

VIETNAM NATIONAL UNIVERSITY, HCMC -  
UNIVERSITY OF SCIENCE

FACULTY OF INFORMATICS TECHNOLOGY  
ADVANCED PROGRAM IN COMPUTER SCIENCE



# DATA STRUCTURES - CS163

---

## Technical Report

### Task 2

---

*Students:*

Le Duc Tung Duong - 23125081

*Instructors:*

Thanh Ho Tuan, M.Sc.

Dung Nguyen Le Hoang, M.Sc.,

Tien Dinh Ba Ph.D.

2024

# Summary

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Improvements In Algorithm</b>	<b>3</b>
2.1	A* Algorithm . . . . .	3
2.1.1	Introduction . . . . .	3
2.1.2	Implementation . . . . .	3
2.1.3	Conclusion . . . . .	5
2.2	Contraction Hierarchies . . . . .	5
2.2.1	Introduction . . . . .	5
2.2.2	Implementation . . . . .	6
2.2.3	Preprocessing . . . . .	6
2.2.4	Query . . . . .	9
2.2.5	Conclusion . . . . .	11
<b>3</b>	<b>Further Improvement - Path Caching</b>	<b>12</b>
3.1	Introduction . . . . .	12
3.2	Implementation . . . . .	12
3.3	Conclusion . . . . .	14
<b>4</b>	<b>Bonus</b>	<b>14</b>
4.1	More comparison . . . . .	14
4.2	Generate .json file for stop information . . . . .	15
4.3	Query along with updating cache . . . . .	15
4.4	Output shortest path into .geojson file . . . . .	17
<b>5</b>	<b>Reference</b>	<b>19</b>

## §1 Introduction

In the previous course *CS162*, I built a graph based on a given set of data of a bus map and implemented *Dijkstra's algorithm* to find the shortest path between two nodes in a positive weighted directed graph. In this task of *CS163*, I will continue to use this graph to implement some improvements to find the shortest path more efficiently.

## §2 Improvements In Algorithm

In the report of task one, I showed that *Dijkstra's algorithm* is not efficient in finding the shortest path between two nodes in a graph with millions of nodes. Thus, in this task, I implement two algorithms which can run faster than *Dijkstra's: A\* algorithm* and *Contraction Hierarchies*, which are built based on *Dijkstra's*.

### §2.1 A\* Algorithm

#### §2.1.1 Introduction

$A^*$  is an algorithm developed in 1968 with the same purpose as *Dijkstra's*, pathfinding. However, it is a combination of a formal path-finding algorithm *Dijkstra's algorithm* and a heuristic approach *Greedy Best-First-Search*. Thus, it determines the step by favoring vertices close to both the source (*Dijkstra's algorithm*) and the target node (*Greedy Best-First-Search*). Although  $A^*$  is built on top of a heuristic approach, which usually does not guarantee the result, it can find the optimal shortest path.

#### §2.1.2 Implementation

In  $A^*$ , the only difference is  $f(n)$ , the cost of each vertex  $n$  and the sum of two values:

- $g(n)$ : the exact lowest cost from the source node to  $n$
- $h(n)$ : the estimated cost to the target node and also the heuristic function. Here, I use Euclidean distance for heuristic value.

Below is the **pseudocode** of  $A^*$  algorithm:

```

1: Initialize open_list as a priority queue
2: Initialize close_list as an empty set
3: function HEURISTIC(current_node, target_node)
4:   return Euclidean distance between current_node and target_node
5: end function
6: while open_list is not empty do
7:    $n \leftarrow$  node in open_list with the smallest  $f(n)$  value
8:   if  $n$  is the target node then
9:     return path to  $n$ 
10:  end if
11:  Add  $n$  to close_list
12:  for each neighbor  $m$  of  $n$  do
13:    if  $m$  is in close_list then
14:      continue
```

```

15:         end if
16:         tentative_g  $\leftarrow g(n) + \text{cost}(n, m)$ 
17:         if  $m$  is not in open_list or tentative_g  $< g(m)$  then
18:             Update  $g(m)$  to tentative_g
19:              $h(m) \leftarrow \text{heuristic}(m, \text{target\_node})$ 
20:              $f(m) \leftarrow g(m) + h(m)$ 
21:             if  $m$  is not in open_list then
22:                 Add  $m$  to open_list
23:             end if
24:         end if
25:     end for
26: end while
27: return failure if open_list is empty and target node is not found
    
```

