

# Disjoint Set Union (DSU) – Hợp Tập Hợp Rời Rạc

Nguyễn Quân Bá Hồng\*

Ngày 24 tháng 8 năm 2025

## Tóm tắt nội dung

This text is a part of the series *Some Topics in Advanced STEM & Beyond*:

URL: [https://nqbh.github.io/advanced\\_STEM/](https://nqbh.github.io/advanced_STEM/).

Latest version:

- *Disjoint Set Union (DSU) – Hợp Tập Hợp Rời Rạc*.

PDF: URL: [https://github.com/NQBH/advanced\\_STEM\\_beyond/blob/main/OLP\\_ICPC/disjoint\\_set\\_union/NQBH\\_disjoint\\_set\\_union.pdf](https://github.com/NQBH/advanced_STEM_beyond/blob/main/OLP_ICPC/disjoint_set_union/NQBH_disjoint_set_union.pdf).

TeX: URL: [https://github.com/NQBH/advanced\\_STEM\\_beyond/blob/main/OLP\\_ICPC/disjoint\\_set\\_union/NQBH\\_disjoint\\_set\\_union.tex](https://github.com/NQBH/advanced_STEM_beyond/blob/main/OLP_ICPC/disjoint_set_union/NQBH_disjoint_set_union.tex).

- .

PDF: URL: [.pdf](#).

TeX: URL: [.tex](#).

## Mục lục

<b>1 Introduction to Disjoint Set Union – Nhập Môn Hợp Tập Hợp Rời Rạc</b>	<b>1</b>
1.1 Data Structure Disjoint Set Union – Cấu trúc dữ liệu Disjoint Set Union	2
1.1.1 Naive implementation of DSU – Cài đặt “ngây thơ” của DSU	3
1.1.2 1st Optimization: Merge according to size/height – Tối ưu 1: Gộp theo kích cỡ/độ cao	4
1.1.3 2nd Optimization: Path compression – Tối ưu 2: Nén đường đi	4
1.2 Time complexity & its proof – Độ phức tạp thời gian & chứng minh	5
1.3 An alternative implementation of DSU – 1 cách cài đặt khác	5
<b>2 Some Applications of DSU – Vài Ứng Dụng của DSU</b>	<b>5</b>
2.1 Save additional information for each set – Lưu thêm thông tin khác cho mỗi tập hợp	5
2.2 Optimize algorithms of finding minimum spanning trees (MST) in graphs – Tối ưu thuật toán tìm cây khung nhỏ nhất trong đồ thị	7
2.3 Reverse query – Đảo ngược truy vấn	7
2.4 Kiểm tra tính chất 2 phía của đồ thị online	10
<b>3 Some Techniques Using DSU’s Properties – Vài Kỹ Thuật Sử Dụng Tính Chất của DSU</b>	<b>10</b>
3.1 Small-to-large merging – Kỹ thuật gộp set	10
3.2 Kỹ thuật DSU trên cây (Sack)	10
<b>4 Miscellaneous</b>	<b>10</b>

## 1 Introduction to Disjoint Set Union – Nhập Môn Hợp Tập Hợp Rời Rạc

### Resources – Tài nguyên.

1. [Algorithms for Competitive Programming/disjoint set union](#).

2. BENJAMIN QI, ANDREW WANG, NATHAN GONG, MICHAEL CAO. [USACO Guide/Disjoint Set Union](#).

**Abstract.** The Disjoint Set Union (DSU) data structure, which allows you to add edges to a graph & test whether 2 vertices of a graph are connected.

– Cấu trúc dữ liệu Disjoint Set Union (DSU), cho phép bạn thêm các cạnh vào đồ thị & kiểm tra xem 2 đỉnh của đồ thị có được kết nối hay không.

---

\*A scientist- & creative artist wannabe, a mathematics & computer science lecturer of Department of Artificial Intelligence & Data Science (AIDS), School of Technology (SOT), UMT Trường Đại học Quản lý & Công nghệ TP.HCM, Hồ Chí Minh City, Việt Nam.  
E-mail: [nguyenquanbahong@gmail.com](mailto:nguyenquanbahong@gmail.com) & [hong.nguyenquanba@umt.edu.vn](mailto:hong.nguyenquanba@umt.edu.vn). Website: <https://nqbh.github.io/>. GitHub: <https://github.com/NQBH>.

Disjoint Set Union (abbr., DSU) là 1 cấu trúc dữ liệu hữu dụng, thường xuất hiện trong các kỳ thi Lập trình Thi Đấu, & có thể được dùng để quản lý 1 cách hiệu 1 tập hợp của các tập hợp.

**Bài toán 1.** Cho 1 đồ thị  $G = (V, E)$  có  $|V| = n \in \mathbb{N}^*$  đỉnh, ban đầu không có cạnh nào,  $E = \emptyset$ . Ta cần xử lý 2 loại truy vấn:

1. Thêm 1 cạnh giữa 2 đỉnh  $x, y \in V$  trong đồ thị, i.e.,  $E = E \cup \{(x, y)\}$  nếu  $G$  là đồ thị vô hướng &  $E = E \cup \{(x, y)\}$  nếu  $G$  là đồ thị có hướng; tuy nhiên DSU thường chỉ áp dụng cho đồ thị vô hướng nên ta chỉ xét trường hợp đồ thị vô hướng cho đơn giản.
2. In ra **yes** nếu như 2 đỉnh  $x, y$  nằm trong cùng 1 thành phần liên thông. In ra **no** nếu ngược lại.

