

# **Skip List and HNSW**

Big Data Management and Governance

---

Prof. Giovanni Simonini, Giovanni Malaguti

University of Modena and Reggio Emilia

November 13, 2025

## Details of the Lab (1)

- We will implement Skip List and HNSW data structures
- Simple implementations, no optimization and we won't cover all the details
- Goal: to obtain two new data structures that support creation, insert and some kind of search operations
- At home, attempt to implement other operations (e.g., key-value semantics, different distance functions, update, delete) and identify necessary code changes

## Details of the Lab (2)

- Clone the repository  
`https://github.com/Stravanni/bdm.git`
- Inside the cloned repository, navigate to  
"bdm/lab/skiplist-hnsw"
- There are three files: `skip_list.py`, `hnsw.py`, and `utils.py`
- We will implement algorithms and data structures inside the first two files from scratch. In `utils.py` there are helper functions for visualization and experiments

## Details of the Lab (3)

Open a shell and move to the lab folder. There create the virtual environment:

```
(skiplist-hnsw)$ python -m venv .venv
```

Activate the environment:

```
(skiplist-hnsw)$ source .venv/bin/activate
```

Install the required packages:

```
(skiplist-hnsw)$ pip install -r requirements.txt
```

# Skip Lists

---

# SkipList: Definition

- Linked List + Binary-search tree/sorted arrays
- Each node has 1+ layers, and each layer might be connected to different nodes
- Layers are assigned following a geometric distribution
- We'll use a simple implementation of Skip Lists, but you can extend it to include key-value pairs inside nodes

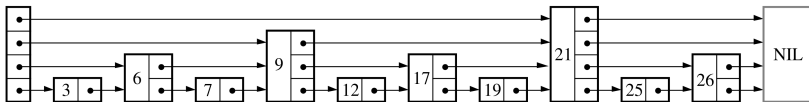


Figure 1: Example of Skip List (from [2])

---

**Algorithm 1** Search

---

**Require:** Skip list  $S$  with header  $head$ , value to search  $q$

**Ensure:** Node with  $q$  if present,  $NIL$  otherwise

```
1:  $current \leftarrow head$ 
2:  $level \leftarrow \text{max level of the skip list}$ 
3: for  $i = level$  to 0 do
4:   while  $current.forward[i] \neq NIL$ 
     and  $current.forward[i].item < q$  do
5:      $current \leftarrow current.forward[i]$ 
6:   end while
7: end for
8:  $current \leftarrow current.forward[0]$ 
9: if  $current \neq NIL$  and  $current.item = q$  then
10:  return  $current$ 
11: else
12:  return  $NIL$ 
13: end if
```

---

# Random Level Generation

---

## Algorithm 2 Random Level Generation

---

**Require:** Probability  $p$ , max level  $maxLevel$

**Ensure:** Random level for the new node

- 1:  $level \leftarrow 0$
  - 2: **while**  $random() < p$  **and**  $level < maxLevel$  **do**
  - 3:      $level \leftarrow level + 1$
  - 4: **end while**
  - 5: **return**  $level$
-



---

## Algorithm 3 Insert

---

**Require:** Skip list  $S$  with header  $head$ , value  $q$

**Ensure:** Skip list updated with the new node

```
1:  $update \leftarrow$  array of pointers with size  $maxLevel + 1$ 
2:  $current \leftarrow head$ 
3:  $level \leftarrow$  current max level
4: for  $i = level$  to 0 do
5:   while  $current.forward[i] \neq NIL$  and  $current.forward[i].item < q$  do
6:      $current \leftarrow current.forward[i]$ 
7:   end while
8:    $update[i] \leftarrow current$ 
9: end for
10:  $current \leftarrow current.forward[0]$ 
11: if  $current \neq NIL$  and  $current.item = q$  then
12:   return
13: end if
14:  $newLevel \leftarrow RandomLevel()$ 
15:  $newNode \leftarrow createNode(q, newLevel)$ 
16: if  $newLevel > currentLevel$  then
17:   for  $i = currentLevel + 1$  to  $newLevel + 1$  do
18:      $update[i] \leftarrow head$ 
19:   end for
20:    $currentLevel \leftarrow newLevel$ 
21: end if
22: for  $i = 0$  to  $newLevel + 1$  do
23:    $newNode.forward[i] \leftarrow update[i].forward[i]$ 
24:    $update[i].forward[i] \leftarrow newNode$ 
25: end for
```

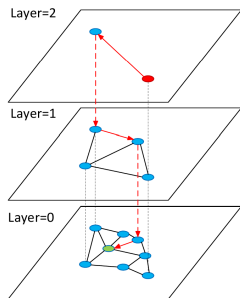
---

# HNSW

---

# From Skip Lists to HNSW

- HNSW[1] is a generalization of Skip List data structure
- Each node can have  $\geq 1$  neighbors at each level, not just one
- The neighbors are a proximity graph of the node



**Figure 2:** An HNSW graph

---

## Algorithm 4 Search

---

**Require:** multilayer graph  $hnws$ , query element  $q$ , number of nearest neighbors to return  $K$ , size of the dynamic candidate list  $efSearch$