Below the source code of  $A^*$  algorithm I implemented:

```

1  def a_star(self, source_node, target_node): # Load stop info before use
2      close_list = set()
3      open_list = []
4      g = {}
5      f = {}
6      h = {}
7      parents = {}
8
9      def heuristic(cur_id, target_id):
10         h[cur_id] = self.get_distance(self.stop_info[str(cur_id)][0],
11         ↪ self.stop_info[str(cur_id)][1],
12         ↪ self.stop_info[str(target_id)][0],
13         ↪ self.stop_info[str(target_id)][1])
14
15         g[source_node] = 0
16         heuristic(source_node, target_node)
17         f[source_node] = g[source_node] + h[source_node]
18         parents[source_node] = source_node
19         heapq.heappush(open_list, (f[source_node], source_node))
20     while open_list:
21         _, n = heapq.heappop(open_list)
22
23         if n in close_list:
24             continue
25
26         if n == target_node:
27             path = self.get_path(n, parents)
28             return f[target_node], path
29
30         close_list.add(n)
31
32         for m in self.G.neighbors(n):
33             g_new = g[n] + self.times_all[n][m]
    
```

```

32         if m not in h:
33             heuristic(m, target_node)
34         f_new = g_new + h[m]
35
36         if m not in f or f_new < f[m]:
37             g[m] = g_new
38             f[m] = f_new
39             parents[m] = n
40             heapq.heappush(open_list, (f_new, m))
41     return 0, []
    
```

### §2.1.3 Conclusion

**Space complexity:**  $O(V)$

**Time complexity:**  $O((E + V) \log V)$  in the worst case

Although the time complexity of  $A^*$  algorithm in the worst case is just the same as *Dijkstra's*,  $A^*$  in this case runs slower. This is because it has to consider two values,  $g(n)$  and  $h(x)$ , with every node. Thus, it takes longer.

The efficiency of  $A^*$  algorithm largely depends on the heuristic function. In a small road network, *Euclidean distance* is good enough to get the shortest path without checking all nodes, except for one case. When there are no shortest path between two nodes, it has to check all nodes, thus, even worse than *Dijkstra's*.

**In summary**, in road networks,  $A^*$  is more efficient than *Dijkstra's* in terms of time, especially in a large graph. With the given graph, there are only nearly 5000 nodes, so *Dijkstra's* and  $A^*$  do not show much difference. Sometimes, the former is even faster.

*e.g.* Given the problem of finding shortest path from stop 6860 to 1568. We compare  $A^*$  and *Dijkstra's*:

Algorithm	Time (s)	Nodes checked
$A^*$	0.0184	2,799
Dijkstra	0.0156	4,397

Although number of nodes  $A^*$  checks is fewer than *Dijkstra's*, it takes a little bit longer to do the query. As I explained, this graph is not really large and  $A^*$  has to check each node longer, thus executing time is slightly longer.

## §2.2 Contraction Hierarchies

### §2.2.1 Introduction

With a graph of millions of nodes, as road networks do, people try to upgrade *Dijkstra's algorithm* to quickly query the shortest path between two nodes. In a road network, these extension algorithms are called *route-planning algorithms*. Among these algorithms, *Contraction Hierarchies* offers a speedup in query with modest preprocessing time.

### §2.2.2 Implementation

Given a directed graph  $G = (V, E)$ , where each edge  $e$  in  $E$  has a weight representing the time it takes to travel on the road, we need to find the shortest path between source node and target node in  $V$ .

In a map, each node has its own importance. Nodes on highways, with many paths going through, usually have higher importance or higher chance to be passed in the shortest path while others, with fewer paths, have lower importance. This is the *hierarchical information* of a graph. We use this information to find the shortest path by getting to higher nodes from the source, then when getting close to the target node, we get to lower nodes until reaching the target.

Thus, *Contraction Hierarchies* has two processes:

- **Preprocess:** Categorize nodes and edges based on the importance.
- **Query:** Run a bidirectional search that only considers increasingly important edges.

