

CS1010X: Programming Methodology I

Chapter 10: Data Structures and Multiple Representations

Lim Dillion
dillionlim@u.nus.edu

2025/26 Special Term 2
National University of Singapore

10 Data Structures and Multiple Representations	3
10.1 From Functional Abstraction to Data Abstraction	4
10.1.1 Why do we need both?	4
10.2 Case Study: The Game of Nim	5
10.3 Data Abstraction for Nim Game State	6
10.3.1 Game state as an ADT	6
10.3.2 Representation 1: Encode as an integer	6
10.3.3 Representation 2: Encode as a tuple	7
10.3.4 Generalising to k piles	7
10.3.5 ADT Support Operations for the Game Loop	8
10.4 Functional Abstraction for Nim	9
10.4.1 Operations in the Process Layer: Human/Computer Moves	9
10.4.2 Announcing the Winner	9
10.4.3 Core Game Loop	10
10.5 Extensibility of Data Structures	11
10.5.1 History design	11
10.5.2 Stack ADT: Specification	11
10.5.3 Stack ADT: Implementation	11
10.5.4 Using stack for Nim undo	12
10.6 Data Structures: Specification vs Implementation	12
10.7 Multiple Representations in Real Software	13
10.8 Case Study: Complex-Arithmetic Package	14
10.8.1 Python math Library	15
10.8.2 Geometry of the two views (rectangular vs. polar)	15
10.8.3 ADT contract for Complex Arithmetic	15
10.8.4 Client arithmetic code (uses only the interface)	16
10.9 Multiple Representations	17
10.9.1 Complex Implementation A: Rectangular	17
10.9.2 Complex Implementation B: Polar	18
10.9.3 Key lesson (and how to test it)	19
10.10 Sets as an Abstract Data Type	20
10.10.1 Specification	20
10.10.2 Set Implementation 1: Unordered List	21
10.10.3 Set Implementation 2: Ordered List	23

10.10.4 Set Implementation 3: Binary Search Tree (BST)	25
10.10.5 Comparing Set Representations	27
10.11 Python sets	28
10.11.1 Creating sets and basic usage	29
10.11.2 Set methods	30
10.11.3 Hashability constraint	31
10.11.4 Connection to the Set ADT	31
10.12 Advanced: Hashing (Out of Syllabus)	32
10.12.1 Hashing	32
10.12.2 The Collision Problem & The Birthday Paradox	32
10.12.3 Collision Resolution Strategies	33
10.12.4 Load Factor (α)	34
10.12.5 Deletion and Tombstones	35
10.12.6 Complexity Summary of Hash Tables	35

Chapter 10: Data Structures and Multiple Representations

Learning Objectives

By the end of this chapter, students should be able to:

- **Explain** the difference between *functional abstraction* (hiding process details) and *data abstraction* (hiding representation details), and **justify** why both are required in non-trivial programs.
- **Decompose** a program into (i) a *process layer* (control flow, algorithms) and (ii) a *data layer* (state representation, operations), using Nim as a case study.
- **Specify** an Abstract Data Type (ADT) via a clear contract: overview, invariants/-constraints, examples, and required operations.
- **Design** an ADT interface using **constructors**, **selectors**, **predicates**, and **printer-operations**.
- **Apply** the *abstraction barrier* principle to keep client code independent of internal representation, and **identify** barrier violations.
- **Implement** a **stack** ADT and **use** it to store history for undo, including correct handling of edge cases (e.g. empty history).
- **Explain** why multiple representations co-exist in real software projects and **compare** strategies to manage them (conversion layers, tagging, dispatch).
- **Build** a **complex-arithmetic package** with a stable interface, and **implement** it using both rectangular and polar representations without changing client arithmetic code.
- **Use** Python's math library (`hypot`, `sin`, `cos`, `atan`, `atan2`) correctly and **state** that angles are measured in **radians**.
- **Implement** the **set** ADT using three representations (unordered lists, ordered lists, binary search trees) and **analyse** time/space trade-offs.
- **Use** Python's built-in `set` type and **relate** it to the set ADT, including key constraints (hashability) and typical performance expectations (average-case).

10.1 From Functional Abstraction to Data Abstraction

In earlier chapters, we relied heavily on **functional abstraction**. A function:

- hides irrelevant details of *how* something is computed,
- exposes only an interface (parameters \rightarrow return value),
- supports reuse and compositional reasoning (build big programs from small pieces).

However, most realistic programs are not just computations over isolated numbers and strings. They manage **state** and **structured information** (records, collections, histories, graphs, etc.). To write such programs cleanly, we need a comparable abstraction tool for **data**.

Abstraction Type	Hides	What does it do?
Functional abstraction	Internal steps of a computation (algorithmic details)	"How do we compute this result?"
Data abstraction	Internal representation of an object (storage details)	"How should this object be represented and manipulated?"

Table 10.1: Functional abstraction vs. data abstraction.

10.1.1 Why do we need both?

Consider a program that simulates a game, processes images, or manages student records. In each case we have:

- **Process concerns:** control flow, algorithms, iteration/recursion, input/output.
- **Data concerns:** what the objects are, what properties they have, how to store and update them.

If we mix these concerns freely, the program becomes:

- hard to extend (a small change breaks many parts),
- hard to optimise (representation changes require rewriting logic),
- hard to maintain (many places depend on the same low-level details).

The goal of this chapter is to show how **good data abstraction**:

- improves modularity,
- enables multiple representations,
- supports new features (like undo),
- and makes large programs manageable.

10.2 Case Study: The Game of Nim

Nim is a common game used in game theory, and its rules are:

- Two players take turns.
- The board consists of piles of coins.
- On your turn, remove any number of coins from exactly one pile.
- The player who takes the last coin wins.

Writing Nim cleanly requires designing:

1. **Functional (process) design:** the game loop.
2. **Data design:** the game state and history structures.

Functional Design. At a high level, we need a loop (or recursion) until the game ends. This is usually known as the game / event loop:

- Human makes a move.
- Update and display state.
- Check game over.
- Computer makes a move.
- Update and display state.
- Check game over.

Data Design. We must store:

- **current state:** number of coins in each pile,
- **history:** earlier states (for undo / replay / debugging),
- **(optional) move record:** which pile and how many were removed.

