

4 - Empirical Study of Union-Find

Jakob Eberhardt

`jakob.eberhardt@estudiantat.upc.edu`

May 7, 2025

Contents

1	Introduction	3
2	Basic Implementation	3
3	QuickUnion without Compression & Problem Statement	4
4	Heuristics for Union	6
4.1	Union by Weight	6
4.2	Union by Rank	7
5	Compression Heuristics for Find	7
5.1	Full Compression	7
5.2	Path Splitting	8
5.3	Path Halving	8
6	Experimental Setup	9
7	Results	10
7.1	Normalised Total Path Length	11
7.2	Normalised Total Pointer Updates	12
7.3	Normalised Cost	13

Listings

1	Virtual Class	3
2	Policy-based Implementation	3
3	Data Initialization	4
4	Union and Find Function Setup	4
5	Find Function without Compression	4
6	QuickUnion Policy	5
7	Union by Weight Policy	6
8	Union by Rank Policy	7
9	Full Compression Policy	7
10	Path Splitting Policy	8
11	Path Halving Policy	8
12	Measuring Metrics	9
13	Index-to-Pair Mapping	9
14	PairPermutation Class	9
15	UnionFind.h	15
16	DisjointSet.h	16
17	UnionPolicies.h	16
18	FindPolicies.h	17
19	main.cpp	18
20	Makefile	21

List of Figures

1	Initial Configuration	5
2	Example Sets after a Union	5
3	Example Sets after as second Union	5
4	Example Sets after a Union by Weight	6
5	Example Sets after a second Union by Weight	6
6	Normalised cost as a function of the fraction of blocks remaining (B/n) for every combination of union strategy (colour) and path-compression variant (marker / line style).	10

7	Normalised Total Path Length without compression.	11
8	Normalised Total Path Length of the Full Compression strategy.	11
9	Normalised Total Path Length of the Path Halving compression strategy.	11
10	Normalised Total Path Length of the Path Splitting compression strategy.	12
11	Normalised Total Pointer Updates required with the Full compression strategy.	12
12	Normalised Total Pointer Updates required with the Path Halving compression strategy.	12
13	Normalised Total Pointer Updates required with the Path Splitting compression strategy.	13
14	Normalised Cost without compression.	13
15	Normalised Cost of the Full Compression strategy.	13
16	Normalised for Path Halving compression.	14
17	Normalised Cost of the Splitting compression strategy.	14

1 Introduction

In this assignment, we study different heuristic strategies for finding and merging elements in disjoint sets. To this end, we set up an experiment to benchmark the respective implementations on synthetic data. To achieve machine-independent results we measure the following metrics:

- Total Path Length (**TPL**): The distance to its representative of every element in the data structure summed up. This measures to which degree the representation of our sets degenerated to a list-like data structure.
- Total Pointer Updates (**TPU**): Measures the number of required pointer updates during a **Find** operation for a given element.
- Total Cost: A heuristic-specific linear combination of TPL and TPU.

In section 2, we introduce the project structure and the policy-independent implementation. Section 3 includes the straight-forward **QuickUnion (QU)** and **No Compression (NC)** implementation and the degeneration problems that go along with it. The heuristics sections 4 and 5 include the different policies that can be applied for coping with these problems. In sections 6 and 7, we discuss the experimental setup and the results of the practical experiment.

2 Basic Implementation

A **Union Find** data structure allows us to store disjoint sets in an efficient manner. In the following section, we introduce the project’s common setup for the experiment. The data structure must feature three operations, adding a new set, merging two sets, and finding the representative member of a set. Hence, we can define a virtual class as an interface for our actual implementation, as can be seen in listing 1. In addition, we add the respective functions to access static properties and run time metrics of our data structure.

```
1 class DisjointSet {
2 public:
3     virtual ~DisjointSet() = default;
4
5     virtual void makeSet(int x) = 0;
6     virtual int find(int x) = 0;
7     virtual void unionSets(int x, int y) = 0;
8     virtual int getParent(int x) const = 0;
9     virtual int countSets() const = 0;
10    virtual int depth(int x) const = 0;
11    virtual long pointerUpdatesDuringFind(int x) const = 0;
12 }
```

Listing 1: The virtual class in `DisjointSet.h` defines the basic interface for our implementation and benchmarking setup

This interface is implemented by the `UnionFind` class as can be seen in listing 2. It is parameterized with two templates, a Union policy (**U**) which determines how two roots are linked, and a Find policy (**F**) which defines how we walk and compress the path to a root. This code design will allow us to easily compile and run every combination of $U \times F$ later on while avoiding redundant code. The data fields of our structure comprise the three vectors. For **QU+NC**, the only relevant data is stored in `parent_` vector. It keeps the index of the parent element for every element. The `size_` and `rank_` vectors will be used for union heuristics.

```
1 template<class U, class F>
2 class UnionFind : public DisjointSet
3 ...
4 private:
5     std::vector<int> parent_, size_, rank_;
```

Listing 2: Policy-based Implementation of the interface in `UnionFind.h`

The data can be initialized using the constructor which will allocate three parallel arrays of size n . It initializes every vertex as an isolated root and sets their `size` and `rank` to one and zero respectively. The `makeSet` function is for testing and ensures proper initialization and resizing if necessary.

```

1 public:
2     explicit UnionFind(int n)
3         : parent_(n), size_(n,1), rank_(n,0)
4     { for (int i=0;i<n;++i) parent_[i]=i; }
5
6     void makeSet(int x) override {
7         if (x >= (int)parent_.size()) {
8             int old = parent_.size();
9             parent_.resize(x+1);
10            size_.resize(x+1,1);
11            rank_.resize(x+1,0);
12            for (int i=old;i<=x;++i) parent_[i]=i;
13        } else {
14            parent_[x]=x; size_[x]=1; rank_[x]=0;
15        }
16    }

```

Listing 3: Constructor and data initialization in `UnionFind.h`

In listing 4, we can see how the `UnionFind` class implements common parts of the `unionSets` and `find` functions before the respective policy-specific code is in-lined. To minimize complexity, we pass all the parameters, independent of the heuristic which is later used, e.g. the dummy variable `upd` which could be used to track pointer updates.

```

1     int find(int x) override {
2         long upd=0;
3         return F::find(x,parent_,upd);
4     }
5
6     void unionSets(int a,int b) override {
7         int r1=find(a), r2=find(b);
8         if (r1==r2) return;
9         U::unite(r1,r2,parent_,size_,rank_);
10    }

```

Listing 4: The policy-independent parts of the `unionSets` and `find` functions in `UnionFind.h`. The parameters are passed to a specific combination of heuristics, e.g. `QuickUnion` without compression.