### §2.2.3 Preprocessing

*Contraction Hierarchies* uses a process of node contraction based on orders of importance to create a hierarchy. After each contraction, some shortcut edges may be added to reduce the number of nodes checked during the query process. When contracting a node  $v$ , if the shortest path between two neighbors  $u$  and  $w$  of  $v$  includes  $v$ , then we add a shortcut between  $u$  and  $w$  with the weight being the cost of the shortest path.

Consider contracting node  $v$ , we have two lists:  $U$  containing  $v$ 's predecessors and  $W$  containing  $v$ 's neighbors (or children). With each node  $u$  in  $U$ , we calculate  $P(w) = \text{weight}(u, v) + \text{weight}(v, w)$  with each  $w$  in  $W$ .

Then, we run the local *Dijkstra* of  $u$  in a graph having no  $v$  until the cost reaches  $P_{\max}$ . If the path we found is longer than the one with  $v$ ,  $\text{dist}(w) > P(w)$ , then a shortcut edge is added between  $u$  and  $w$  with the weight  $P(w)$ .

Below is the process of adding shortcuts when contracting  $v$ :

```

1  for u in list(self.G.predecessors(v)):
2      if u in self.order_of:
3          continue
4      P = {}
5      for w in list(self.G.neighbors(v)):
6          if w in self.order_of:
7              continue
8          P[w] = self.times_all[u][v] + self.times_all[v][w]
9      if not P:
10         continue
11     P_max = max(P.values())
12
13     D = self.local_dijkstra_without_v(u, v, P_max)
14     for w in list(self.G.neighbors(v)):
15         if w in self.order_of:
```

```

16         continue
17
18     if D[w] > P[w]:
19         if self.G.has_edge(u, w):
20             self.G.get_edge_data(u, w)[0]['shortcut_node'] = v
21         else:
22             self.G.add_edge(u, w, shortcut_node=v)
23             self.times_all[u][w] = P[w]

```

Below is the local *Dijkstra's algorithm* used in above code with the source node  $u$  in a graph having no  $v$ :

```

1  def local_dijkstra_without_v(self, u, v, P_max):
2      vertices = list(self.G.nodes)
3      visited = set()
4      pq = [(0, u)]
5      D = {v: float('inf') for v in vertices}
6      visited.add(v)
7      D[u] = 0
8      while pq:
9          cost, n = heapq.heappop(pq)
10         if n in visited:
11             continue
12         if cost > P_max:
13             break
14         visited.add(n)
15
16         for neighbor in list(self.G.neighbors(n)):
17             if neighbor in self.order_of:
18                 continue
19             old_cost = D[neighbor]
20             new_cost = D[n] + self.times_all[n][neighbor]
21             if new_cost < old_cost:
22                 D[neighbor] = new_cost
23                 heapq.heappush(pq, (new_cost, neighbor))
24     return D

```

We do this with each node in order of importance. This order is just used to create shortcuts then we use these shortcuts to search. Thus, *the fewer shortcuts, the better the search is*. I implemented two ways of determining the order: *degree* and *edge difference*.

- **Degree:** With this approach, to get the degree of  $v$ , we just need to calculate the number of edges containing  $v$  by adding the number of predecessors with the number of neighbors. Then, sort all the nodes with ascending degree and get the order.

Here, I use *Min Heap* to store the degree and the node along with it. Then, I gradually take each element in the heap to store in two predefined lists, which store orders of all the nodes.

```

1     def get_node_order(self):
2         node_pq = []
3         for v in list(self.G.nodes()):
4             val = len(list(self.G.neighbors(v))) +
                    ↪ len(list(self.G.predecessors(v)))
5             heapq.heappush(node_pq, (val, v))
6         i = 1
7         while node_pq:
8             _, v = heapq.heappop(node_pq)
9             self.node_order[i] = v
10            self.order_of[v] = i
11            i += 1
    
```

- **Edge Difference:** The edge difference of node  $v$  is the difference between the number of shortcuts added when contracting only  $v$  in the graph and its degree (number of edges containing  $v$  in the beginning).