A common mistake is to jump directly into the game loop code and "figure out data later". It is often better to decide the representation of data, since it will have significant implications on how the process is designed. Therefore, to accurately design a data structure to contain the game data, we must:

1. Decide what information needs to be represented (state, history).
2. Write an interface (constructors/selectors/operations) for that representation.
3. Only then write the control flow (game loop) using the interface.

This keeps the process code clean and resilient to representation changes.

10.3 Data Abstraction for Nim Game State

10.3.1 Game state as an ADT

Conceptually, a Nim game state for two piles is the mathematical object:

$$(pile_1, pile_2)$$

with operations such as:

- create an initial state,
- read the size of a pile,
- remove coins from a chosen pile (producing a new state),
- compute total coins, check if the game is over,
- display the state.

We define a **game state ADT** interface:

- **Constructor:** `make_game_state(n, m)`
- **Selectors:** `size_of_pile(state, p)` for pile p
- **Operation:** `remove_coins_from_pile(state, k, p)`

The client code should only use these functions. It does not need to know the internal representation.

10.3.2 Representation 1: Encode as an integer

We can store (n, m) as $10n + m$. Thus, we can write the above operations as:

```
1 def make_game_state(n, m):
2     return 10*n + m
3
4 def size_of_pile(state, pile_number):
5     if pile_number == 1:
6         return state // 10
7     else:
8         return state % 10
9
10 def remove_coins_from_pile(state, num_coins, pile_number):
11     if pile_number == 1:
12         return state - 10*num_coins
13     else:
14         return state - num_coins
```

Hidden assumptions. However, there are some hidden assumptions with our representation above, that are not obvious from the interface:

- It naturally supports only **two** piles (it is not straightforward to extend this to n piles).
- It breaks once pile sizes exceed 9 if we insist on base-10 decoding.

This is a classic reason to prefer clearer representations when future changes are likely.

10.3.3 Representation 2: Encode as a tuple

```
1 def make_game_state(n, m):
2     return (n, m)
3
4 def size_of_pile(state, p):
5     return state[p-1]
6
7 def remove_coins_from_pile(state, num_coins, p):
8     n, m = state
9     if p == 1:
10        return make_game_state(n - num_coins, m)
11    else:
12        return make_game_state(n, m - num_coins)
```

Advantages. This structure has a few advantages over the previous representation:

- The structure is explicit: a state is (n, m) .
- Extending to more piles is natural: use a tuple/list of length k .
- The representation communicates fewer hidden constraints.

10.3.4 Generalising to k piles

If we want k piles, the ADT stays the same conceptually, but the constructor and selector become:

- `make_game_state(piles)` where `piles` is a list/tuple of pile sizes,
- `size_of_pile(state, p)` returns element $p - 1$.

Notice that the tuple representation is more flexible than the integer representation of piles.

Invariants should be stated

If your representation requires constraints, state them explicitly. For Nim state, typical invariants are:

- pile sizes are integers,
- pile sizes are non-negative,
- pile index p is valid (between 1 and number of piles).

10.3.5 ADT Support Operations for the Game Loop

The game loop depends on operations that:

- **query** the current game state (e.g. total coins, game-over),
- **present** the state (printer),
- **validate** moves (to keep rules checks out of the loop).

These are part of the **data abstraction layer** because they are properties of the *game state object*, and should remain correct even if the internal representation changes. Therefore, each function below should only use the ADT interface (selectors/operations), so the game loop remains representation-independent, and does not break if the representation changes.

To support both 2-pile and k -pile versions cleanly, we include a new selector:

```
1 def num_piles(state):
2     """Returns the number of piles in this game state."""
3     return 2
```

If you later generalise the representation to a tuple/list of length k , then `num_piles` becomes the only function that needs to know that detail.

We can also write some derived queries for the total size and game over states, as well as a way to display the game state.

```
1 def total_size(state):
2     """Total number of coins across all piles."""
3     total = 0
4     for p in range(1, num_piles(state) + 1):
5         total += size_of_pile(state, p)
6     return total
7 def is_game_over(state):
8     """True iff there are no coins left."""
9     return total_size(state) == 0
10 def display_game_state(state):
11     """Print the sizes of all piles (1-indexed)."""
12     print()
13     for p in range(1, num_piles(state) + 1):
14         print("Pile " + str(p) + ": " + str(size_of_pile(state, p)))
15     print()
```

These functions keep rule checks out of the game loop and out of the move policies:

```
1 def is_valid_pile(state, p):
2     return 1 <= p <= num_piles(state)
3 def is_valid_move(state, num_coins, p):
4     """A move is valid if it removes 1..size coins from valid pile."""
5     if not is_valid_pile(state, p):
6         return False
7     if num_coins <= 0:
8         return False
9     return num_coins <= size_of_pile(state, p)
```


10.4 Functional Abstraction for Nim

Now that we have a clean game-state ADT interface, we can implement the game loop as **pure process code**.

10.4.1 Operations in the Process Layer: Human/Computer Moves

The functions `human_move` and `computer_move` are **process-level policies**:

- `human_move` includes user interaction (I/O) and validation,
- `computer_move` encodes a strategy (AI policy).

They should **not** depend on representation details; they should call only ADT operations.

We can get input by using:

```
1 def prompt(msg):
2     return input(msg)
```

Then, we can implement the human move function:

```
1 def human_move(state):
2     """Ask the human for a move and return the updated state."""
3     while True:
4         p = int(prompt("Which pile will you remove from? "))
5         n = int(prompt("How many coins do you want to remove? "))
6         if is_valid_move(state, n, p):
7             return remove_coins_from_pile(state, n, p)
8         print("Invalid move. Try again.")
```

The computer uses a simple policy: pick the first non-empty pile and remove 1 coin.

```
1 def computer_move(state):
2     """A very simple strategy: remove 1 coin from the first non-empty
3     pile."""
4     for p in range(1, num_piles(state) + 1):
5         if size_of_pile(state, p) > 0:
6             print("Computer removes 1 coin from pile " + str(p))
7             return remove_coins_from_pile(state, 1, p)
8     # Should not happen if is_game_over is checked correctly.
9     return state
```