3 QuickUnion without Compression & Problem Statement

If we do not employ any compression strategy while finding the set representative, the `find` function simply chases the pointer until we encounter the set representative. Hence, we never update and pointer and we only read every node once, `followMult` is therefore one.

```

1 struct NC {
2     static int find(int x, std::vector<int>& p, long& upd)
3     {
4         while (x != p[x]) x = p[x];
5         return x;
6     }
7     static long updatesGivenDepth(int depth) { return 0; }
8     static double followMult() { return 1.0; }
9 };

```

Listing 5: The `find` implementation without any additional compression considerations in `FindPolicies.h`

As mentioned, upon initialization, each element belongs to its own set, meaning its parent value points to its own index (`parent[i] == i`). Therefore, we always start with a data structure that looks like figure 1.

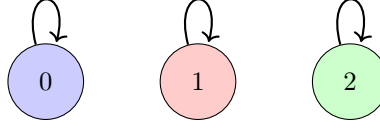


Figure 1: Initial configuration: three separate nodes (0-2). Each node points to itself (self-loop), indicating that each is initially its own set.

The policies for how two representatives are linked upon a union operation are defined in the `UnionPolicies.h` file. For many use cases, it is essential how roots are attached in order to keep the resulting tree balanced. For QuickUnion, however, we always assign the second operand as the parent of the first, as can be seen in listing 6.

```

1 struct QU {
2     static void unite(int r1, int r2,
3                       std::vector<int>& parent,
4                       std::vector<int>& size,
5                       std::vector<int>& rank)
6     { parent[r1] = r2; }
7 };

```

Listing 6: The QuickUnion policy defined in `UnionPolicies.h`

Following the example from figure 1, if we apply `union(0,1)` seen in figure 2 and `union(1,2)` in figure 3, we can see how a sequence of operations like this can cause an unbalance data structure or even a degeneration to a list.

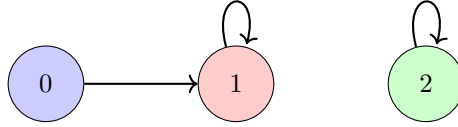


Figure 2: After `union(0,1)`: Element zero's parent becomes element one, forming a chain fragment ($0 \rightarrow 1$). Element one remains the representative of the set containing nodes zero and one, while element two is still a separate set.

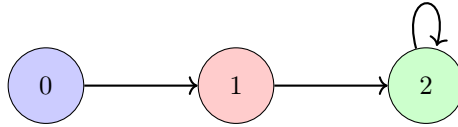


Figure 3: After `union(1,2)`: Element one's parent becomes element two, resulting in the complete chain: $0 \rightarrow 1 \rightarrow 2$. Element two is now the representative of the set.

In the following, we will discuss and compare heuristics for union and compression techniques during find operations to prevent this behavior.

4 Heuristics for Union

Every union operation will cause two find operations. The time required for each find operation is proportional to the amount of parent pointers we have to follow. Hence, we aim to keep the trees that represent our sets as flat as possible, independent of the input order. Additionally, we want to keep the amount of pointers that have to be updated to flatten a tree as low as possible.

4.1 Union by Weight

For the Union by weight heuristic, we use the `size` array to track how many elements are represented by the given root. We try to keep the trees shallow by putting the smaller tree under the larger one by comparing their size S_i .

```
1 struct UW {
2     static void unite(int r1,int r2,
3                       std::vector<int>& parent,
4                       std::vector<int>& size,
5                       std::vector<int>& rank)
6     {
7         if (size[r1] < size[r2]) std::swap(r1,r2);
8         parent[r2] = r1;
9         size[r1] += size[r2];
10    }
11};
```

Listing 7: The Union by weight policy defined in `UnionPolicies.h`

The height of the tree which remains a root may only grow by a maximum factor of two since

$$S' = S_1 + S_2 \leq S_1 + S_1 = 2 \cdot S_1$$

For example, for the sequence `union(0,1) → union(1,2)`, which previously resulted in a height of $n - 1$ depicted in figure 2 and 3, we now maintain a height of one, as can be seen in figure 4 and 5.

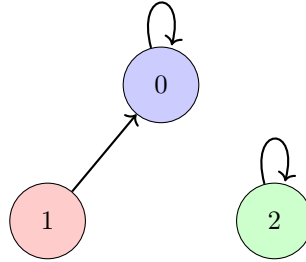


Figure 4: After `union(0,1)` with the union by weight policy. If the two subtrees weigh the same, the right operand becomes the root or respectively the representative of the set. Element 2 is still alone.

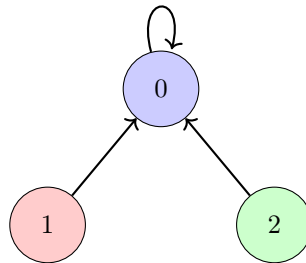


Figure 5: After `union(1,2)` with Union by weight: the smaller subtree (root 2) is attached to the larger one (root 0).

In the worst case, where we always have to merge equal-size trees like in this sequence: $\text{union}(0,1) \rightarrow \text{union}(2,3) \rightarrow \text{union}(4,5) \rightarrow \text{union}(6,7) \rightarrow \text{union}(0,2) \rightarrow \text{union}(4,6) \rightarrow \text{union}(0,4)$ we get a tree with the maximum height of $\Theta(\log n)$

4.2 Union by Rank

In Union by rank the rank of a root starts at zero and is monotone, meaning we never decrease it, e.g. during path compression. Each increment happens precisely when two equal-ranked trees merge, which at least doubles the size of the resulting tree, so after k increments the tree contains $\geq 2^k$ elements. This ties the rank to the maximum possible size of the set which is n . Hence for a set of size n we have $\text{rank} \leq \lfloor \log_2 n \rfloor$ and, because the actual height is never larger than the rank,

$$\text{height} \leq \text{rank} \leq \lfloor \log_2 n \rfloor$$

which gives us the same $\Theta(\log n)$ height guarantee as union by size while requiring less bookkeeping, because we just have to occasionally update a rank.

```

1 struct UR {
2     static void unite(int r1,int r2,
3                       std::vector<int>& parent,
4                       std::vector<int>& size,
5                       std::vector<int>& rank)
6     {
7         if (rank[r1] < rank[r2]) std::swap(r1,r2);
8         parent[r2] = r1;
9         if (rank[r1] == rank[r2]) ++rank[r1];
10    }
11 };

```

Listing 8: The Union by rank policy defined in `UnionPolicies.h`

5 Compression Heuristics for Find

We can employ different pointer compression heuristics during find operations to guarantee certain properties after a call. Depending on the use case or the targeted hardware, we can modify the forest more or less aggressively, as the following examples will showcase.