**Ensure:**  $K$  nearest neighbors to  $q$

- 1:  $W \leftarrow \emptyset$
  - 2:  $ep \leftarrow$  enter point for  $hnsw$
  - 3:  $L \leftarrow$  level of  $ep$
  - 4: **for**  $l_c \leftarrow L \dots 0$  **do**
  - 5:      $W \leftarrow \text{SearchLayer}(q, ep, ef = 1, l_c)$
  - 6:      $ep \leftarrow$  get nearest element from  $W$  to  $q$
  - 7: **end for**
  - 8:  $W \leftarrow \text{SearchLayer}(q, ep, ef = efSearch, l_c = 0)$
  - 9: **return**  $\text{SelectNeighborsSimple}(q, W, K)$
-

# Search: how to search candidates on each layer?

---

## Algorithm 5 SearchLayer

---

**Require:** query element  $q$ , enter points  $ep$ , number of nearest neighbors to  $q$  to return  $ef$ , layer number  $l_c$

**Ensure:**  $ef$  nearest neighbors to  $q$

```
1:  $v \leftarrow ep$  // visited nodes. Note that the passed  $ep$  is a set of points, not a single node!
2:  $W \leftarrow ep$  // set of final results (updated during search)
3:  $C \leftarrow ep$  // set of candidate nodes for search steps
4: while  $|C| > 0$  do
5:    $c \leftarrow$  extract nearest element from  $C$  to  $q$ 
6:    $f \leftarrow$  get furthest element from  $W$  to  $q$ 
7:   if  $\text{distance}(q, c) > \text{distance}(q, f)$  then
8:     break
9:   end if
10:  for each  $e \in \text{neighbourhood}(c)$  at layer  $l_c$  do
11:    if  $e \notin v$  then
12:       $v \leftarrow v \cup e$ 
13:       $f \leftarrow$  get furthest element from  $W$  to  $q$ 
14:      if  $\text{distance}(q, e) < \text{distance}(q, f)$  or  $|W| < ef$  then
15:         $C \leftarrow C \cup e$ 
16:         $W \leftarrow W \cup e$ 
17:        if  $|W| > ef$  then
18:          remove furthest element from  $W$  to  $q$ 
19:        end if
20:      end if
21:    end if
22:  end for
23: end while
24: return  $W$ 
```

---

## Search: how to select neighbors (simple version)?

- This is the simplest version of a selection method for neighbors
- It may lead disconnected graphs, where few nodes are not properly linked to the others
- It's simpler and faster than other heuristics (perfect for us)

---

### Algorithm 6 Select Neighbors Simple

---

**Require:** base element  $q$ , candidate elements  $C$ , number of neighbors to return  $K$

**Ensure:**  $K$  nearest elements to  $q$

1: **return**  $K$  nearest elements from  $C$  to  $q$

---

---

## Algorithm 7 Insert



---

**Require:** multilayer graph  $hns$ , new element  $q$ , number of connections per node  $M_{max}$ , size of dynamic candidate list  $efConstruction$ , normalization factor  $m_L$

**Ensure:** update  $hns$  inserting new element  $q$

```
1:  $W \leftarrow \emptyset$  // list for the currently found nearest neighbors
2:  $ep \leftarrow$  get enter point for  $hns$ 
3:  $L \leftarrow$  level of  $ep$  // current top-layer for  $hns$ 
4:  $l \leftarrow \text{RandomLevel}(m_L)$ 
5: for  $l_c \leftarrow L \dots l + 1$  do
6:    $W \leftarrow \text{SearchLayer}(q, ep, ef = 1, l_c)$ 
7:    $ep \leftarrow$  get the nearest element from  $W$  to  $q$ 
8: end for
9: for  $l_c \leftarrow \min(L, l) \dots 0$  do
10:   $W \leftarrow \text{SearchLayer}(q, ep, efConstruction, l_c)$ 
11:   $neighbors \leftarrow \text{SelectNeighborsSimple}(q, W, M)$ 
12:  add bidirectional connections from  $neighbors$  to  $q$  at layer  $l_c$ 
13:  for each  $e \in neighbors$  do
14:     $eConn \leftarrow \text{neighbourhood}(e)$  at layer  $l_c$ 
15:    if  $|eConn| \geq M_{max}$  then
16:       $eNewConn \leftarrow \text{SelectNeighborsSimple}(e, eConn, M_{max})$ 
17:      set neighbourhood( $e$ ) at layer  $l_c$  to  $eNewConn$ 
18:    end if
19:  end for
20:   $ep \leftarrow W$ 
21: end for
22: if  $l > L$  then
23:   set enter point for  $hns$  to  $q$ 
24: end if
```

---

-  MALKOV, Y. A., AND YASHUNIN, D. A.  
Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs, 2018.
-  PUGH, W.  
Skip lists: a probabilistic alternative to balanced trees.  
*Commun. ACM* 33, 6 (June 1990), 668–676.