10.4.2 Announcing the Winner

In our `play(state, player)` design, `player` is the *next* player to move. If the game is already over at the start of a turn, the next player loses.

```
1 def announce_winner(player):
2     if player == "human":
3         print("You lose. Better luck next time.")
4     else:
5         print("You win. Congratulations.")
```

10.4.3 Core Game Loop

Now, the loop should rely on operations that we have previously implemented, such as:

- `display_game_state(state)`
- `is_game_over(state)`
- `human_move(state)` (returns updated state)
- `computer_move(state)` (returns updated state)
- `announce_winner(player)`

If these services respect the abstraction barrier, the loop never needs to know how state is stored. As a result, using our data abstraction above, we can write the game loop either recursively or iteratively.

```
1 def play(state, player):
2     display_game_state(state)
3     if is_game_over(state):
4         announce_winner(player)
5     elif player == "human":
6         play(human_move(state), "computer")
7     elif player == "computer":
8         play(computer_move(state), "human")
9     else:
10        print("Unknown player:", player)
```

```
1 def play_iterative(state):
2     player = "human"
3     while True:
4         display_game_state(state)
5         if is_game_over(state):
6             announce_winner(player)
7             break
8         if player == "human":
9             state = human_move(state)
10            player = "computer"
11        else:
12            state = computer_move(state)
13            player = "human"
```

Abstraction Barriers

The game loop is an example of **functional abstraction**: it should operate only on the *interface* of the game state (and related ADTs like stacks). If we change the representation of the game state (e.g. from tuples to lists), the loop should not change.

10.5 Extensibility of Data Structures

We shall see how good design allows for easy extension when requirements change. Suppose we wanted to support an undo operation, which lets a **human player** to undo their last move. Undo requires remembering previous states, so we need to modify our data representation.

10.5.1 History design

We maintain a history of game states such that:

- Before each human move, record the current state.
- When undo is requested:
 - pop the most recent state from history,
 - make it the current state again,
 - keep it the human's turn.

The natural structure is a **stack**, since we want to easily access the **previous** state.

10.5.2 Stack ADT: Specification

A stack, as covered in the previous chapter, is a **Last-In-First-Out** (LIFO) collection. We noted that we can directly use lists as stacks, but we shall write a more explicit stack ADT here.

- `make_stack()` returns a new empty stack.
- `push(s, item)` adds `item` to stack `s`.
- `pop(s)` removes and returns the most recently pushed item; returns `None` if empty.
- `is_empty(s)` returns `True` iff empty.

10.5.3 Stack ADT: Implementation

We implement the stack ADT using a list as the underlying representation. The abstraction barrier principle still applies: client code should only use `make_stack`, `push`, `pop`, and `is_empty`.

```
1 def make_stack():
2     """Return a new empty stack."""
3     return [] # internal representation: Python list
4 def is_empty(s):
5     """Return True iff stack s is empty."""
6     return len(s) == 0
7 def push(s, item):
8     """Push item onto stack s (mutates s)."""
9     s.append(item)
10    # Return value is optional; we return s for convenience.
11    return s
12 def pop(s):
13    """Pop and return the top item of stack s (mutates s).
14    Return None if s is empty."""
15    if is_empty(s):
16        return None
17    return s.pop()
```

10.5.4 Using stack for Nim undo

Now that we have the stack ADT, we can use it to implement the undo functionality. The human player enters 0 to indicate undo. Notice that other than the functions below, we do not need to modify any other functions.

```
1 game_stack = make_stack()
2
3 def human_move(state):
4     p = prompt("Which pile will you remove from? (0 for undo)")
5     if int(p) == 0:
6         return handle_undo(state)
7
8     n = prompt("How many coins do you want to remove?")
9     push(game_stack, state) # save current state before changing it
10    return remove_coins_from_pile(state, int(n), int(p))
11
12 def handle_undo(state):
13     old_state = pop(game_stack)
14     if old_state is not None:
15         print("Undo!")
16         display_game_state(old_state)
17         return old_state # human's turn continues in the loop
18     else:
19         print("No more previous moves!")
20         return state
```

10.6 Data Structures: Specification vs Implementation

Having seen the above case study, we find that a good data structure design separates:

1. **Specification (contract):** what it represents, what operations it provides, and what guarantees hold.
2. **Implementation (representation):** how it is stored and how operations are executed efficiently.

A good specification typically contains:

- **Overview:** conceptual meaning of the structure.
- **Invariants / constraints:** what must always be true.
- **Interface:** list of required operations.
- **Contracts:** algebraic laws / expected relationships.
- **Examples:** concrete sample inputs and outputs.

On the other hand, implementation typically requires:

- choosing a representation (there are often multiple choices),
- implementing constructors/selectors/predicates/operations,
- ensuring all contracts and invariants are satisfied.

10.7 Multiple Representations in Real Software

A good data structure design not only allows for extensibility with changing requirements, but the ability to handle multiple representations of data. For more complex data types, **multiple representations are possible**. A set, for example, can be represented as:

- unordered list without duplicates,
- ordered list without duplicates,
- binary search tree without duplicates,
- (in Python) a hash table (built-in `set`).

Multiple representations can co-exist in real software since:

- Large projects live a long time; requirements and performance needs change.
- Different operations need different strengths (fast membership vs fast union).
- Teams may work independently and pick different representations.

Therefore, in order to manage multiple representations, we make use of:

1. **Abstraction barrier:** insist all clients use the interface, not representation details.
2. **Conversion layer:** convert between representations at module boundaries (e.g. tree \leftrightarrow sorted list).
3. **Tagged data + dispatch (next chapter):** store a tag indicating representation type and choose the appropriate implementation at runtime.

The key engineering trade-off

Multiple representations can improve performance and flexibility, but increase complexity:

- more code paths,
- more invariants to maintain,
- more testing required.

Abstraction barriers are what keep this complexity manageable.

10.8 Case Study: Complex-Arithmetic Package

Complex numbers have two common mathematical views:

- **Rectangular form:** $z = x + iy$ (store the Cartesian coordinates (x, y))
- **Polar form:** $z = re^{iA}$ (store the magnitude-angle pair (r, A))

In real software, different teams might prefer different representations:

- Addition/subtraction are simplest in rectangular form.
- Multiplication/division are simplest in polar form.

Our goal is to write *client arithmetic code* (add/mul/etc.) that works **without knowing which representation is used**. This is exactly what **data abstraction** and the **abstraction barrier** are for.