5.1 Full Compression

Full compression is the most aggressive heuristic. We follow the parent pointers to the root, then we make every node on that path point directly to the representative. Hence, the former tree structure is flattened to a star at the cost of a second traversal to update the pointers.

```

1 struct FC {
2     static int find(int x, std::vector<int>& p, long& upd)
3     {
4         if (p[x] != x) {
5             p[x] = find(p[x], p, upd);
6             ++upd;
7         }
8         return p[x];
9     }
10    static long updatesGivenDepth(int depth) { return depth; }
11    static double followMult() { return 2.0; }
12 };

```

Listing 9: The Full Compression policy defined in `FindPolicies.h`

5.2 Path Splitting

In path splitting, we split one long edge into two shorter edges. Namely, we skip over $\langle x, \text{parent}(x) \rangle$ and add $\langle x, \text{grandparent}(x) \rangle$ along with $\langle \text{parent}(x), \text{grandparent}(x) \rangle$ upon a single top-to-root pass. Hence, we get a height ≤ 2 on the visited path.

```
1 struct PS {
2     static int find(int x, std::vector<int>& p, long& upd)
3     {
4         while (p[x] != p[p[x]]) {
5             int parent = p[x];
6             p[x] = p[parent];
7             ++upd;
8             x = parent;
9         }
10        return p[x];
11    }
12    static long updatesGivenDepth(int depth) { return depth/2; }
13    static double followMult() { return 1.0; }
14};
```

Listing 10: The Path Splitting policy defined in FindPolicies.h

5.3 Path Halving

In path halving we shortcut every second edge encountered on the top-to-root scan by making each visited node point directly to its grandparent, e.g., $\langle x, \text{parent}(x) \rangle$ becomes $\langle x, \text{grandparent}(x) \rangle$. Because only alternate vertices are rewired, the depth of every node on the walked path is reduced from d to $\lceil d/2 \rceil$, *halving* the distance to the root. The whole path is processed in a single pass, so we obtain this depth reduction with exactly one traversal and about half as many pointer writes as full path compression.

```
1 struct PH {
2     static int find(int x, std::vector<int>& p, long& upd)
3     {
4         while (p[x] != p[p[x]]) {
5             p[x] = p[p[x]];
6             ++upd;
7             x = p[x];
8         }
9         return p[x];
10    }
11    static long updatesGivenDepth(int depth) { return (depth+1)/2; }
12    static double followMult() { return 1.0; }
13};
```

Listing 11: The Path Halving policy defined in FindPolicies.h

6 Experimental Setup

During experimentation, we use three instance sizes, small ($n = 1000$), medium ($n = 5000$) and large ($n = 10000$). The experimental setup is defined in the `main.cpp` file. The function `measureMetrics` scans all n vertices after each sampling point so that the recorded totals TPL and TPU reflect the entire forest at that moment for measurements at every Δ unions.

```
1 struct Metrics { long tpl = 0; long tpu = 0; };
2
3 Metrics measureMetrics(const DisjointSet& uf, int n)
4 {
5     Metrics m;
6     for (int v = 0; v < n; ++v) {
7         m.tpl += uf.depth(v);
8         m.tpu += uf.pointerUpdatesDuringFind(v);
9     }
10    return m;
11 }
```

Listing 12: The `measureMetrics` in `main.cpp`

The helper function `indexToPair` converts a rank $k \in 0, \dots, \binom{n}{2} - 1$ into the unique unordered pair (i, j) , turning the set of all pairs into a contiguous integer range that can be shuffled cheaply and without storing explicit pairs.

```
1 inline std::pair<int,int> indexToPair(long long k)
2 {
3     long long i = static_cast<long long>((1 + std::sqrt(1 + 8.0*k)) / 2);
4     while (i*(i-1)/2 > k) --i;
5     long long j = k - (i*(i-1))/2;
6     return {static_cast<int>(i), static_cast<int>(j)};
7 }
```

Listing 13: The `indexToPair` mapping function in `main.cpp`

The class `PairPermutation` builds a vector of those ranks, shuffles it with `std::shuffle` and a fixed seed to obtain a reproducible uniform random permutation, and then delivers one fresh rank per call to `next()`, guaranteeing that no pair is ever processed twice.

```
1 class PairPermutation {
2     public:
3         PairPermutation(long long m, unsigned seed)
4             : perm_(m), idx_(0)
5         {
6             for (long long i = 0; i < m; ++i) perm_[i] = i;
7             std::shuffle(perm_.begin(), perm_.end(), std::mt19937_64(seed));
8         }
9
10        bool next(long long &k)
11        {
12            if (idx_ >= perm_.size()) return false;
13            k = perm_[idx_++];
14            return true;
15        }
16        private:
17            std::vector<long long> perm_;
18            size_t idx_;
19 };
```

Listing 14: The `PairPermutation` Class in `main.cpp`

7 Results

From figure 6 we can see that Union-by-Rank and Union-by-Weight in combination with a compression heuristic yield the best results. QuickUnion without any compression performs the worst, as can be seen in figure 7. This showcases the importance of implementing a union heuristic beyond plain root attachment. Without any further compression, we can save orders of magnitude in terms of the expected total path length in the forest and hence a massive saving in cost, as can be seen in figure 14. For QuickUnion, we can also observe how increasing the instance size influences the shape of the TPL and Cost curves, as they grow faster as the instance size grows. Meanwhile, Union-by-weight and Union-by-rank maintain the same shape, indifferent of the input size. A surprising result of the experiment is the sharp drop in TPL and Cost when any kind of compression heuristic is employed, but most significantly when using Full Compression. As the amount of distinct blocks decreases, we likely traverse the same sets over and over again, which results in aggressive compression. The same behavior can be seen in the required pointer updates in figure 12, as many elements are already attached to the ultimate root, meaning they will never require an update again during the benchmark.

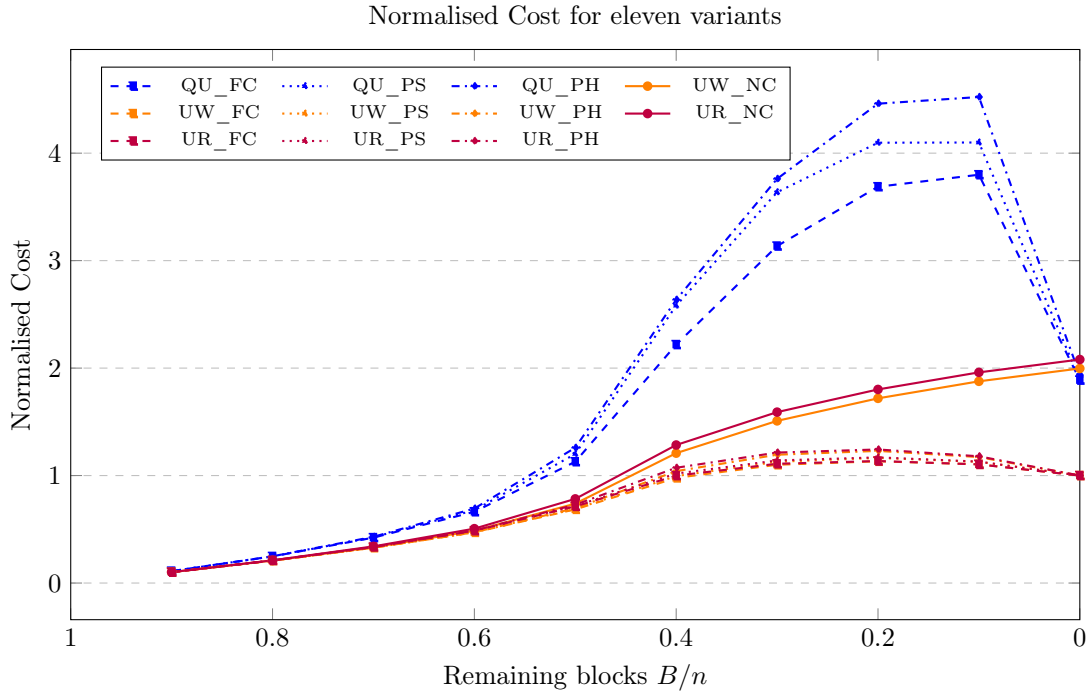


Figure 6: Normalised cost as a function of the fraction of blocks remaining (B/n) for every combination of union strategy (colour) and path-compression variant (marker / line style).