To get the number of shortcuts added when contracting  $v$ , we run the local *Dijkstra* in a graph with no  $v$ , as introduced above. Then, we sort all the nodes based on the edge difference and get the initial order.

```

1     def edge_difference(self, v):
2         dif = - len(list(self.G.neighbors(v))) -
                    ↪ len(list(self.G.predecessors(v)))
3         for u in list(self.G.predecessors(v)):
4             if u in self.order_of:
5                 continue
6             P = {}
7             for w in list(self.G.neighbors(v)):
8                 if w in self.order_of:
9                     continue
10            P[w] = self.times_all[u][v] + self.times_all[v][w]
11            if not P:
12                continue
13            P_max = max(P.values())
14
15            D = self.local_dijkstra_without_v(u, v, P_max)
16
17            for w in list(self.G.neighbors(v)):
18                if w in self.order_of:
19                    continue
20                if D[w] > P[w]:
21                    dif += 1
22            return dif
    
```

However, during the process of preprocessing, when a node is contracted, the edge



difference of others may be affected. Thus, we need to precompute the edge difference until we get the minimum edge difference of the current graph.

```

1  while node_pq:
2      # Calculate edge difference again to update the pq and get the next node
3      _, v = heapq.heappop(node_pq)
4
5      new_dif = self.edge_difference(v)
6      if node_pq and new_dif > node_pq[0][0]:
7          heapq.heappush(node_pq, (new_dif, v))
8          continue
9
10     # Continue with the process of adding shortcuts when contracting v

```

### §2.2.4 Query

There are many approaches for bidirectional search but here I implemented *bidirectional Dijkstra*.

- From the source node, I go forward to higher-ordered nodes of the current node  $v$ , which are also the neighbors (or children) of  $v$ .
- From the target node, I go backward to higher-ordered nodes of the current node  $v$  which are also the predecessors of  $v$ .

```

1  def bidirectional_dijkstra(self, source_node, target_node):
2      vertices = list(self.G.nodes())
3      visited_start = set()
4      visited_end = set()
5      parents1 = {}
6      parents2 = {}
7      dist1 = {v: float('inf') for v in vertices}
8      dist2 = {v: float('inf') for v in vertices}
9
10     parents1[source_node] = source_node
11     parents2[target_node] = target_node
12     dist1[source_node] = 0
13     dist2[target_node] = 0
14     pq_start = [(0, source_node)]
15     pq_end = [(0, target_node)]
16     while pq_start or pq_end:
17         if pq_start:
18             _, current_vertex = heapq.heappop(pq_start)
19             if current_vertex in visited_start:
20                 continue
21             visited_start.add(current_vertex)
22

```

```

23         for neighbor in self.G.neighbors(current_vertex):
24             if self.order_of[neighbor] <= self.order_of[current_vertex]:
25                 continue
26
27             new_cost = dist1[current_vertex] +
↪ self.times_all[current_vertex][neighbor]
28             if new_cost < dist1[neighbor]:
29                 parents1[neighbor] = current_vertex
30                 dist1[neighbor] = new_cost
31                 heapq.heappush(pq_start, (new_cost, neighbor))
32
33         if pq_end:
34             _, current_vertex = heapq.heappop(pq_end)
35             if current_vertex in visited_end:
36                 continue
37             visited_end.add(current_vertex)
38
39         for neighbor in self.G.predecessors(current_vertex):
40             if self.order_of[neighbor] <= self.order_of[current_vertex]:
41                 continue
42
43             new_cost = dist2[current_vertex] +
↪ self.times_all[neighbor][current_vertex]
44             if new_cost < dist2[neighbor]:
45                 parents2[neighbor] = current_vertex
46                 dist2[neighbor] = new_cost
47                 heapq.heappush(pq_end, (new_cost, neighbor))

```

After that, we store the all the nodes that exist in both searches in  $L$ . From this, we find the one that satisfy:

$$\text{dist}(s, t) = \min\{\text{dist}(s, v) + \text{dist}(v, t)\} \quad \text{over all } v \in L$$

```

1     L = [v for v in self.G.nodes if dist1[v] != float('inf') and dist2[v] !=
↪ float('inf')]
2     if not L:
3         return 0, []
4
5     shortest_time = math.inf
6     common_node = 0
7     for v in L:
8         if shortest_time > dist1[v] + dist2[v]:
9             shortest_time = dist1[v] + dist2[v]
10            common_node = v

