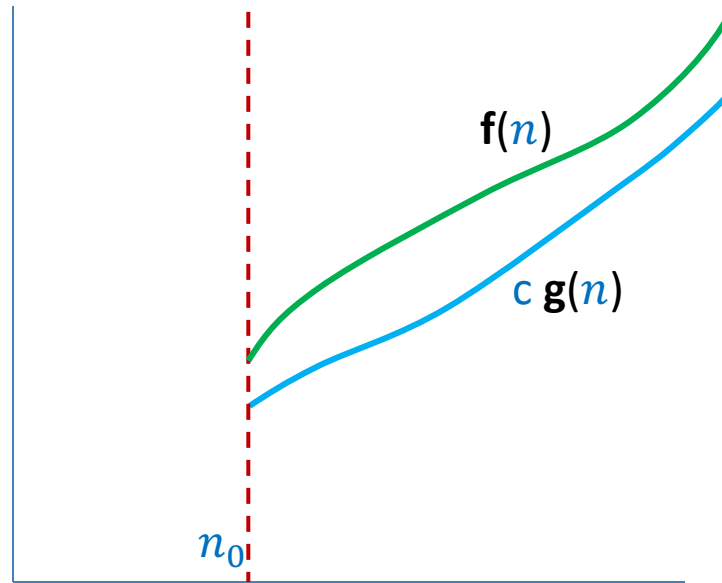
Order notation extended

Definition: Let $f(n)$ and $g(n)$ be any two increasing functions of n .

$f(n)$ is said to be

if there exist constants c and n_0 such that

$$f(n) \geq c g(n) \quad \text{for all } n > n_0$$



$$\frac{n^2}{100} = \Omega(10000 n \log n)$$

$$f(n) = \Omega(g(n))$$

Order notation extended

Observations:

- $f(n) = O(g(n))$ if and only if $g(n) = \Omega(f(n))$

One more Notation:

If $f(n) = O(g(n))$ and $g(n) = O(f(n))$, then

$$g(n) = \Theta(f(n))$$

Examples:

- $\frac{n^2}{100} = \Theta(10000 n^2)$

- Time complexity of Quick Sort is $\Omega(n \log n)$
- Time complexity of Merge sort is $\Theta(n \log n)$