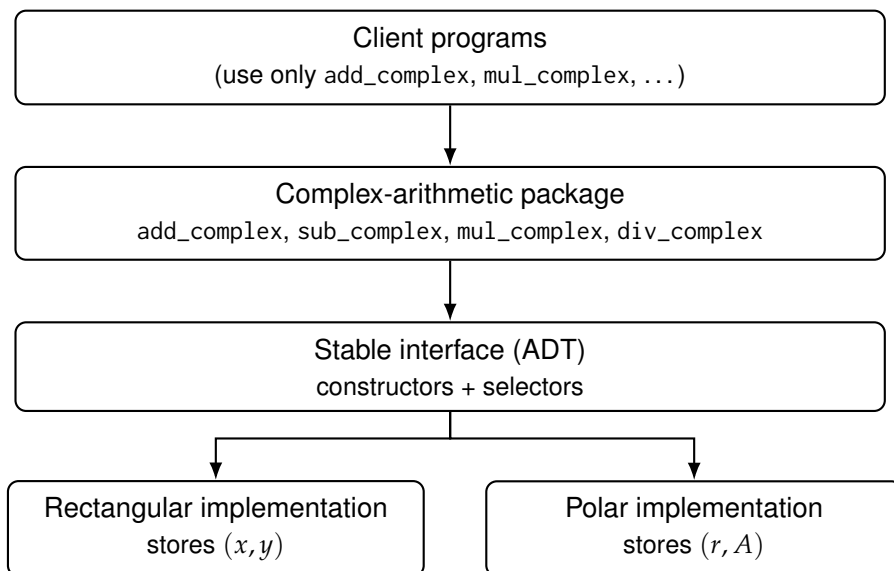


Figure 10.1: The abstraction barrier: arithmetic (clients) depend only on the interface, not on the representation.

What the barrier buys you

- You can **swap implementations** (rectangular vs. polar) without changing arithmetic code.
- You can **optimize** one layer (e.g. faster magnitude) without breaking users.
- You can evolve your system: as requirements change, you can change the representation safely.

10.8.1 Python math Library

In order to use the math library, you need to do `import math`. Useful functions are:

- `math.hypot(x, y)` returns $\sqrt{x^2 + y^2}$.
- `math.sin(a)`, `math.cos(a)` for $\sin a$ and $\cos a$ respectively.
- `math.atan(t)` computes $\arctan(t)$; `math.atan2(y, x)` computes angle of vector (x, y) .

Prefer `atan2` over `atan(y/x)`

`atan(y/x)` fails when $x = 0$ and can lose quadrant information. `atan2(y, x)` is the standard robust choice.

10.8.2 Geometry of the two views (rectangular vs. polar)

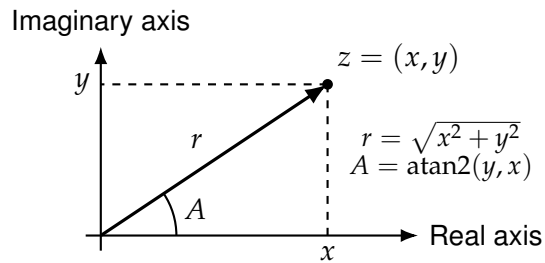


Figure 10.2: Rectangular coordinates (x, y) and polar coordinates (r, A) describe same point.

Angles and units

Python trigonometric functions use **radians**, not degrees. Use `math.radians(deg)` and `math.degrees(rad)` to convert.

10.8.3 ADT contract for Complex Arithmetic

We commit to the following interface, regardless of representation:

- **Constructors:**
 - `make_from_real_imag(x, y)`
 - `make_from_mag_ang(r, a)`
- **Selectors:**
 - `real_part(z)`, `imag_part(z)`
 - `magnitude(z)`, `angle(z)`

Representation invariants

Your implementation should ensure reasonable invariants such as:

- `magnitude(z) ≥ 0`,
- angles are kept in a standard range (e.g. $(-\pi, \pi]$).

These are *implementation choices*; clients should not rely on how you enforce them.

10.8.4 Client arithmetic code (uses only the interface)

All arithmetic is written *only* in terms of constructors and selectors.

```
1 def add_complex(z1, z2):
2     """Return z1 + z2."""
3     return make_from_real_imag(real_part(z1) + real_part(z2),
4                                imag_part(z1) + imag_part(z2))
5
6 def sub_complex(z1, z2):
7     """Return z1 - z2."""
8     return make_from_real_imag(real_part(z1) - real_part(z2),
9                                imag_part(z1) - imag_part(z2))
10
11 def mul_complex(z1, z2):
12     """Return z1 * z2."""
13     return make_from_mag_ang(magnitude(z1) * magnitude(z2),
14                              angle(z1) + angle(z2))
```

Why these formulas?

- Addition in rectangular form: $(x_1 + iy_1) + (x_2 + iy_2) = (x_1 + x_2) + i(y_1 + y_2)$.
- Multiplication in polar form: $(r_1 e^{iA_1})(r_2 e^{iA_2}) = (r_1 r_2) e^{i(A_1 + A_2)}$.

So we deliberately choose the constructor that makes the operation simplest.

Abstraction barrier: what clients must *not* do

A client must not assume `z[0]` is the real part. That is only true for one possible representation and will break for the other.

10.9 Multiple Representations

10.9.1 Complex Implementation A: Rectangular

Rectangular representation stores (x, y) directly. That is, for $z = x + yi$, it is stored as

$$z \equiv (x, y)$$

```
1  import math
2
3  def make_from_real_imag(x, y):
4      """Construct a complex number from rectangular coordinates."""
5      return (x, y)
6
7  def real_part(z):
8      """Return the real part x."""
9      return z[0]
10
11 def imag_part(z):
12     """Return the imaginary part y."""
13     return z[1]
14
15 def magnitude(z):
16     """Return r = sqrt(x^2 + y^2)."""
17     return math.hypot(real_part(z), imag_part(z))
18
19 def angle(z):
20     """Return A = atan2(y, x)."""
21     return math.atan2(imag_part(z), real_part(z))
22
23 def make_from_mag_ang(r, a):
24     """Construct from polar coordinates by converting to (x, y)."""
25     return (r * math.cos(a), r * math.sin(a))
```

Key trade-off

This representation makes `real_part` and `imag_part` constant-time and trivial, but `magnitude` and `angle` require computation.