## 1.1 Data Structure Disjoint Set Union – Cấu trúc dữ liệu Disjoint Set Union

Cho tiện, ta đánh số  $n$  đỉnh của đồ thị  $G$  bởi  $1, 2, \dots, n$  (trường hợp  $n$  đỉnh được dán nhãn bởi  $v_1, v_2, \dots, v_n$  hoàn toàn tương tự vì ta có thể làm việc trên chỉ số  $i$  của  $v_i$ ), khi đó  $V = [n]$ . Giả sử  $G$  có  $c := \text{num\_connected\_component} \in \mathbb{N}^*$  (số thành phần liên thông)  $C_1, C_2, \dots, C_c$  với  $\{C_i\}_{i=1}^c$  là 1 phân hoạch của  $V = [n]$ , i.e.:

$$\bigcup_{i=1}^c C_i = [n], \quad C_i \cap C_j = \emptyset, \quad \forall i, j \in [c], \quad i \neq j.$$

Nếu ta coi mỗi đỉnh của đồ thị  $G = (V, E)$  là 1 phần tử & mỗi thành phần liên thông (connected component) trong đồ thị là 1 tập hợp, truy vấn 1 sẽ trở thành gộp 2 tập hợp lần lượt chứa phần tử  $x, y$  thành 1 tập hợp mới & truy vấn 2 trở thành việc hỏi 2 phần tử  $x, y$  có nằm trong cùng 1 tập hợp không.

Để tiện tính toán & lý luận về mặt toán học cho riêng cấu trúc dữ liệu DSU, sau đây là 1 định nghĩa lai Toán–Tin mang tính cá nhân của tác giả [NQBH], hoàn toàn không chính thống trong Lý thuyết Đồ thị:

**Định nghĩa 1** (Chỉ số thành phần liên thông). Cho đồ thị vô hướng  $G = (V, E)$  với  $V = [n]$ . Gọi  $C(i) \subset [n]$  là thành phần liên thông của  $G = (V, E)$  chứa đỉnh  $i \in [n]$  & gọi chỉ số của thành phần liên thông chứa đỉnh  $i$  là  $\text{cid}(i)$ , i.e.,  $i \in C_{\text{cid}(i)} \equiv C(i)$ , với hàm  $\text{cid} : [n] \rightarrow [c]$  được gọi là hàm chỉ số liên thông.

Với định nghĩa 1, ta có ngay

$$\begin{cases} i \in C(i) = C_{\text{cid}(i)} \subset [n], \\ i, j \text{ are connected, } \forall j \in C(i). \end{cases} \quad \forall i \in [n].$$

Ở đây, ta coi mỗi đỉnh của đồ thị tự liên thông với chính nó theo nghĩa đỉnh đó đến được (reachability) chính đỉnh đó thông qua 1 đường đi có độ dài 0, được gọi là 1 *đường đi tầm thường*, chứ không phải theo nghĩa khuyên (loop).

**Lemma 1** (A characterization of connectedness – 1 đặc trưng hóa của tính liên thông). Cho đồ thị vô hướng  $G = (V, E)$ . (i) 2 đỉnh trên 1 đồ thị  $G$  không liên thông với nhau, i.e., không có đường đi nào trên  $G$  nối chúng khi & chỉ khi chúng thuộc 2 thành phần liên thông khác nhau, i.e.,

$$i, j \text{ are not connected} \Leftrightarrow C(i) \neq C(j) \Leftrightarrow C(i) \cap C(j) = \emptyset \Leftrightarrow \text{cid}(i) \neq \text{cid}(j), \quad \forall i, j \in [n].$$

(ii) 2 đỉnh trên đồ thị  $G$  liên thông với nhau, i.e., có đường đi trên  $G$  nối chúng khi & chỉ khi chúng cùng thuộc 1 thành phần liên thông, i.e.:

$$i, j \text{ are connected} \Leftrightarrow C(i) = C(j) \Leftrightarrow C(i) \cap C(j) \neq \emptyset \Leftrightarrow \text{cid}(i) = \text{cid}(j), \quad \forall i, j \in [n].$$

Với truy vấn 1, giả sử 2 đỉnh  $i, j \in [n]$  chưa có cạnh nối chúng trực tiếp, i.e.,  $\{i, j\} \notin E$ . Có 2 trường hợp xảy ra:

- Trường hợp 1: Giả sử  $i, j$  cùng thuộc 1 thành phần liên thông, theo Lem. 1, có  $C(i) = C(j)$ ,  $\text{cid}(i) = \text{cid}(j)$  nên ta chỉ cần thêm cạnh  $\{i, j\}$  vào tập cạnh  $E$ :  $E \leftarrow E \cup \{(i, j)\}$  hay `edge_list.append(i, j)` mà không cần cập nhật  $c$  tập liên thông  $\{C_i\}_{i=1}^c$  hay hàm chỉ số liên thông  $\text{cid}(\cdot)$ .
- Trường hợp 2: Giả sử  $i, j$  thuộc 2 thành phần liên thông khác nhau, theo 1, có  $C(i) \neq C(j)$ ,  $\text{cid}(i) \neq \text{cid}(j)$ . Khi đó, việc nối cạnh  $i, j$  với nhau, i.e., cập nhật tập cạnh bằng cách thêm cạnh  $\{i, j\}$  vào tập cạnh  $E$ :  $E \leftarrow E \cup \{(i, j)\}$  hay `edge_list.append(i, j)`, nhưng điều khác biệt ở đây nằm ở chỗ: ta cũng cần phải cập nhật tập các thành phần liên thông & hàm chỉ số liên thông như sau:

1. Cập nhật thành phần liên thông chứa cả  $i$  &  $j$  bằng cách lấy hợp của 2 tập hợp  $C(i), C(j)$ , i.e.,  $C_{\text{new}}(i) = C_{\text{new}}(j) := C(i) \cup C(j)$ .