7.1 Normalised Total Path Length

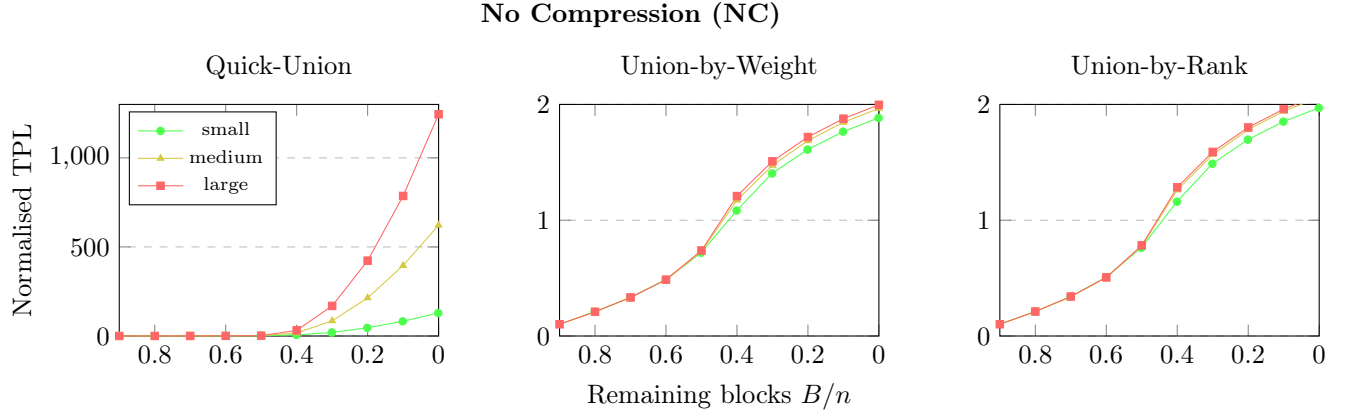


Figure 7: Normalised Total Path Length without compression.

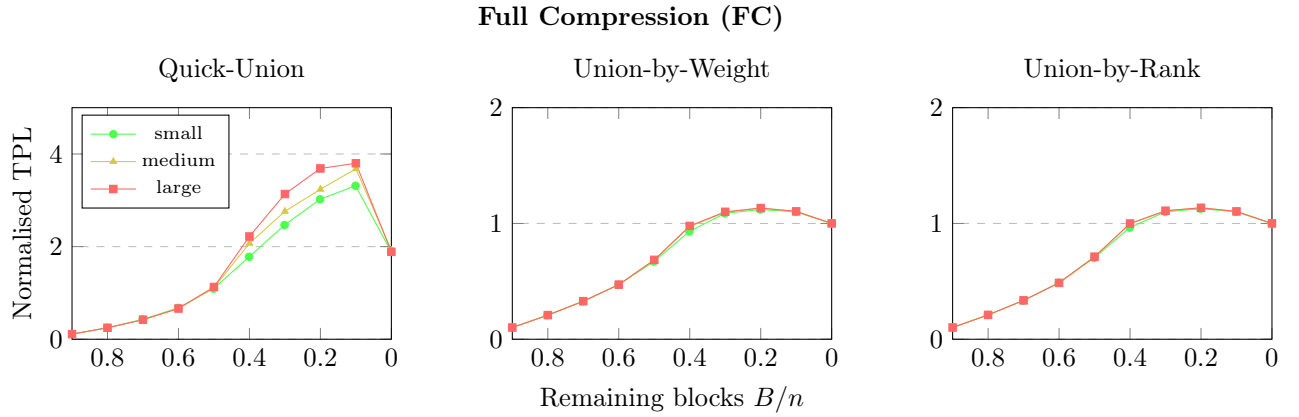


Figure 8: Normalised Total Path Length of the Full Compression strategy.

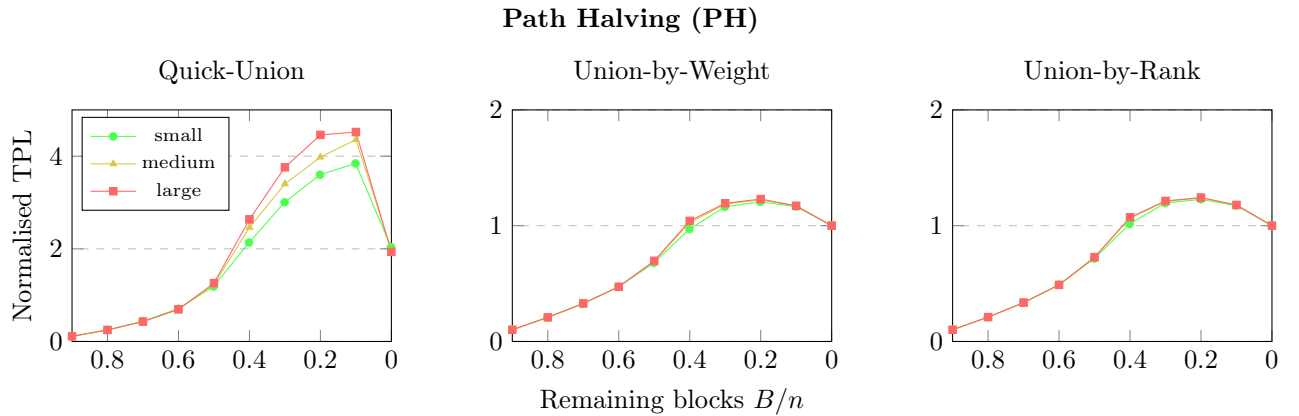


Figure 9: Normalised Total Path Length of the Path Halving compression strategy.