10.9.2 Complex Implementation B: Polar

Polar representation stores (r, A) directly. That is, for $z = re^{iA}$, it is stored as

$$z \equiv (r, A)$$

Normalizing is not strictly required for correctness, but it keeps angles in a standard range. Therefore, we can write the code as:

```
1  import math
2
3  def normalize_angle(a):
4      """Return an equivalent angle in (-pi, pi]."""
5      return math.atan2(math.sin(a), math.cos(a))
6
7  def make_from_mag_ang(r, a):
8      """Construct a complex number from polar coordinates."""
9      # Invariant: r >= 0 (if r < 0, we can flip the angle by pi)
10     if r < 0:
11         r = -r
12         a = a + math.pi
13     return (r, normalize_angle(a))
14
15  def magnitude(z):
16      """Return r."""
17     return z[0]
18
19  def angle(z):
20      """Return A."""
21     return z[1]
22
23  def real_part(z):
24      """Return x = r cos(A)."""
25     return magnitude(z) * math.cos(angle(z))
26
27  def imag_part(z):
28      """Return y = r sin(A)."""
29     return magnitude(z) * math.sin(angle(z))
30
31  def make_from_real_imag(x, y):
32      """Construct from rectangular coordinates by converting to (r, A).
33         """
34     r = math.hypot(x, y)
35     a = math.atan2(y, x)
36     return make_from_mag_ang(r, a)
```

Key trade-off

This representation makes magnitude and angle constant-time and trivial, but `real_part` and `imag_part` require computation.

10.9.3 Key lesson (and how to test it)

Both implementations satisfy the *same interface*. Therefore, **arithmetic code does not change** when we swap representations.

For example, if we were to write the code below, regardless of which representation we have used, it will give us the same result.

```
1 # Choose one implementation by importing/binding its constructors/  
  selectors,  
2 # then run the same client code.  
3  
4 z1 = make_from_real_imag(3, 4)  
5 z2 = make_from_mag_ang(2, 1.0)  
6  
7 z3 = add_complex(z1, z2)  
8 z4 = mul_complex(z1, z2)  
9  
10 print(real_part(z3), imag_part(z3))  
11 print(magnitude(z4), angle(z4))
```

What should remain stable?

- The *meaning* of the results should be stable (same mathematical complex number).
- The *exact tuple contents* are allowed to differ across representations.

10.10 Sets as an Abstract Data Type

10.10.1 Specification

A **set** is an **unordered** collection of objects **without duplicates**.

- $\{3, 1, 2\} = \{1, 2, 3\}$
- duplicates are not allowed (e.g. $\{3, 3, 1, 2\}$ is not a valid set)

Therefore, a set needs the following interfaces:

- Constructors:
 - `make_set()`
 - `adjoin_set(x, S)`
- Predicates:
 - `is_empty_set(S)`
 - `is_element_of_set(x, S)`
- Operations:
 - `union_set(S, T)`
 - `intersection_set(S, T)`
- Printer:
 - `print_set(S)`

We have the following contracts that a set must follow. For any sets S, T and object x :

- `is_element_of_set(x, adjoin_set(x, S))` is True.
- Membership in union:

```
1 is_element_of_set(x, union_set(S, T))
2 # is equivalent to
3 is_element_of_set(x, S) or is_element_of_set(x, T)
```

- No object is a member of the empty set.

Adjoining Sets

The adjoin of two sets is their *union*.

For any two sets A and B , their adjoin / union, denoted by $A \cup B$, is the set containing all elements that belong to A , or to B , or to both. In set-builder notation, we write:

$$A \cup B = \{x \mid x \in A \text{ or } x \in B\}$$

10.10.2 Set Implementation 1: Unordered List

We can represent a set as a Python list with no duplicates in any order.

```
1 # Example representation:
2 S = [3, 1, 2] # valid: unordered, no duplicates
```

```
1 def make_set():
2     return []
3
4 def is_empty_set(S):
5     return len(S) == 0
6
7 def is_element_of_set(x, S):
8     for e in S:
9         if e == x:
10             return True
11     return False
12
13 def adjoin_set(x, S):
14     if is_element_of_set(x, S):
15         return S
16     return S + [x]
17
18 def union_set(S, T):
19     U = S
20     for x in T:
21         U = adjoin_set(x, U)
22     return U
23
24 def intersection_set(S, T):
25     I = make_set()
26     for x in S:
27         if is_element_of_set(x, T):
28             I = adjoin_set(x, I)
29     return I
```

Now, we analyse the time complexity of operations in this implementation of a set.

Let $n = |S|$ and $m = |T|$.

Python list concatenation is not constant-time

In Python, the expression $S + [x]$ does *not* append in place. It constructs a **new list** and copies all elements of S , then adds x . Therefore,

$$S + [x] \text{ costs } O(|S|).$$

This copying cost is a key reason why some "functional style" list-based implementations are slower than in-place versions.

Membership: `is_element_of_set(x, S)`. The function scans the list until it finds x (or reaches the end).

Best-case: $O(1)$ (first element), Worst-case: $O(n)$ (absent or last), Typical/average: $\Theta(n)$.

Adjoin: `adjoin_set(x, S)`.

- First, it performs a membership test: $O(n)$ worst-case.
- If $x \notin S$, it returns $S + [x]$, which copies S : $O(n)$.

Thus, whether x is present or absent,

`adjoin_set(x, S)` is $O(n)$ in the worst case.

Union: `union_set(S, T)`. The loop processes each of the m elements of T , calling `adjoin_set(x, U)` each time.

- **Best-case:** every $x \in T$ is already in S . Then U stays size n , and each adjoin is dominated by membership testing: $O(n)$. Hence total time $= m \cdot O(n) = O(mn)$.
- **Worst-case:** no element of T is in S . Then U grows through sizes $n, n+1, \dots, n+m-1$. Each adjoin costs $O(|U|)$, so the total time is:

$$O(n) + O(n+1) + \dots + O(n+m-1) = O\left(\sum_{i=0}^{m-1} (n+i)\right) = O(mn + m^2).$$

Therefore,

`union_set(S, T)` is $O(mn)$ best-case and $O(mn + m^2)$ worst-case.