2. Cập nhật hàm chỉ số liên thông:  $\text{cid}_{\text{new}}(i) = \text{cid}_{\text{new}}(j) = \min\{\text{cid}(i), \text{cid}(j)\}$  (cũng có thể lấy  $\max$  thay vì  $\min$  nhưng theo quy ước của DSU data structure, ta nên lấy  $\min$  để đảm bảo tính nhất quán<sup>1</sup>) & sau đó ta có thể đánh chỉ số các thành phần liên thông lại nếu cần vì sau khi nối  $i, j$  với nhau, số thành phần liên thông  $c$  đã giảm đi 1, i.e.,  $c \leftarrow c - 1$  hay  $-c$ .

Với truy vấn 2, với 2 đỉnh  $i, j \in [n]$ ,  $i \neq j$ , ta có thể kiểm tra 2 đỉnh này có cùng nằm trong 1 thành phần liên thông hay không bằng cách so sánh  $C(i)$  &  $C(j)$  hoặc  $\text{cid}(i)$  &  $\text{cid}(j)$ , nhờ Lem. 1. Việc thực hiện truy vấn này khá hiển nhiên do bản chất của cấu trúc dữ liệu DSU chính là mã sự liên thông thành các thành phần liên khác nhau để tiện quản lý.

Về mặt cài đặt, để giải Bài toán 1, ta sẽ xây dựng 1 cấu trúc dữ liệu có 3 thao tác sau:

1. **make\_set(v)** tạo ra 1 tập hợp mới chỉ chứa phần tử  $v$ :  $\text{output of make\_set}(v) = \{v\}$ .
2. **union\_sets(a, b)** gộp tập hợp chứa phần tử  $a$  & tập hợp chứa phần tử  $b$  thành 1, i.e., bước cập nhật thành phần liên thông chung bằng cách lấy hợp của 2 thành phần liên thông tương ứng:  $C_{\text{new}}(i) = C_{\text{new}}(j) := C(i) \cup C(j)$  đã để cập.
3. **find\_set(v)** cho biết *đại diện* (representative) của tập hợp có chứa phần tử  $v$  (phần tử đại diện cho  $v$  không nhất thiết phải là  $v$ , e.g., người đỡ đầu cho 1 tập hợp các đứa trẻ trong 1 trại trẻ mồ côi). Đại diện này là sẽ 1 phần tử của tập hợp đó & có thể thay đổi sau mỗi lần gọi thao tác **union\_sets** (e.g., nếu chọn phần tử đại diện là đỉnh có giá trị nhỏ nhất thì nếu tìm được đỉnh mới nhỏ hơn phần tử đại diện hiện tại, thì phải cập nhật phần tử đại diện mới này). Ta có thể sử dụng đại diện đó để kiểm tra 2 phần tử có nằm trong cùng 1 tập hợp hay không.  $a, b$  nằm trong cùng 1 tập hợp nếu như đại diện của 2 tập chứa chúng là giống nhau & không nằm trong cùng 1 tập hợp nếu ngược lại.

Ta có thể xử lý 3 thao tác này 1 cách hiệu quả với các tập hợp được biểu diễn dưới dạng các cây (tree-based representation), mỗi phần tử là 1 đỉnh & mỗi cây tương ứng với 1 tập hợp. Gốc của mỗi cây sẽ là đại diện của tập hợp đó.

Ban đầu, mỗi phần tử thuộc 1 tập hợp riêng biệt:  $\{1\}, \{2\}, \dots, \{n\}$ , nên mỗi đỉnh là 1 cây riêng biệt, cũng là 1 thành phần liên thông riêng biệt.  $|E|$  bước tiếp theo, ở bước thứ  $i \in [|E|]$ , ta gộp 2 tập hợp chứa phần tử  $a_i, b_i \in [n]$ ,  $a_i \neq b_i$ . Sau  $|E|$  bước, ta được  $c$  thành phần liên thông  $\{C_i\}_{i=1}^c$ . Với cách cài đặt này, ta sẽ lưu 1 mảng **parent** với **parent[v]** là cha của phần tử  $v$ .

### 1.1.1 Naive implementation of DSU – Cài đặt “ngây thơ” của DSU

Nếu 1 cây được đánh số sao cho nhãn của node cha luôn nhỏ hơn nhãn của node con thì ta có thể dễ dàng định nghĩa khái niệm *gốc* của 1 cây như sau:

**Định nghĩa 2.** Cho 1 đồ thị vô hướng  $G = (V, E)$ ,  $V = [n]$ . Gốc của 1 cây chứa 1 đỉnh  $i \in [n]$  được định nghĩa bởi công thức

$$r(i) := \min\{j \in [n]; i, j \text{ are connected}\}.$$

Để tạo 1 tập hợp mới gồm phần tử  $v$  bởi **make\_set(v)**, ta chỉ cần tạo 1 cây có gốc là  $v$ , với **parent[v] = v** (giống như kiểu phim *Predestination* (2014) – tạm dịch: *Tiền Định/Định Mệnh*, mình là cha & là mẹ của chính mình, cũng là con của chính mình luôn). Để gộp 2 tập hợp lần lượt chứa 2 phần tử  $a, b \in [n]$  bởi **union\_sets(a, b)**, ta sẽ tìm gốc  $r(a) \in [n]$  của cây có chứa phần tử  $a$  & gốc  $r(b) \in [n]$  của cây có chứa phần tử  $b$ . Nếu 2 giá trị này giống nhau  $r(a) = r(b) = r \in [n]$ , thì  $r \in C(a) \cap C(b) \Rightarrow C(a) \cap C(b) \neq \emptyset$  nên theo Lem. 1(ii),  $C(a) = C(b)$  hay  $a, b$  đã liên thông sẵn, nên ta sẽ không làm gì do 2 phần tử  $a, b$  đã nằm trong cùng 1 tập hợp chứa gốc chung  $r$ . Còn nếu không, i.e.,  $r(a) \neq r(b)$ , ta sẽ đặt gốc cây này là cha của gốc cây còn lại. Cụ thể hơn, nếu muốn khớp với Định nghĩa 2, ta áp dụng truy vấn 1 để thêm 1 cạnh giữa 2 đỉnh  $r(a), r(b) \in [n]$ :  $E \leftarrow E \cup \{(r(a), r(b))\}$ . Để thấy điều này sẽ gộp 2 cây lại thành 1. Hơn nữa, gốc chung mới vừa được cập nhật của cây chứa cả 2 đỉnh  $a, b$  chính là  $\min\{r(a), r(b)\}$ .

Để tìm ký hiệu của 1 tập hợp có chứa phần tử  $v$  bởi **find\_set(v)**, ta đơn giản nhảy lên các tổ tiên của đỉnh  $v$  cho đến khi ta đến gốc của cây. Thao tác này có thể dễ dàng được cài đặt bằng đệ quy.

```

1 void make_set(int v) {
2     parent[v] = v; // tạo ra cây mới có gốc là đỉnh v
3 }
4
5 int find_set(int v) {
6     if (v == parent[v]) return v; // trả về đỉnh v nếu như đỉnh v là gốc của cây
7     return find_set(parent[v]); // đệ quy lên cha của đỉnh v
8 }
9
10 void union_set(int a, int b) {
11     a = find_set(a); // tìm gốc của cây có chứa đỉnh a
12     b = find_set(b); // tìm gốc của cây có chứa đỉnh b
13     if (a != b) parent[b] = a; // gộp 2 cây nếu như 2 phần tử ở 2 cây khác nhau
14 }
```

Vì đây là cách cài đặt ngây thơ, ta có thể dễ dàng tạo ra 1 ví dụ sao cho khi sử dụng cách cài đặt này, cây sẽ trở thành 1 đoạn thẳng gồm  $n$  phần tử. Khi đó, độ phức tạp của thao tác **find\_set** sẽ là  $O(n)$  – hiển nhiên không thể chấp nhận được, nên ta tìm hiểu 2 phương pháp tối ưu thuật toán sau.

<sup>1</sup>Consistency is 1 of the kings in natural sciences.

### 1.1.2 1st Optimization: Merge according to size/height – Tối ưu 1: Gộp theo kích cỡ/độ cao

Phương pháp tối ưu này sẽ thay đổi thao tác `union_sets`, i.e., thay đổi cách xét trong 2 cây đang gộp, gốc của cây nào sẽ là cha của gốc của cây còn lại. Có khá nhiều cách để xét điều này, nhưng 2 cách được sử dụng nhiều nhất chính là gộp theo kích cỡ & gộp theo độ cao của cây. Giả sử mỗi cây có 1 giá trị, lần lượt là kích cỡ & độ cao của cây theo 2 cách gộp. Ở cả 2 cách gộp, ta sẽ luôn đặt gốc của cây có giá trị lớn hơn là cha của gốc của cây có giá trị nhỏ hơn.

Thao tác `union_sets` được tối ưu gộp theo kích cỡ:

```
1 // merge according to size
2 void make_set_size(int v) {
3     parent[v] = v;
4     size[v] = 1; // ban đầu tập hợp chứa v có kích cỡ là 1
5 }
6
7 void union_set_size(int a, int b) {
8     a = find_set(a);
9     b = find_set(b);
10    if (a != b) {
11        if (size[a] < size[b]) swap(a, b); // đặt biến a là gốc của cây có kích cỡ lớn hơn
12        parent[b] = a;
13        size[a] += size[b]; // cập nhật kích cỡ của cây mới gộp lại
14    }
15 }
```

Thao tác `union_sets` được tối ưu gộp theo độ cao:

```
1 // merge according to size
2 void make_set_height(int v) {
3     parent[v] = v;
4     rank[v] = 0; // gốc của cây có độ cao là 0
5 }
6
7 void union_set_height(int a, int b) {
8     a = find_set(a);
9     b = find_set(b);
10    if (a != b) {
11        if (rank[a] < rank[b]) swap(a, b); // đặt biến a là gốc của cây có độ cao lớn hơn
12        parent[b] = a;
13        if (rank[a] == rank[b]) ++rank[a]; // nếu 2 cây có cùng 1 độ cao, độ cao của cây mới
14        // sau khi gộp sẽ tăng 1
15    }
16 }
```

Chỉ cần sử dụng phương pháp tối ưu này, với 1 trong 2 cách cài đặt, độ phức tạp của thao tác `find_set` sẽ trở thành  $O(\log n)$ . Tuy nhiên, ta có thể tối ưu hóa hơn nữa khi kết hợp với phương pháp tối ưu thứ 2.