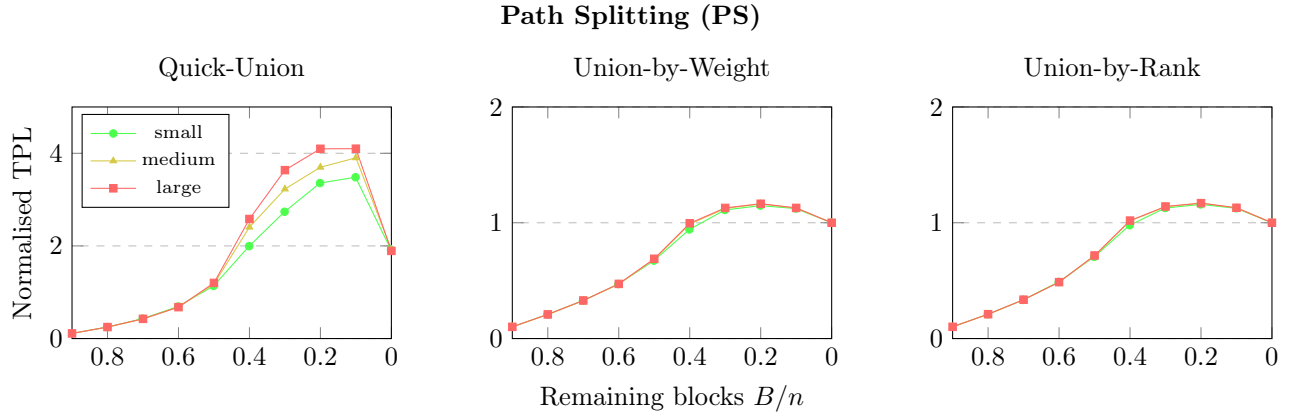


Figure 10: Normalised Total Path Length of the Path Splitting compression strategy.

7.2 Normalised Total Pointer Updates

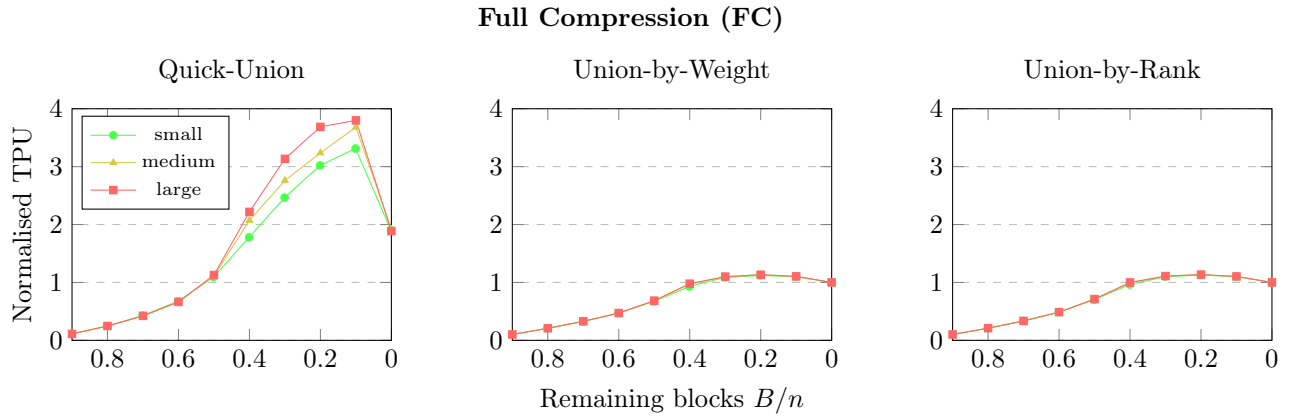


Figure 11: Normalised Total Pointer Updates required with the Full compression strategy.

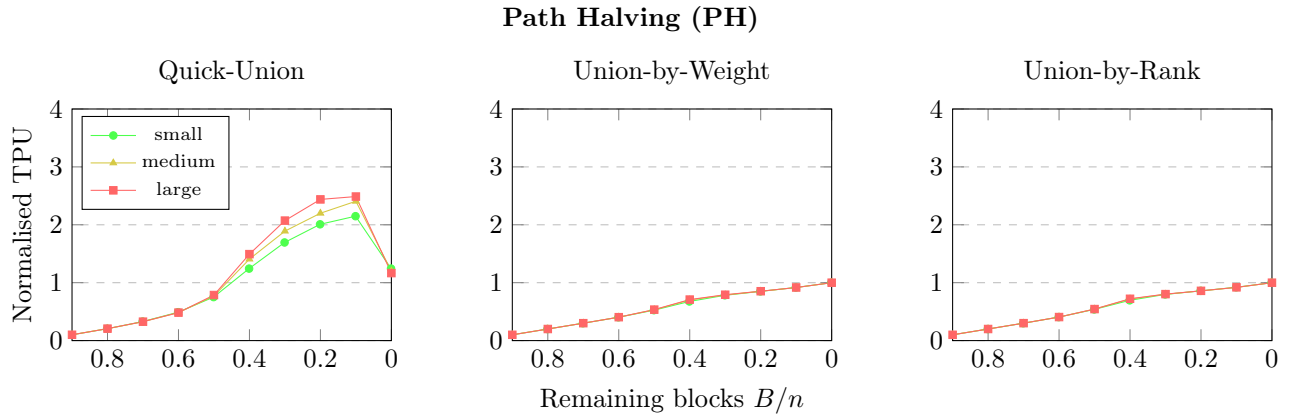


Figure 12: Normalised Total Pointer Updates required with the Path Halving compression strategy.

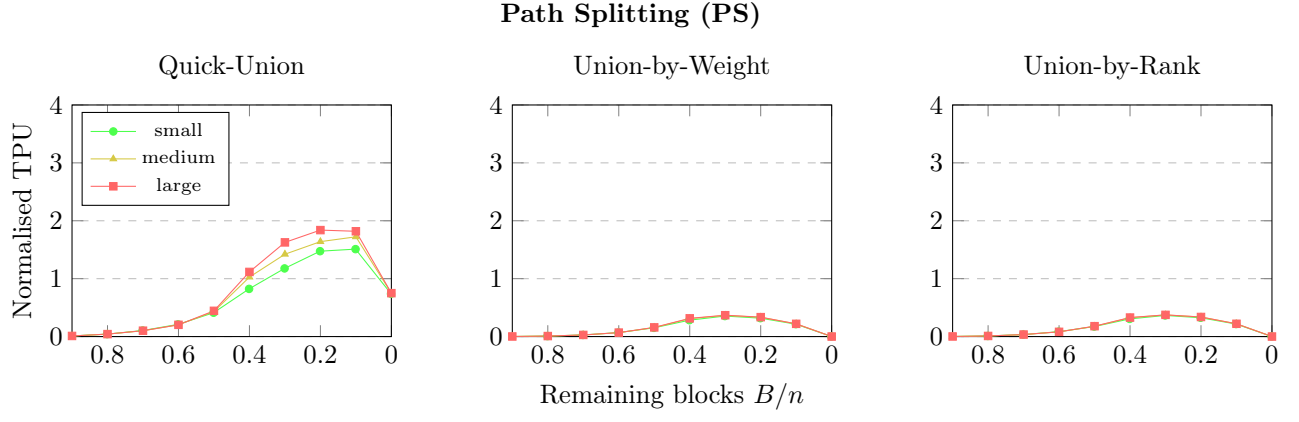


Figure 13: Normalised Total Pointer Updates required with the Path Splitting compression strategy.