```

After that, we trace back the path by combining it with a recursive function to get the shortcut node stored in an edge. Thus, we get a path from the source node to the highest-ordered node that can be reached and a path from this node to the target node.

```

1  def generate_shortcut(start_node, end_node):
2      shortcut_node = self.G.get_edge_data(start_node,
↪   end_node)[0]['shortcut_node']
3      if shortcut_node != 0:
4          return generate_shortcut(start_node, shortcut_node) + [shortcut_node] +
↪   generate_shortcut(shortcut_node, end_node)
5      else:
6          return []
7
8      shortest_path = []
9      path1 = []
10     cur_node = common_node
11
12     while parents1[cur_node] != cur_node:
13         tmp_node = parents1[cur_node]
14         path = []
15         if self.G.get_edge_data(tmp_node, cur_node)[0]['shortcut_node'] != 0:
16             path = generate_shortcut(tmp_node, cur_node)
17         path1 = path + path1
18         path1 = [tmp_node] + path1
19         cur_node = tmp_node
20
21     cur_node = common_node
22     path2 = []
23     while parents2[cur_node] != cur_node:
24         path2.append(cur_node)
25         tmp_node = parents2[cur_node]
26         path = []
27         if self.G.get_edge_data(cur_node, tmp_node)[0]['shortcut_node'] != 0:
28             path = generate_shortcut(cur_node, tmp_node)
29         path2 += path
30         cur_node = tmp_node
31     path2.append(cur_node)
32
33     shortest_path = path1 + path2

```

### §2.2.5 Conclusion

The complexity of *Contraction Hierarchies* is very complicated to precisely calculate. So, I will take example to show how efficient it is.

e.g. Continue with two nodes in 2.1.3., from 6860 to 1568 , I will do some comparison:

Algorithm	Preprocess Time (s)	Query Time (s)	Nodes checked
A*	0.0	0.0184	2,799
Dijkstra	0.0	0.0156	4,397
Contraction Hierarchies	17.4383	0.0032	226

By taking a modest process time, the number of nodes needed to be checked is significantly reduced to 226 nodes (nearly **20 times** smaller than *Dijkstra's* and over **14 times** smaller than  $A^*$ ). Thus, query time is also reduced to **3.2 milliseconds** (nearly 6 times faster than  $A^*$  and nearly 5 times with Dijkstra's).

## §3 Further Improvement - Path Caching

### §3.1 Introduction

Suppose we have district  $U$  and  $V$ . From a source node  $u$  in  $U$  to a target node  $v$  in  $V$ , there can be a fixed route which all shortest paths from  $u$  to  $v$  go through. Our target is to store this route and use it when doing query.

### §3.2 Implementation

By calculating all shortest paths from  $U$  to  $V$ , we can take the LCS (Longest Common String) of all of these paths. I implemented two approaches to find LCS:

- **Approach 1:** Manually check each node in  $path2$  with  $path1$

```

1     def update_fixed_point(self, path1, path2, time1, time2):
2         new_path = []
3         cnt = {v: 1 for v in path1}
4
5         for u in path2:
6             if cnt.get(u, 0) == 1:
7                 i = path1.index(u)
8                 j = path2.index(u)
9
10                while i < len(path1) and j < len(path2) and path1[i] ==
↪ path2[j]:
11                    new_path.append(path1[i])
12                    i += 1
13                    j += 1
14                break
15
16        if not new_path:
17            if time2 < time1:
18                return time2, path2
19            else:
20                return time1, path1
21
22        new_time = 0
23        for u, v in zip(new_path, new_path[1:]):
24            new_time += self.times_all[u][v]
25        return new_time, new_path

```

- **Approach 2:** *Counter()* function in *collections* library

```

1     def update_intersection(self, path1, path2, time1, time2):
2         result = collections.Counter(path1) & collections.Counter(path2)
3         new_path = list(result.elements())
4
5         new_time = 0
6         for u, v in zip(new_path, new_path[1:]):
7             new_time += self.times_all[u][v]
8
9         return new_time, new_path