**Question 1.** Có cách nào hybrid để trộn giữa 2 cách hợp theo kích cỡ & độ cao không?

### 1.1.3 2nd Optimization: Path compression – Tối ưu 2: Nén đường đi

Phương pháp tối ưu nén đường đi nhằm tăng tốc thao tác `find_set`. Giả sử ta gọi hàm `find_set(v)` với 1 đỉnh  $v \in [n]$  bất kỳ, ta tìm được  $p$  là gốc của cây, đồng thời cũng là giá trị của mọi hàm `find_set(u)` với  $u$  là 1 đỉnh nằm trên đường đi từ  $v$  đến  $p$ .<sup>2</sup> Cách tối ưu ở đây chính là làm cho đường đi đến gốc của các đỉnh  $u$  ngắn đi bằng cách gán trực tiếp cha của các đỉnh  $u$  này thành  $p$ .

```
1 // path compression
2 int find_set_path_compression(int v) {
3     if (v == parent[v]) return v; // trả về đỉnh v nếu như đỉnh v là gốc của cây
4     int p = find_set(parent[v]); // đệ quy lên cha của đỉnh v
5     parent[v] = p; // nén đoạn từ v lên gốc của cây
6     return p;
7 }
```

1 cách cài đặt khác của thao tác `find_set` thường được sử dụng nhiều trong Competitive Programming do tính ngắn gọn của nó:

<sup>2</sup>NQBH: Trong bài viết gốc [VNOI Wiki/disjoint set union](#), tác giả viết nhầm thành “với  $u$  là 1 đỉnh nằm trên đường đi từ  $u$  đến  $p$ ” sửa lại thành “với  $u$  là 1 đỉnh nằm trên đường đi từ  $v$  đến  $p$ ” & vài chỗ khác nhầm ký hiệu đỉnh  $v$  thành  $u$ .

```

1 // brief find set
2 int find_set_brief(int v) {
3     return v == parent[v] ? v : parent[v] = find_set_brief(parent[v]);
4 }

```

## 1.2 Time complexity & its proof – Độ phức tạp thời gian & chứng minh

## 1.3 An alternative implementation of DSU – 1 cách cài đặt khác

### Resources – Tài nguyên.

1. LÊ MINH HOÀNG. *Giải Thuật & Lập Trình*.
2. AtCoder Library: <https://github.com/atcoder/ac-library>.

Thay vì cài đặt cấu trúc dữ liệu DSU bằng 2 mảng `parent`, `size`, chỉ 1 mảng `lab` được sử dụng. Nếu `lab[v] < 0` thì  $v$  là gốc của 1 cây `-lab[v]` là số lượng đỉnh của cây đó. Còn nếu `lab[v] > 0` thì `lab[v]` là cha của đỉnh  $v$ .

```

1 // alternative implementation of DSU
2 void make_set_lab(int v) {
3     lab[v] = -1;
4 }
5
6 int find_set_lab(int v) {
7     return lab[v] < 0 ? v : lab[v] = find_set_lab(lab[v]);
8 }
9
10 void union_set_lab(int a, int b) {
11     a = find_set_lab(a);
12     b = find_set_lab(b);
13     if (a != b) {
14         if (lab[a] > lab[b]) swap(a, b);
15         lab[a] += lab[b];
16         lab[b] = a;
17     }
18 }

```

## 2 Some Applications of DSU – Vài Ứng Dụng của DSU

### 2.1 Save additional information for each set – Lưu thêm thông tin khác cho mỗi tập hợp

Ngoài việc lưu các thông tin về cấu trúc cây, ta có thể lưu các hàm có tính chất giao hoán & kết hợp của từng tập hợp. E.g., ta có thể lưu tổng các phần tử/giá trị phần tử nhỏ nhất hoặc lớn nhất, trung bình cộng, etc. của từng tập hợp. Khi đó, các thao tác của DSU sẽ được cài đặt như sau:

```

1 // save additional information for each set in DSU
2 void make_set_size_more_info(int v) {
3     parent[v] = v;
4     size[v] = 1;
5     Min[v] = value[v]; // value[v] là giá trị của phần tử thứ v
6     Max[v] = value[v];
7     sum[v] = value[v];
8     avg[v] = value[v];
9 }
10
11 int find_set_size_more_info(int v) {
12     return v == parent[v] ? v : parent[v] = find_set_size_more_info(parent[v]);
13 }
14
15 void union_set_size_more_info(int a, int b) {
16     a = find_set_size_more_info(a);
17     b = find_set_size_more_info(b);
18     if (a != b) {
19         if (size[a] < size[b]) swap(a, b);
20         parent[b] = a;
21         size[a] += size[b];

```

```

22         sum[a] += sum[b];
23         Min[a] = min(Min[a], Min[b]);
24         Max[a] = max(Min[a], Min[b]);
25         avg[a] = static_cast<double>(sum[a]) / size(a);
26     }
27 }