Intersection: `intersection_set(S, T)`. The loop runs over the n elements of S . For each $x \in S$, it checks membership in T , which costs $O(m)$ worst-case:

$$n \cdot O(m) = O(nm).$$

Additionally, when x is in T , we adjoin it into I . In the worst case, all n elements of S are also in T , so I grows from size 0 to size n , and the total cost of building I by repeated adjoins is:

$$O(0) + O(1) + \dots + O(n-1) = O(n^2).$$

So,

`intersection_set(S, T)` is $O(nm)$ best-case and $O(nm + n^2)$ worst-case.

Summary.

Operation	Time (worst-case)
<code>is_element_of_set(x, S)</code>	$O(n)$
<code>adjoin_set(x, S)</code>	$O(n)$
<code>union_set(S, T)</code>	$O(mn + m^2)$
<code>intersection_set(S, T)</code>	$O(nm + n^2)$

10.10.3 Set Implementation 2: Ordered List

We represent a set as a Python list that is:

- **sorted** in increasing order, and
- contains **no duplicates**.

This representation requires elements to be comparable (support `<`, `==`).

```
1 # Example representation:
2 S = [1, 2, 3, 10] # valid: ordered, no duplicates

1 def make_set():
2     return []
3
4 def is_empty_set(S):
5     return len(S) == 0
6
7 def is_element_of_set(x, S):
8     """Return True iff x is in the ordered-list set S."""
9     for e in S:
10         if e == x:
11             return True
12         if e > x:
13             # Early termination: once we pass x, it cannot appear
14             # later
15             return False
16     return False
17
18 def adjoin_set(x, S):
19     """Return a new ordered-list set equal to S union {x}."""
20     i = 0
21     while i < len(S) and S[i] < x:
22         i += 1
23     # If x already present, return S unchanged
24     if i < len(S) and S[i] == x:
25         return S
26     # Otherwise, insert x at position i
27     return S[:i] + [x] + S[i:]
28
29 def union_set(S, T):
30     """Return union as a new ordered-list set."""
31     i, j = 0, 0
32     U = []
33     while i < len(S) and j < len(T):
34         if S[i] == T[j]:
35             U.append(S[i]); i += 1; j += 1
36         elif S[i] < T[j]:
37             U.append(S[i]); i += 1
38         else:
39             U.append(T[j]); j += 1
40     # Append remaining tail (at most one of these is non-empty)
41     return U + S[i:] + T[j:]
```

```

41 def intersection_set(S, T):
42     """Return intersection as a new ordered-list set."""
43     i, j = 0, 0
44     I = []
45     while i < len(S) and j < len(T):
46         if S[i] == T[j]:
47             I.append(S[i]); i += 1; j += 1
48         elif S[i] < T[j]:
49             i += 1
50         else:
51             j += 1
52     return I

```

Now, we analyse the time complexity of operations. Let $n = |S|$ and $m = |T|$.

Important list-slicing / concatenation costs

In Python, the following operations create **new lists** and copy elements:

- $S[:i]$ copies i elements: $O(i)$; $S[i:]$ copies $(n - i)$ elements: $O(n - i)$,
- $A + B$ copies $|A| + |B|$ elements: $O(|A| + |B|)$,
- $U + S[i:] + T[j:]$ copies all appended elements (linear in output size).

These costs matter whenever we build results without mutating in place.

Membership: `is_element_of_set(x, S)`. We scan from the beginning until:

- we find x , or
- we see an element $e > x$ (early termination), or
- we reach the end.

Best-case: $O(1)$, Worst-case: $O(n)$, Typical: often $< n$, early stopping (but still $\Theta(n)$).

Adjoin: `adjoin_set(x, S)`. There are two phases:

- find insertion index i by scanning: $O(n)$ worst-case,
- build the new list $S[:i] + [x] + S[i:]$: copying n existing elements overall, so $O(n)$.

Thus,

`adjoin_set(x, S)` is $O(n)$ worst-case.

Union: `union_set(S, T)`. The merge loop advances i or j each step, so it runs at most $n + m$ iterations. Each iteration does $O(1)$ work (comparisons and append). Appending the tails at the end copies at most $n + m$ elements overall. Therefore,

`union_set(S, T)` is $\Theta(n + m)$.

Intersection: `intersection_set(S, T)`. Similarly, the merge loop advances i or j each step, for at most $n + m$ iterations. Each iteration does $O(1)$ work, and the output size is at most $\min(n, m)$. Therefore,

$$\text{intersection_set}(S, T) \text{ is } \Theta(n + m).$$

Summary.

Operation	Time (worst-case)
<code>is_element_of_set(x, S)</code>	$O(n)$
<code>adjoin_set(x, S)</code>	$O(n)$
<code>union_set(S, T)</code>	$\Theta(n + m)$
<code>intersection_set(S, T)</code>	$\Theta(n + m)$

10.10.4 Set Implementation 3: Binary Search Tree (BST)

In the previous chapter, we already wrote the core BST operations (`bst_search`, `bst_insert`) and analysed their costs in terms of the tree height h . Here, we simply **reuse that BST implementation** to implement the **Set ADT**.

We represent a set as a BST with **no duplicates**:

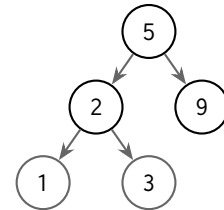
- the empty set is `None`,
- a non-empty node is a tuple (`left`, `value`, `right`).