```

Because we have to calculate shortest paths of all pairs of nodes, I use *Dijkstra's algorithm* to do this work. From a node  $u$ , it can visit all other nodes in modest time with  $O((E + V) \log V)$ . Besides, I store a boolean variable *have\_fixed\_point* to indicate whether there is a fixed route between two districts or not.

```

1     def generate_caching_json(self):
2         with open("input_data/stop_coordinate.json", "r", encoding='utf-8') as f:
3             stop_zone = json.load(f)
4         with open("input_data/path_caching_demo.json", "r", encoding='utf-8') as f:
5             fixed_point = json.load(f)
6
7         vertices = list(self.G.nodes)
8         for source_node in vertices:
9             shortest_time, parents = self.dijkstra(source_node)
10            source_zone = stop_zone[str(source_node)][2]
11
12            for target_node in vertices:
13                path = self.get_path(target_node, parents)
14                target_zone = stop_zone[str(target_node)][2]
15
16                if fixed_point[source_zone][target_zone]['have_fixed_point'] == False:
17                    continue
18                elif fixed_point[source_zone][target_zone]['have_fixed_point'] == True:
19                    fixed_point[source_zone][target_zone]['time'],
20                    ↪ fixed_point[source_zone][target_zone]['path'] =
21                    ↪ self.update_intersection(fixed_point[source_zone][target_zone]['path']
22                    ↪ , path, fixed_point[source_zone][target_zone]['time'],
23                    ↪ shortest_time[target_node])
24
25                if fixed_point[source_zone][target_zone]['time'] == 0:
26                    fixed_point[source_zone][target_zone]['have_fixed_point'] = False
27            else:
28                fixed_point[source_zone][target_zone]['have_fixed_point'] = False
29                for v in path:
30                    if stop_zone[str(v)][2] != source_zone and stop_zone[str(v)][2]
31                    ↪ != target_zone:

```

```

27         fixed_point[source_zone][target_zone]['have_fixed_point'] =
        ↪ True
28         fixed_point[source_zone][target_zone]['time'] =
        ↪ shortest_time[target_node]
29         fixed_point[source_zone][target_zone]['path'] = path
30         break
31
32     with open("input_data/path_caching_demo.json", "w", encoding='utf-8') as f:
33         json_object = json.dumps(fixed_point, indent=4, ensure_ascii=False)
34         f.write(json_object)

```

### §3.3 Conclusion

Although path caching is a feature to optimize searching process, it is not suitable with all other algorithms, especially it needs to find shortest paths of two segments.

*e.g.* Continue with example in 2.1.3. and 2.2.5., from 6860 to 1568, I add path caching to each algorithm to see their behaviour:

Algorithm	Preprocess Time (s)	Query Time (s)	Nodes checked
A*	0.0	0.0184	2,799
A* (Path Caching)	0.0	0.0073	2,706
Dijkstra	0.0	0.0156	4,397
Dijkstra (Path Caching)	0.0	0.0192	8,794
CH	17.4383	0.0032	226
CH (Path Caching)	17.0597	0.0054	331

\*CH: *Contraction Hierarchies*

With algorithms which need to check all other nodes like *Dijkstra* or all higher-ordered nodes like *Contraction Hierarchies*, we need to do this process twice for two segments. Thus, number of nodes checked is bigger (331 nodes > 226 nodes - *Contraction Hierarchies* or 8794 nodes > 4397 nodes - *Dijkstra*), resulting slower query time.

Meanwhile, A\* go straight to the node we need to find. Combined with path caching, it can reduce number of nodes and speed up query process **over 2 times**.

**In summary**, only algorithms going straight to the target node like A\* uses path caching efficiently.

## §4 Bonus

### §4.1 More comparison

To have an overview about each algorithm with and without caching, I query 15\*4397 or 65955 pairs and calculate executing time of each.

The result is the same as before. In conclusion, path caching is suitable with A\* than *Contraction Hierarchies*. However, the latter is faster in all cases because it has preprocessed.