```

Tương tự như thông tin về độ lớn `size` của cây hay độ cao `rank` của cây, ta sẽ lưu các hàm này tại gốc của từng cây.

```

1  int find_sum(int v) { // trả về tổng của các phần tử trong tập hợp chứa v
2      v = find_set_size_more_info(v);
3      return sum[v];
4  }
5
6  int find_min(int v) { // trả về giá trị nhỏ nhất của các phần tử trong tập hợp chứa v
7      v = find_set_size_more_info(v);
8      return Min[v];
9  }
10
11 int find_max(int v) { // trả về giá trị lớn nhất của các phần tử trong tập hợp chứa v
12     v = find_set_size_more_info(v);
13     return Max[v];
14 }
15
16 double find_avg(int v) { // trả về giá trị trung bình của các phần tử trong tập hợp chứa v
17     v = find_set_size_more_info(v);
18     return avg[v];
19 }

```

**Bài toán 2** (Xếp hàng). Cho  $n \in \mathbb{N}^*$  người đang xếp hàng ở các vị trí từ 1 đến  $n$ . Viết chương trình xử lý 2 truy vấn:

1. Người đứng ở vị trí thứ  $i$  rời khỏi hàng.
2. Tìm người gần nhất về bên phải vị trí  $p$  mà chưa rời khỏi hàng.

*1st solution: DSU.* Gọi  $Q_{\text{curr}} \subset [n]$  là hàng đợi hiện tại gồm những người còn đứng trong hàng. Với mỗi vị trí  $i \in [n]$ , ta sẽ có 1 con trỏ  $\text{ptr}(i)$ . Nếu người đứng ở vị trí  $i$  vẫn đang đứng trong hàng, con trỏ trỏ vào vị trí đó, i.e.,  $\text{ptr}(i) = i$ , nếu không thì con trỏ này sẽ trỏ vào con trỏ ở vị trí ngay bên phải. Ta có công thức đệ quy cho hàm con trỏ  $\text{ptr} : [n] \rightarrow [n]$  như sau:

$$\text{ptr}(i) = \begin{cases} i & \text{if } i \in Q_{\text{curr}}, \\ \text{ptr}(i+1) & \text{if } i \notin Q_{\text{curr}}. \end{cases}$$

Dễ dàng dự đoán & chứng minh bằng quy nạp công thức toán học cho hàm con trỏ  $\text{ptr} : [n] \rightarrow [n]$  như sau:

$$\text{ptr}(i) = \min\{j \in [n]; i \leq j, j \in Q_{\text{curr}}\} = \min\{j \in [n]; i \leq j, \text{ptr}(j) = j\},$$

có ý nghĩa là con trỏ ở vị trí  $i$  trỏ vào người gần nhất tính từ vị trí  $i$  về phía bên phải.

Giả sử ban đầu có  $n$  người trong hàng đợi,  $V := [n]$ ,  $Q_0 = V = [n]$ ,  $\text{ptr} = \text{id}_V$ , i.e.,  $\text{ptr}(i) = i$ ,  $\forall i \in [n]$ . Tại các bước tiếp theo, nếu người  $i \in [n]$  nào đó rời hàng,  $Q_{\text{curr}} := Q_{\text{lastest}} \setminus \{i\}$ , ta sẽ thực hiện truy vấn 1 của Bài toán 1 bằng cách nối 2 đỉnh  $i$  &  $i+1$  lại với nhau, i.e.,  $E \leftarrow E \cup \{\{i, i+1\}\}$ , ta sẽ gán  $\text{ptr}(i) = \text{ptr}(i+1)$ : dễ thấy để tìm người gần nhất bên phải mà chưa rời khỏi hàng, ta đi dần dần sang phải cho đến khi gặp 1 vị trí có con trỏ trỏ đến chính nó.

Ta có thể sử dụng cấu trúc dữ liệu DSU để lưu trữ các thông tin trên & sử dụng phương pháp tối ưu nén để đạt được độ phức tạp trung bình  $O(\log n)$  với mỗi truy vấn. Vị trí ta cần tìm chính là vị trí có thứ tự lớn nhất trong mỗi tập hợp (ứng với mỗi thành phần liên thông). Ta có thể lưu phần tử lớn nhất trong 1 tập hợp như đã cài đặt, qua đó đạt được độ phức tạp trung bình  $O(\alpha(n))$  với mỗi truy vấn.

```

1  #include <iostream>
2  #include <vector>
3  using namespace std;
4
5  vector<int> parent, Size, Max;
6
7  void make_set_size_max(int v) {
8      parent[v] = v;
9      Size[v] = 1;
10     Max[v] = v;
11 }

```

```

12
13 int find_set(int v) {
14     return v == parent[v] ? v : parent[v] = find_set(parent[v]);
15 }
16
17 void union_set_size_max(int a, int b) {
18     a = find_set(a);
19     b = find_set(b);
20     if (a != b) {
21         if (Size[a] < Size[b]) swap(a, b);
22         parent[b] = a;
23         Size[a] += Size[b];
24         Max[a] = max(Max[a], Max[b]);
25     }
26 }
27
28 void leave(int v) { // người thứ v rời khỏi hàng
29     union_set_size_max(v, v + 1);
30 }
31
32 int find_next(int p) { // trả về thứ tự của người gần nhất bên phải vị trí p mà chưa rời khỏi hàng
33     p = find_set(p);
34     return Max[p];
35 }
36
37 int main() {
38     int n, num_query, type_query, node; // query has 2 types:
39     // type 1: 1 i: remove person i from current queue &
40     // type 2: 2 p: find closest person from p to the right that is still in the queue
41     cin >> n >> num_query;
42     parent.resize(n);
43     Size.resize(n);
44     Max.resize(n);
45     for (int i = 1; i <= n; ++i) make_set_size_max(i);
46     for (int i = 0; i < num_query; ++i) {
47         cin >> type_query >> node;
48         if (type_query == 1) leave(node); // remove person i from current queue
49         else cout << find_next(node) << '\n';
50     }
51 }