7.3 Normalised Cost

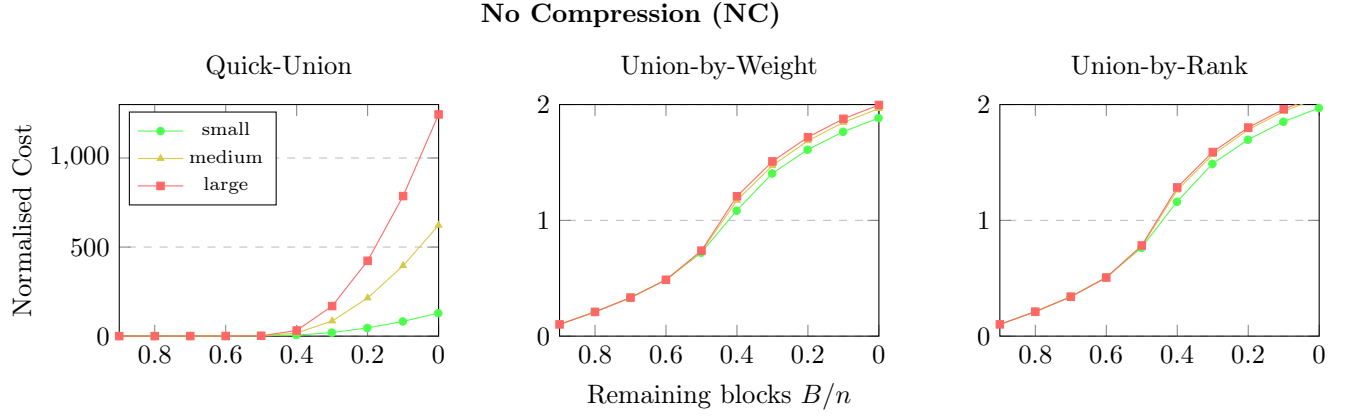


Figure 14: Normalised Cost without compression.

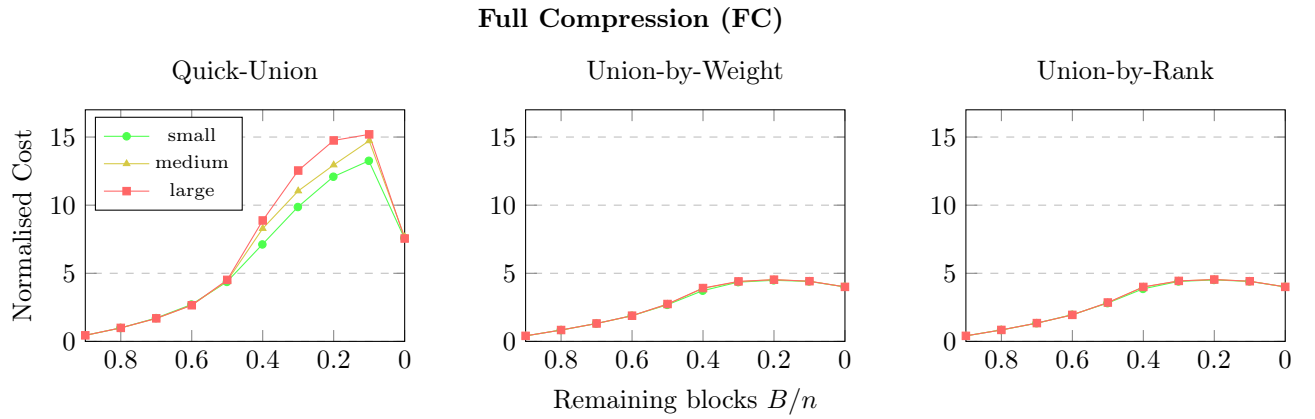


Figure 15: Normalised Cost of the Full Compression strategy.

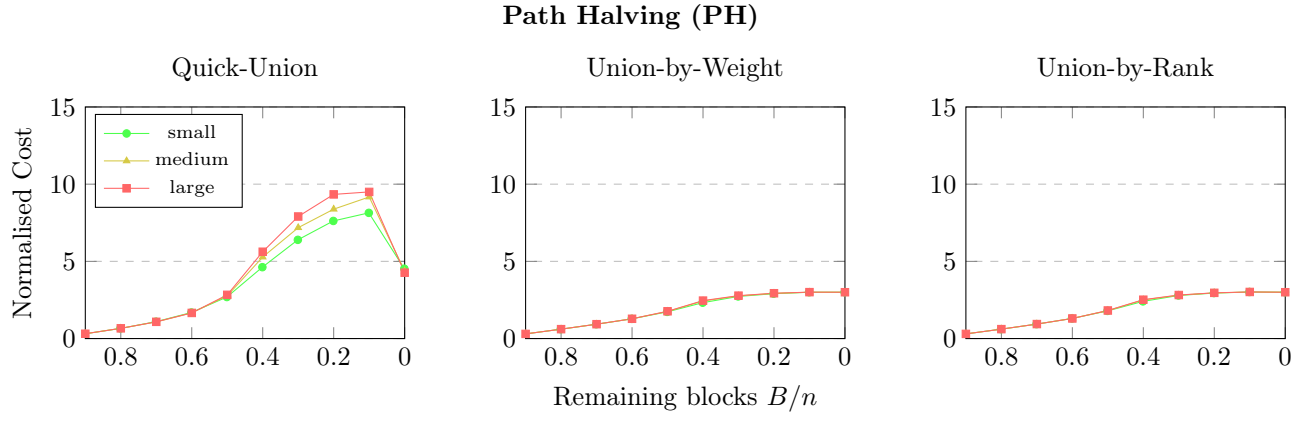


Figure 16: Normalised for Path Halving compression.

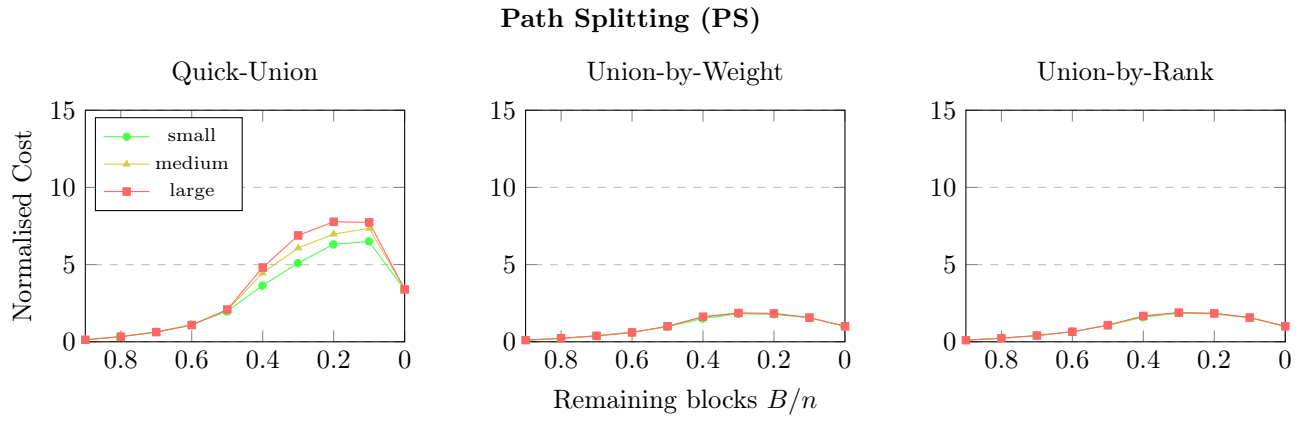


Figure 17: Normalised Cost of the Splitting compression strategy.