Algorithm	Query Time (s)	Query Time (Path Caching) (s)
A*	0.0114	0.0088
Contraction Hierarchies	0.0017	0.0018

## §4.2 Generate .json file for stop information

To quickly get the coordinate and district (or zone) of a stop, I write a function to do the queries with all stops.

```

1  def get_stop_info(self):
2      stop_query = StopQuery()
3      stop_query.read_stop()
4      vertices = list(self.G.nodes())
5      coordinate = {}
6      cnt = 1
7      for u in vertices:
8          print(f".....Loading {cnt}/4397 stop coordinates.....")
9          cnt+=1
10         nodes = stop_query.searchByABC(StopId = u)
11         coordinate[u] = [nodes[0].get_lng(), nodes[0].get_lat(), nodes[0].get_zone()]
12     return coordinate

```

Then, I write information of all stops into a file.

```

1  def write_stop_info_json(self):
2      stop_info = self.read_stop_info()
3      with open("input_data/stop_coordinate.json", "w", encoding='utf-8') as f:
4          json_object = json.dumps(stop_info, indent=4, ensure_ascii=False)
5          f.write(json_object)
6      print("ALL STOP COORDINATES HAVE BEEN WRITTEN INTO FILE!")

```

Now, coordinate and district of all nodes are stored in the file *stop\_coordinate.json*. Each time we need stop information, we just need to call the function below to load all information into a dict.

```

1  def load_stop_info_json(self):
2      with open("input_data/stop_coordinate.json", "r", encoding='utf-8') as f:
3          self.stop_info = json.load(f)

```

## §4.3 Query along with updating cache

Suppose having a huge graph with hundreds of millions of nodes and millions of edges, the preprocessing process of finding fixed routes may take too many resources. In this case, I think we can do the query and update the cache at the same time.

Idea:

- **Step 1:** We need to find the shortest path from source node  $u$  to target node  $v$ . Suppose between 2 districts that 2 nodes belong to has no fixed route, we do apply

any search algorithm to find the path from  $u$  to  $v$ . Otherwise, there is a fixed route from  $a$  to  $b$ .

- **Step 2:** We find the shortest paths from  $u$  to  $a$  and from  $b$  to  $v$ . Adding three parts, together we have the path from  $u$  to  $v$ . However, this is not guaranteed the optimal solution, so we do a search from  $u$  to  $v$  alone. If it is shorter, we update the cache.

Here, the search I use is A\*, but it can be changed into any others.

```

1  def get_shortest_path_astar_pathcaching_updating(self, source_node, target_node):
2      self.load_stop_info_json()
3      with open("input_data/path_caching.json", "r", encoding='utf-8') as f:
4          fixed_point = json.load(f)
5          time = 0
6          path = []
7          source_zone = self.stop_info[str(source_node)][2]
8          target_zone = self.stop_info[str(target_node)][2]
9          if fixed_point[source_zone][target_zone]['have_fixed_point'] == False:
10             self.a_star(source_node, target_node)
11          elif fixed_point[source_zone][target_zone]['have_fixed_point'] == True:
12             new_time, new_path = self.a_star(source_node, target_node)
13             if new_time is not None:
14                 tmp_time = fixed_point[source_zone][target_zone]['time']
15                 tmp_path = fixed_point[source_zone][target_zone]['path']
16                 time1, path1 = self.a_star(source_node, tmp_path[0])
17                 time2, path2 = self.a_star(tmp_path[-1], target_node)
18                 cur_time = time1 + tmp_time + time2
19                 cur_path = path1 + tmp_path + path2
20
21             if new_time >= cur_time:
22                 time = cur_time
23                 path = cur_path
24             else:
25                 fixed_point[source_zone][target_zone]['time'],
26                 ↪ fixed_point[source_zone][target_zone]['path'] =
27                 ↪ self.update_intersection(tmp_path, new_path, tmp_time, new_time)
28                 if fixed_point[source_zone][target_zone]['time'] == 0:
29                     fixed_point[source_zone][target_zone]['have_fixed_point'] = False
30                 time = new_time
31                 path = new_path
32             else:
33                 new_time, new_path = self.a_star(source_node, target_node)
34                 fixed_point[source_zone][target_zone]['have_fixed_point'] = False
35                 if source_zone != "Ngoại thành" and target_zone != "Ngoại thành":
36                     for v in new_path:
37                         if self.stop_info[str(v)] != source_zone and self.stop_info[str(v)]
38                         ↪ != target_zone:
39                             fixed_point[source_zone][target_zone]['have_fixed_point'] = True
40                             fixed_point[source_zone][target_zone]['time'] = new_time

```



```

38         fixed_point[source_zone][target_zone]['path'] = new_path
39         break
40     time = new_time
41     path = new_path
42
43     print(time, path)
44     with open("input_data/path_caching.json", "w", encoding='utf-8') as f:
45         json_object = json.dumps(fixed_point, indent=4, ensure_ascii=False)
46         f.write(json_object)
47     if not path:
48         print(f'No p found from {source_node} to {target_node}!')
49         return 0, path
50     return time, path

```

Although it can be slower than A\* alone several milliseconds, it is so small that users may not care about. When the data is enough, we can do the query with path caching immediately as I introduced in section 3.