```

□

## 2.2 Optimize algorithms of finding minimum spanning trees (MST) in graphs – Tối ưu thuật toán tìm cây khung nhỏ nhất trong đồ thị

### Resources – Tài nguyên.

1. [VNOI Wiki/minimum spanning tree – bài toán tìm cây khung nhỏ nhất trong đồ thị.](#)

Sử dụng DSU, ta có thể tối ưu độ phức tạp của thuật toán tìm cây khung nhỏ nhất của đồ thị từ  $O(m \log n + n^2)$  xuống  $O(m \log n)$ .

## 2.3 Reverse query – Đảo ngược truy vấn

Do tính chất 1 chiều của cấu trúc dữ liệu DSU: chỉ thêm đỉnh hoặc cạnh của đồ thị chứ không xóa được, ở 1 số bài ta phải đảo ngược thứ tự của các truy vấn trong bài để giải.

**Problem 1** (CodeForces/Intel Code Challenge Elimination Round (Div. 1 + Div. 2, combined)/C. destroying array). *You are given an array consisting of  $n$  nonnegative integers  $a_1, a_2, \dots, a_n$ . You are going to destroy integers in the array 1 by 1. Thus, you are given the permutation of integers from 1 to  $n$  defining the order elements of the array are destroyed.*

*After each element is destroyed, you have to find out the segment of the array, s.t. it contains no destroyed elements & the sum of its elements is maximum possible. The sum of elements in the empty segment is considered to be 0.*

**Input.** *The 1st line of the input contains a single integer  $n \in [10^5]$  – the length of the array. The 2nd line contains  $n$  integers  $a_1, a_2, \dots, a_n$ ,  $a_i \in [0, 10^9]$ ,  $\forall i \in [n]$ . The 3rd line contains a permutation of integers from 1 to  $n$  – the order used to destroy elements.*



**Output.** Print  $n$  lines. The  $i$ th line should contain a single integer – the maximum possible sum of elements on the segment containing no destroyed elements, after 1st  $i$  operators are performed.

Sample.

destroy_array.inp	destroy_array.out
4 1 3 2 5 3 4 1 2	5 4 3 0
5 1 2 3 4 5 4 2 3 5 1	6 5 5 1 0
8 5 5 4 4 6 6 5 5 5 2 8 7 1 3 4 6	18 16 11 8 8 6 6 0

**Bài toán 3** (Destroying array). Cho mảng gồm  $n \in \mathbb{N}^*$  số tự nhiên  $a_1, a_2, \dots, a_n \in \mathbb{N}$  & 1 hoán vị  $\pi \in S_n$  của  $[n]$ .  $n$  phần tử sẽ lần lượt bị phá hủy theo thứ tự hoán vị trên. Sau mỗi lần 1 phần tử bị phá hủy, in ra dãy con liên tiếp có tổng lớn nhất mà không có phần tử nào đã bị phá hủy. Tổng của 1 đoạn con rỗng là 0.

Giới hạn.  $n \in [10^5], a_i \in [0, 10^9], \forall i \in [n]$ .

1st solution: DSU. Source: <https://wiki.vnoi.info/algo/data-structures/disjoint-set-union> Do các phần tử là các số tự nhiên, nếu nếu sau khi 1 số phần tử bị phá hủy, dãy bị chia thành  $k$  đoạn con liên tiếp thì đáp án sẽ là 1 trong  $k$  đoạn này. Đảo ngược thứ tự của các truy vấn, ta có thể thấy bài toán trở nên dễ dàng hơn rất nhiều: Hồi sinh 1 số bị phá hủy trở về ban đầu & in ra đoạn con có tổng lớn nhất. Khi đó ta dùng cấu trúc dữ liệu DSU để xử lý các đoạn con liên tiếp.

Khi 1 số được hồi sinh, ta sẽ kiểm tra vị trí ngay bên trái số đó (chứ không phải tất cả vị trí bên trái số đó), nếu có số nào đã được hồi sinh từ trước thì ta sẽ thêm cạnh giữa số đó & số ngay bên trái số đó. Tương tự với số ngay bên phải. Ở mỗi thời điểm, các thành phần liên thông trong DSU sẽ thể hiện cho 1 đoạn con liên tiếp. Việc lưu trữ tổng của 1 thành phần liên thông đã được cài đặt ở trên.

C++ implementation.