Appendix

```
1  #ifndef UNION_FIND_H
2  #define UNION_FIND_H
3
4  #include "DisjointSet.h"
5
6  template<class U, class F>
7  class UnionFind : public DisjointSet
8  {
9  public:
10     explicit UnionFind(int n)
11         : parent_(n), size_(n,1), rank_(n,0)
12     { for (int i=0;i<n;++i) parent_[i]=i; }
13
14     void makeSet(int x) override {
15         if (x >= (int)parent_.size()) {
16             int old = parent_.size();
17             parent_.resize(x+1);
18             size_.resize(x+1,1);
19             rank_.resize(x+1,0);
20             for (int i=old;i<=x;++i) parent_[i]=i;
21         } else {
22             parent_[x]=x; size_[x]=1; rank_[x]=0;
23         }
24     }
25
26     int find(int x) override {
27         long upd=0;
28         return F::find(x,parent_,upd);
29     }
30
31     void unionSets(int a,int b) override {
32         int r1=find(a), r2=find(b);
33         if (r1==r2) return;
34         U::unite(r1,r2,parent_,size_,rank_);
35     }
36
37     int getParent(int x) const override { return parent_[x]; }
38     int countSets() const override {
39         int c=0; for (size_t i=0;i<parent_.size();++i) if (parent_[i]==(int)i) ++c; return c;
40     }
41     int depth(int x) const override {
42         int d=0; while (x!=parent_[x]){x=parent_[x];++d;} return d;
43     }
44     long pointerUpdatesDuringFind(int x) const override {
45         return F::updatesGivenDepth(depth(x));
46     }
47
48     static double followMultiplier() { return F::followMult(); }
49
50 private:
51     std::vector<int> parent_, size_, rank_;
52 };
53
54 #endif
```

Listing 15: UnionFind.h


```

1  #ifndef DISJOINT_SET_H
2  #define DISJOINT_SET_H
3
4  class DisjointSet {
5  public:
6      virtual ~DisjointSet() = default;
7
8      virtual void makeSet(int x)                = 0;
9      virtual int find(int x)                    = 0;
10     virtual void unionSets(int x,int y)         = 0;
11
12     virtual int getParent(int x)                const = 0;
13     virtual int countSets()                     const = 0;
14     virtual int depth(int x)                    const = 0;
15     virtual long pointerUpdatesDuringFind(int x) const = 0;
16 };
17
18 #endif

```

Listing 16: DisjointSet.h

```

1  #ifndef UNION_POLICIES_H
2  #define UNION_POLICIES_H
3
4  #include <vector>
5
6  struct QU {
7      static void unite(int r1,int r2,
8                      std::vector<int>& parent,
9                      std::vector<int>& size,
10                     std::vector<int>& rank)
11      { parent[r1] = r2; }
12 };
13
14  struct UW {
15      static void unite(int r1,int r2,
16                      std::vector<int>& parent,
17                      std::vector<int>& size,
18                      std::vector<int>& rank)
19      {
20          if (size[r1] < size[r2]) std::swap(r1,r2);
21          parent[r2] = r1;
22          size[r1] += size[r2];
23      }
24 };
25
26  struct UR {
27      static void unite(int r1,int r2,
28                      std::vector<int>& parent,
29                      std::vector<int>& size,
30                      std::vector<int>& rank)
31      {
32          if (rank[r1] < rank[r2]) std::swap(r1,r2);
33          parent[r2] = r1;
34          if (rank[r1] == rank[r2]) ++rank[r1];
35      }
36 };
37 #endif

```

Listing 17: UnionPolicies.h

```

1  #ifndef FIND_POLICIES_H
2  #define FIND_POLICIES_H
3
4  #include <vector>
5
6  struct NC {
7      static int find(int x, std::vector<int>& p, long& upd)
8      {
9          while (x != p[x]) x = p[x];
10         return x;
11     }
12     static long updatesGivenDepth(int depth) { return 0; }
13     static double followMult() { return 1.0; }
14 };
15
16 struct FC {
17     static int find(int x, std::vector<int>& p, long& upd)
18     {
19         if (p[x] != x) {
20             p[x] = find(p[x], p, upd);
21             ++upd;
22         }
23         return p[x];
24     }
25     static long updatesGivenDepth(int depth) { return depth; }
26     static double followMult() { return 2.0; }
27 };
28
29 struct PS {
30     static int find(int x, std::vector<int>& p, long& upd)
31     {
32         while (p[x] != p[p[x]]) {
33             int parent = p[x];
34             p[x] = p[parent];
35             ++upd;
36             x = parent;
37         }
38         return p[x];
39     }
40     static long updatesGivenDepth(int depth) { return depth/2; }
41     static double followMult() { return 1.0; }
42 };
43
44 struct PH {
45     static int find(int x, std::vector<int>& p, long& upd)
46     {
47         while (p[x] != p[p[x]]) {
48             p[x] = p[p[x]];
49             ++upd;
50             x = p[x];
51         }
52         return p[x];
53     }
54     static long updatesGivenDepth(int depth) { return (depth+1)/2; }
55     static double followMult() { return 1.0; }
56 };
57 #endif

```

Listing 18: FindPolicies.h

```

1  #include "DisjointSet.h"
2  #include <vector>
3  #include <random>
4  #include <iostream>
5  #include <iomanip>
6  #include <cmath>
7  #include <algorithm>
8  #include <fstream>
9  #include "util/json.hpp"
10 #include "UnionPolicies.h"
11 #include "FindPolicies.h"
12 #include "UnionFind.h"
13
14 using json = nlohmann::json;
15 struct Metrics { long tpl = 0; long tpu = 0; };
16
17 Metrics measureMetrics(const DisjointSet& uf, int n)
18 {
19     Metrics m;
20     for (int v = 0; v < n; ++v) {
21         m.tpl += uf.depth(v);
22         m.tpu += uf.pointerUpdatesDuringFind(v);
23     }
24     return m;
25 }
26
27 inline std::pair<int,int> indexToPair(long long k)
28 {
29     long long i = static_cast<long long>((1 + std::sqrt(1 + 8.0*k)) / 2);
30     while (i*(i-1)/2 > k) --i;
31     long long j = k - (i*(i-1))/2;
32     return {static_cast<int>(i), static_cast<int>(j)};
33 }
34
35 class PairPermutation {
36 public:
37     PairPermutation(long long m, unsigned seed)
38         : perm_(m), idx_(0)
39     {
40         for (long long i = 0; i < m; ++i) perm_[i] = i;
41         std::shuffle(perm_.begin(), perm_.end(), std::mt19937_64(seed));
42     }
43
44     bool next(long long &k)
45     {
46         if (idx_ >= perm_.size()) return false;
47         k = perm_[idx_++];
48         return true;
49     }
50 private:
51     std::vector<long long> perm_;
52     size_t idx_;
53 };
54
55 using Factory = std::function<std::unique_ptr<DisjointSet>(int)>;
56
57 struct Variant {
58     std::string name;
59     double followMult;
60     Factory make;
61 };