## §4.4 Output shortest path into .geojson file

To view the path more easily, I implement a function to view all stops and path on the map.

```

1  def output_shortest_path_json(self, paths):
2      # geojson file to view the path on the map
3      stop_query = StopQuery()
4      stop_query.read_stop()
5      geojson_data = {
6          "type": "FeatureCollection",
7          "features": [
8              {
9                  "type": "Feature",
10                 "geometry": {
11                     "type": "LineString",
12                     "coordinates": []
13                 },
14             }
15         ]
16     }
17
18     for (a, b) in zip(paths, paths[1:]):
19         coordinates = self.G.get_edge_data(a, b)[0]['coordinates']
20         for coordinate in coordinates:
21             geojson_data['features'][0]['geometry']['coordinates'].append(coordinate)
22
23     stop_objects = [stop_query.searchByABC(StopId=stopid) for stopid in paths]
24     for stop_coordinate, stop_id in zip(stop_objects, paths):
25         feature = {

```

```

26         "type": "Feature",
27         "geometry": {
28             "type": "Point",
29             "coordinates": [stop_coordinate[0].Lng, stop_coordinate[0].Lat]
30         },
31         "properties": {
32             "StopId": stop_id,
33         }
34     }
35     geojson_data['features'].append(feature)
36
37 json_object = json.dumps(geojson_data, indent = 4)
38 with open("shortest_path_coordinates.geojson", "w") as outfile:
39     outfile.write(json_object)
40
41     # json file for all the edges gone through and start, end stop of each
42     stop_query = StopQuery()
43     stop_query.read_stop()
44     edges_data = []
45     first_edge = False
46     for (a, b) in zip(paths, paths[1:]):
47         if first_edge is False:
48             edge = self.G.get_edge_data(a, b)[0]
49             first_edge = True
50             prev_route_id = edge['route_id']
51             prev_route_var_id = edge['route_var_id']
52         else:
53             edges = self.G.get_edge_data(a, b)
54             find = False
55             for key, item in edges.items():
56                 if item['route_id'] == prev_route_id and item['route_var_id'] ==
57                 ↪ prev_route_var_id:
58                     edge = item
59                     find = True
60                     break
61             if find is False:
62                 edge = edges[0]
63                 prev_route_id = edge['route_id']
64                 prev_route_var_id = edge['route_var_id']
65             start_stop = stop_query.searchByABC(StopId=a, RouteId=str(edge['route_id']),
66             ↪ RouteVarId=str(edge['route_var_id']))
67             end_stop = stop_query.searchByABC(StopId=b, RouteId=str(edge['route_id']),
68             ↪ RouteVarId=str(edge['route_var_id']))
69             edge_data = {
70                 'From': a,
71                 'To': b,
72                 'RouteId': edge['route_id'],
73                 'RouteVarId': edge['route_var_id'],

```

```
71         'Distance': edge['distance'],
72         'Time': edge['time'],
73         'StartStop': start_stop[0].get_stop(),
74         'EndStop': end_stop[0].get_stop()
75     }
76     edges_data.append(edge_data)
77     with open("shortest_path_stops.json", "w", newline='', encoding='utf-8') as f:
78         for edge in edges_data:
79             json_object = json.dumps(edge, ensure_ascii=False)
80             f.write(json_object + '\n')
```

## §5 Reference

[John Lazarsfeld - Tufts Comp 150 Alg Project - Contraction Hierarchies](#)

[Introduction to A\\* From Amit's Thoughts on Pathfinding](#)

[Geeks for Geeks - A\\* Search Algorithm](#)

This is the link to [my Github Repository](#).