1. VNOI's C++: destroying array: <https://wiki.vnoi.info/algo/data-structures/disjoint-set-union>:

```

1  #include <iostream>
2  #include <vector>
3  using namespace std;
4  #define ll long long
5
6  const int N = 1e5 + 5;
7  int n, a[N], p[N];
8  ll ans, res[N];
9  bool flag[N];
10
11 struct DSU {
12     vector<ll> parent, sz, sum;
13     DSU(ll n) : parent(n), sz(n), sum(n) {};
14
15     void make_set(ll v) {
16         parent[v] = v;
17         sz[v] = 1;
18         sum[v] = a[v];
19     }
20
21     ll find_set(ll v) {
22         return v == parent[v] ? v : parent[v] = find_set(parent[v]);
23     }

```



```

24
25 void join_sets(ll a, ll b) {
26     a = find_set(a);
27     b = find_set(b);
28     if (a != b) {
29         if (sz[a] < sz[b]) swap(a, b);
30         parent[b] = a;
31         sz[a] += sz[b];
32         sum[a] += sum[b];
33     }
34 }
35 };
36
37 int main() {
38     ios_base::sync_with_stdio(false); cin.tie(NULL);
39     cin >> n;
40     for (int i = 1; i <= n; ++i) cin >> a[i];
41     for (int i = 1; i <= n; ++i) cin >> p[i];
42     DSU g(n + 5);
43     for (int i = 1; i <= n; i++) g.make_set(i);
44     for (int i = n; i >= 1; --i) {
45         flag[p[i]] = true;
46         if (p[i] > 1 && flag[p[i] - 1]) g.join_sets(p[i], p[i] - 1);
47         if (p[i] < n && flag[p[i] + 1]) g.join_sets(p[i], p[i] + 1);
48         ans = max(ans, g.sum[g.find_set(p[i])]);
49         res[i - 1] = ans;
50     }
51     for (int i = 1; i <= n; ++i) cout << res[i] << '\n';
52 }

```

## 2. NQBH's C++: destroying array:

```

1  #include <iostream>
2  #include <vector>
3  using namespace std;
4
5  vector<int> a, p, parent, sz;
6  vector<long long> res, sum;
7  vector<bool> activated_pos;
8
9  void make_set_size_sum(int i) {
10     parent[i] = i;
11     sz[i] = 1;
12     sum[i] = a[i];
13 }
14
15 int find_set_size_sum(int v) {
16     return v == parent[v] ? v : parent[v] = find_set_size_sum(parent[v]);
17 }
18
19 void union_set_size_sum(int a, int b) {
20     a = find_set_size_sum(a);
21     b = find_set_size_sum(b);
22     if (a != b) {
23         if (sz[a] < sz[b]) swap(a, b);
24         parent[b] = a;
25         sz[a] += sz[b];
26         sum[a] += sum[b];
27     }
28 }
29
30 int main() {
31     int n;
32     cin >> n;
33     vector<bool> activated_pos(n + 1, false);

```

```

34     a.resize(n + 1);
35     p.resize(n + 1);
36     parent.resize(n + 1);
37     sz.resize(n + 1);
38     res.resize(n + 1);
39     sum.resize(n + 1);
40     res[n] = 0; // last result is always 0
41     for (int i = 1; i <= n; ++i) cin >> a[i];
42     for (int i = 1; i <= n; ++i) cin >> p[i];
43     activated_pos[p[n]] = true;
44     long long ans = a[p[n]]; // initialize current max sum
45     res[n - 1] = ans; // last element being destroyed
46     make_set_size_sum(p[n]);
47     for (int i = n - 1; i >= 1; --i) {
48         activated_pos[p[i]] = true;
49         make_set_size_sum(p[i]);
50         if (p[i] > 1 && activated_pos[p[i] - 1]) union_set_size_sum(p[i], p[i] - 1);
51         if (p[i] < n && activated_pos[p[i] + 1]) union_set_size_sum(p[i], p[i] + 1);
52         ans = max(ans, sum[find_set_size_sum(p[i])]);
53         res[i - 1] = ans;
54     }
55     for (int i = 1; i <= n; ++i) cout << res[i] << '\n';
56 }

```

□

## 2.4 Kiểm tra tính chất 2 phía của đồ thị online

**Bài toán 4.** Cho 1 đồ thị có  $n \in \mathbb{N}^*$  đỉnh, ban đầu không có cạnh nào. Xử lý các truy vấn thêm cạnh vào đồ thị. Hỏi sau truy vấn nào thì đồ thị không còn là đồ thị 2 phía?

## 3 Some Techniques Using DSU's Properties – Vài Kỹ Thuật Sử Dụng Tính Chất của DSU

Ngoài ra tính chất của DSU còn được sử dụng trong một số kĩ thuật khá phổ biến.

### 3.1 Small-to-large merging – Kỹ thuật gộp set

Giả sử ta cần lưu trực tiếp các phần tử của 1 tập hợp bằng 1 cấu trúc dữ liệu như set/map, thì liệu có cách nào đủ hiệu quả để thực hiện thao tác `union_sets` không? Có 1 kỹ thuật gộp set (small-to-large merging) để thực hiện yêu cầu đó.

### 3.2 Kỹ thuật DSU trên cây (Sack)

Kỹ thuật DSU trên cây (Sack) là 1 thuật toán sử dụng ý tưởng của kỹ thuật gộp set để giải quyết 1 số bài toán truy vấn trên cây 1 cách hiệu quả.

**Bài toán 5.** Cho 1 cây có  $n \in \mathbb{N}^*$  đỉnh với gốc là đỉnh 1, đỉnh thứ  $i$  được tô màu  $c_i$ . Cho  $q \in \mathbb{N}^*$  truy vấn có dạng  $v \ c$ , với mỗi truy vấn in ra số lượng đỉnh có màu  $c$  trong cây con gốc  $v$ .

## 4 Miscellaneous

**Question 2** (DSU vs. GNNs, GCNs, GATs). Liệu DSU data structure có thể có lợi để cài đặt Graph Neural Networks (GNNs), Graph Convolutional Network (GCNs), Graph Attention Networks (GATs) không?