Because the BST invariant is strict ($<$ on the left, $>$ on the right), duplicates are naturally excluded: attempting to insert an existing value simply returns the original tree.

```

1 # Example representation (using the
  previous chapter's representation):
2 S = ((None, 1, None),
3      2,
4      (None, 3, None))
5 T = (S, 5, (None, 9, None))

```



We assume `bst_search` and `bst_insert` are exactly as defined in the previous chapter.

```

1 def make_set():
2     return None
3
4 def is_empty_set(S):
5     return S is None
6
7 def is_element_of_set(x, S):
8     # Reuse BST searching
9     return bst_search(S, x)
10
11 def adjoin_set(x, S):
12     # Reuse BST insertion (duplicates are ignored)
13     return bst_insert(S, x)

```

Union and intersection. A simple (and very clear) strategy is: *traverse one tree and insert into another*. We use an in-order traversal because it is the canonical BST traversal (left, node, right), but any full traversal works.

```

1  def union_set(S, T):
2      """Return union as a BST-set."""
3      def insert_all(tree, acc):
4          if tree is None:
5              return acc
6          left, val, right = tree
7          acc = insert_all(left, acc)
8          acc = adjoin_set(val, acc)
9          acc = insert_all(right, acc)
10         return acc
11
12     return insert_all(T, S)
13
14 def intersection_set(S, T):
15     """Return intersection as a BST-set."""
16     def keep_common(tree, other, acc):
17         if tree is None:
18             return acc
19         left, val, right = tree
20         acc = keep_common(left, other, acc)
21         if is_element_of_set(val, other):
22             acc = adjoin_set(val, acc)
23         acc = keep_common(right, other, acc)
24         return acc
25
26     return keep_common(S, T, make_set())

```

BST performance depends on height

All fundamental BST operations follow a single root-to-leaf path. Thus:

$$\text{bst_search} = O(h), \quad \text{bst_insert} = O(h),$$

where h is the height of the tree.

- Balanced BST: $h = \Theta(\log n)$.
- Skewed BST: $h = \Theta(n)$ (worst case).

In this section, we state costs primarily in terms of height.

Complexity analysis. Let $n = |S|$, $m = |T|$, and let $k = |S \cap T|$. Let h_S and h_T denote the heights of the respective BSTs.

- **Membership:** `is_element_of_set(x, S)` calls `bst_search`, so it is $O(h_S)$.
- **Adjoin:** `adjoin_set(x, S)` calls `bst_insert`, so it is $O(h_S)$.

Union: `union_set(S, T)`. We traverse all m nodes of T (visiting each once), and for each visited value we insert into the accumulating tree. If the height of the accumulating tree is h_U , each insertion costs $O(h_U)$, so:

$$\text{union_set}(S, T) \text{ is } O(m \cdot h_U) \text{ (plus } O(m) \text{ traversal).}$$

Balanced typical case: $O(m \log(n + m))$. Worst-case skew: $O(m(n + m))$.

Intersection: `intersection_set(S, T)`. We traverse all n nodes of S . For each value, we test membership in T (cost $O(h_T)$). If it is present, we insert it into the result tree (there are k such values). Hence:

$$\text{intersection_set}(S, T) \text{ is } O(n \cdot h_T + k \cdot h_I)$$

where h_I is the height of the result tree. Balanced typical case: $O(n \log m + k \log k)$. Worst-case skew: $O(nm + k^2)$.

Why express costs in terms of h ?

Unlike ordered lists (which guarantee $O(n + m)$ merging), a plain BST does not automatically stay balanced. To guarantee $\Theta(\log n)$ height, we need self-balancing trees (AVL / Red-Black), which are beyond our scope here.

Summary.

Operation	Time (height-based)
<code>is_element_of_set(x, S)</code>	$O(h_S)$
<code>adjoin_set(x, S)</code>	$O(h_S)$
<code>union_set(S, T)</code>	$O(m \cdot h_U)$
<code>intersection_set(S, T)</code>	$O(n \cdot h_T + k \cdot h_I)$

10.10.5 Comparing Set Representations

Representation	Membership	Union/Intersect	Notes
Unordered list	$O(n)$	$O(nm)$ (simple)	Simplest; good for small sets.
Ordered list	$O(n)$ (early stop)	$O(n + m)$	Requires comparability; great for frequent union/intersection.
BST	$O(h)$ (avg. $\log n$)	varies	Fast membership/insertion when balanced; worst-case $O(n)$.

Table 10.2: Trade-offs among set implementations.

10.11 Python sets

Python provides a built-in `set` type: an **unordered** collection of **unique** elements. Conceptually, a set supports operations such as: membership, insertion, removal, union, and intersection.

A Python set guarantees:

- **Uniqueness:** each element appears at most once.
- **Fast membership (average-case):** `x in S` is typically $O(1)$ on average.

A Python set does *not* guarantee:

- **Order:** do not rely on iteration order for correctness.

Unordered does *not* mean "random"

Python sets are **unordered** in the sense that their iteration/printing order is **not specified by the set abstraction** and is **not guaranteed** to follow insertion order.

However, *unordered does not mean random*:

- Within a single run, a set's iteration order is typically **deterministic** given its current contents (because it is driven by the internal hash table layout).
- Small changes (adding/removing elements), different runs, or different Python versions may produce a different order.

Rule: You may *observe* a consistent order, but you must **not rely on it for correctness**. If you need a stable order, explicitly request one: `sorted(S)`.

Mathematical Set Operations

Operation	Notation	Set-Builder Definition	Description
Union	$A \cup B$	$\{x \mid x \in A \vee x \in B\}$	Elements belonging to at least one of the sets (OR).
Intersection	$A \cap B$	$\{x \mid x \in A \wedge x \in B\}$	Elements belonging strictly to both sets (AND).
Difference	$A \setminus B$	$\{x \mid x \in A \wedge x \notin B\}$	Elements in A , excluding those in B .
Symmetric Diff.	$A \Delta B$	$\{x \mid x \in (A \setminus B) \cup (B \setminus A)\}$	Elements in either A or B , but not in both (XOR).

10.11.1 Creating sets and basic usage

The most common operations for working with Python `set` objects are summarised in Table 10.3. Here $n = |S|$ is the number of elements currently in the set, and $m = |T|$ is the size of a second set (or other iterable) used in binary operations.