```

```

62
63     const std::vector<Variant> variants = {
64         {"QU_NC", 1.0, [] (int n){ return std::make_unique<UnionFind<QU,NC>>(n); }},
65         {"QU_FC", 2.0, [] (int n){ return std::make_unique<UnionFind<QU,FC>>(n); }},
66         {"QU_PS", 1.0, [] (int n){ return std::make_unique<UnionFind<QU,PS>>(n); }},
67         {"QU_PH", 1.0, [] (int n){ return std::make_unique<UnionFind<QU,PH>>(n); }},
68         {"UW_NC", 1.0, [] (int n){ return std::make_unique<UnionFind<UW,NC>>(n); }},
69         {"UW_FC", 2.0, [] (int n){ return std::make_unique<UnionFind<UW,FC>>(n); }},
70         {"UW_PS", 1.0, [] (int n){ return std::make_unique<UnionFind<UW,PS>>(n); }},
71         {"UW_PH", 1.0, [] (int n){ return std::make_unique<UnionFind<UW,PH>>(n); }},
72         {"UR_NC", 1.0, [] (int n){ return std::make_unique<UnionFind<UR,NC>>(n); }},
73         {"UR_FC", 2.0, [] (int n){ return std::make_unique<UnionFind<UR,FC>>(n); }},
74         {"UR_PS", 1.0, [] (int n){ return std::make_unique<UnionFind<UR,PS>>(n); }},
75         {"UR_PH", 1.0, [] (int n){ return std::make_unique<UnionFind<UR,PH>>(n); }},
76     };
77
78 void runExperiment(int n, int delta, int T,
79                  bool csv, double followMult, double epsilon,
80                  const Factory& makeUF)
81 {
82     const int steps = (n - 1) / delta + 1;
83     std::vector<long long> accTPL(steps,0), accTPU(steps,0);
84     std::vector<int> accCnt(steps,0);
85
86     const unsigned baseSeed = 42;
87     for (int t = 0; t < T; ++t)
88     {
89
90         auto uf = makeUF(n);
91
92         PairPermutation perm(1LL*n*(n-1)/2, baseSeed + t);
93         long long k;
94         int nextThresh = n - delta + 1;
95         int slot = 0;
96
97         while (perm.next(k)) {
98             auto [i,j] = indexToPair(k);
99             uf->unionSets(i,j);
100
101             if (uf->countSets() <= nextThresh) {
102                 Metrics m = measureMetrics(*uf, n);
103                 accTPL[slot] += m.tpl;
104                 accTPU[slot] += m.tpu;
105                 accCnt[slot] += 1;
106
107                 nextThresh = std::max(1, nextThresh - delta);
108                 ++slot;
109                 if (uf->countSets() == 1) break;
110             }
111         }
112     }
113
114     if (!csv) {
115         std::cout << std::left
116             << std::setw(18) << "Number of Blocks"
117             << std::setw(15) << "AvgTPL"
118             << std::setw(15) << "AvgTPU"
119             << std::setw(15) << "Cost"
120             << std::setw(15) << "TPL/n"
121             << std::setw(15) << "TPU/n"
122             << std::setw(15) << "Cost/n" << '\n'
123             << std::string(108, '-') << '\n';

```

```

124     }
125
126     for (int s = 0; s < steps; ++s) {
127         if (accCnt[s] < 5) continue;
128
129         int    blocks = (s < steps-1) ? n - delta + 1 - s*delta : 1;
130         double tpl   = double(accTPL[s]) / accCnt[s];
131         double tpu    = double(accTPU[s]) / accCnt[s];
132         double cost   = followMult * tpl + epsilon * tpu;
133
134         if (csv) {
135             if (s == 0)
136                 std::cout << "Blocks,AvgTPL,AvgTPU,Cost,TPL_per_n,TPU_per_n,Cost_per_n\n";
137             std::cout << blocks << ',,'
138                 << tpl   << ',,'
139                 << tpu    << ',,'
140                 << cost   << ',,'
141                 << tpl/n  << ',,'
142                 << tpu/n  << ',,'
143                 << cost/n << '\n';
144         } else {
145             std::cout << std::left
146                 << std::setw(18) << blocks
147                 << std::setw(15) << tpl
148                 << std::setw(15) << tpu
149                 << std::setw(15) << cost
150                 << std::setw(15) << tpl/n
151                 << std::setw(15) << tpu/n
152                 << std::setw(15) << cost/n << '\n';
153         }
154     }
155 }
156
157 int main(int argc, char* argv[])
158 {
159
160     int    n          = 1000;
161     int    delta      = 100;
162     int    T          = 20;
163     bool   csv        = false;
164     double epsilon     = 2.0;
165
166     if (argc == 2) {
167         std::ifstream in(argv[1]);
168         if (!in) { std::cerr << "Cannot open config file " << argv[1] << '\n'; return 1; }
169         json cfg; in >> cfg;
170         if (cfg.contains("n"))          n          = cfg["n"];
171         if (cfg.contains("delta"))      delta      = cfg["delta"];
172         if (cfg.contains("T"))          T          = cfg["T"];
173         if (cfg.contains("csv"))        csv        = cfg["csv"];
174         if (cfg.contains("epsilon"))    epsilon    = cfg["epsilon"];
175     }
176
177     for (const auto& v : variants) {
178         std::cout << "\n=== " << v.name << " ===\n";
179         runExperiment(n, delta, T, csv, v.followMult, epsilon, v.make);
180     }
181     return 0;
182 }

```

Listing 19: main.cpp

```

1 CXX      = g++
2 CXXFLAGS = -std=c++17 -O2 -Wall -Iinclude
3 TARGET   = benchmark
4
5 SRCS = main.cpp
6 OBJS = $(SRCS:.cpp=.o)
7
8 all: $(TARGET)
9
10 $(TARGET): $(OBJS)
11     $(CXX) $(CXXFLAGS) -o $@ $^
12
13 %.o: %.cpp
14     $(CXX) $(CXXFLAGS) -c $< -o $@
15
16 clean:
17     rm -f $(TARGET) $(OBJS)
18
19 run: $(TARGET)
20     ./$$(TARGET)

```

Listing 20: Makefile