Operation	Syntax / Example	Time	Notes
Construct empty	<code>set()</code>	$O(1)$	Create an empty set.
Set literal	<code>{1, 2, 3}</code>	$O(n)$	Duplicates are discarded during construction.
From iterable	<code>set(iterable)</code>	$O(n)$ avg.	Inserts each element from the iterable; duplicates are ignored.
Membership	<code>x in S</code>	$O(1)$ avg.	Hash-table lookup; worst case can degrade, but average is constant-time.
Size	<code>len(S)</code>	$O(1)$	Size is stored in the set header.
Iteration	<code>for x in S: ...</code>	$O(n)$	Visits each element once. Order is consistent for a given run/state but not guaranteed .
Insert	<code>S.add(x)</code>	$O(1)$ avg.	Mutates in place. No effect if <code>x</code> already present.
Remove (error)	<code>S.remove(x)</code>	$O(1)$ avg.	Mutates in place. Raises <code>KeyError</code> if <code>x</code> not present.
Remove (no error)	<code>S.discard(x)</code>	$O(1)$ avg.	Mutates in place. Does nothing if <code>x</code> absent.
Pop arbitrary	<code>S.pop()</code>	$O(1)$ avg.	Removes and returns an arbitrary element; raises <code>KeyError</code> if empty.
Clear all	<code>S.clear()</code>	$O(n)$	Removes all elements.
Union	<code>S T</code>	$O(n + m)$ avg.	Returns a new set. Method form: <code>S.union(T)</code> .
Intersection	<code>S & T</code>	$O(n + m)$ avg.	Returns a new set. Method form: <code>S.intersection(T)</code> .
Difference	<code>S - T</code>	$O(n + m)$ avg.	Returns a new set. Method form: <code>S.difference(T)</code> .
Symmetric diff	<code>S ^ T</code>	$O(n + m)$ avg.	Returns a new set. Method form: <code>S.symmetric_difference(T)</code> .
In-place updates	<code>S = T</code>	$O(n + m)$ avg.	Mutates left operand. Analogous: <code>&=</code> , <code>-=</code> , <code>^=</code> .

Table 10.3: Basic set operations: syntax, (average) time complexity, and behaviour.

In-place updates vs. returning new sets

Operator forms such as `S | T` return a **new** set. The augmented assignments `S |= T`, `S &= T`, `S -= T`, `S ^= T` update the left-hand set **in place**. Choose in-place updates when you want to mutate an existing set object; choose the non-in-place operators when you want to keep the original sets unchanged.

10.11.2 Set methods

A set is an unordered collection with no duplicates, implemented as a hash table.

Method	Example	Time	Effect
<code>add(x)</code>	<code>s.add(10)</code>	$O(1)$ avg.	Insert element <code>x</code> into the set (no effect if already present).
<code>remove(x)</code>	<code>s.remove(10)</code>	$O(1)$ avg.	Remove <code>x</code> ; raise <code>KeyError</code> if <code>x</code> not in the set.
<code>discard(x)</code>	<code>s.discard(10)</code>	$O(1)$ avg.	Remove <code>x</code> if present; do nothing if absent (no error).
<code>pop()</code>	<code>s.pop()</code>	$O(1)$ avg.	Remove and return an arbitrary element from the set.
<code>clear()</code>	<code>s.clear()</code>	$O(n)$	Remove all elements from the set.
<code>union(t)</code>	<code>s.union(t)</code>	$O(s + t)$	Return a new set with all elements that appear in <code>s</code> , <code>t</code> , or both.
<code>intersection(t)</code>	<code>s.intersection(t)</code>	$O(s + t)$	Return a new set containing only elements common to both sets.
<code>difference(t)</code>	<code>s.difference(t)</code>	$O(s + t)$	Return a new set of elements in <code>s</code> but not in <code>t</code> .
<code>symmetric_difference(t)</code>	<code>s.symmetric_difference(t)</code>	$O(s + t)$	Return a new set of elements that are in exactly one of <code>s</code> or <code>t</code> (but not both).
<code>issubset(t)</code>	<code>s.issubset(t)</code>	$O(s)$ to $O(s + t)$	Return <code>True</code> if every element of <code>s</code> is also in <code>t</code> .
<code>issuperset(t)</code>	<code>s.issuperset(t)</code>	$O(t)$ to $O(s + t)$	Return <code>True</code> if <code>s</code> contains every element of <code>t</code> .
<code>isdisjoint(t)</code>	<code>s.isdisjoint(t)</code>	$O(\min(s , t))$	Return <code>True</code> if the two sets have no elements in common.

Average-case vs. worst-case

The " $O(1)$ average" costs assume a **good hash function** and that the table is not too full. In adversarial or pathological cases, collisions can make operations slower. Most real programs rely on the average-case behaviour.

10.11.3 Hashability constraint

Elements of Python sets must be **hashable** (roughly: usable as dictionary keys). That means:

- the object has a hash value that does not change while it is in the set,
- equality is well-defined for hashing (if $a == b$ then their hashes must match).

```
1 {(1, 2), (3, 4)}    # OK: tuples of hashable items
2 {[1, 2], [3, 4]}    # ERROR: lists are unhashable (mutable)
```

Immutability

Immutable values (integers, strings, and tuples containing only immutable items) are usually hashable. Mutable containers (like lists) are not hashable, because their contents can change, which would break the hash table's bookkeeping.

10.11.4 Connection to the Set ADT

In this chapter, we studied sets via multiple representations (unordered lists, ordered lists, BSTs) to highlight the core lesson of **data abstraction**:

- The **specification** of a set (membership, union, intersection, ...) stays stable.
- The **representation choice** changes constraints (ordering, comparability, hashability) and performance (linear vs. logarithmic vs. average constant-time).

Python's built-in **set** is one particularly efficient representation: a **hash-table-based** implementation.

10.12 Advanced: Hashing (Out of Syllabus)

Hash tables are the foundational data structure behind sets and dictionaries (hashmaps) in almost every modern programming language.

10.12.1 Hashing

A hash table relies on a function $h : U \rightarrow \{0, 1, \dots, M - 1\}$ that maps a universe of keys U to an integer index in an array of size M . To be effective, this function must satisfy three critical properties:

1. **Determinism:** For any inputs x, y , if $x = y$, then $h(x)$ must equal $h(y)$.
2. **Uniformity:** The function should map inputs as evenly as possible over the range $[0, M - 1]$ to minimize clustering.
3. **Efficiency:** Computation should be $O(1)$ (or $O(k)$ for string length k) with a minimal constant factor.

10.12.2 The Collision Problem & The Birthday Paradox

A **collision** occurs when two distinct keys map to the same index:

$$h(x) = h(y) \quad \text{for } x \neq y$$

By the Pigeonhole Principle, collisions are inevitable if $|U| > M$. However, the **Birthday Paradox** proves that collisions occur much earlier than intuition suggests.

Mathematical Insight: The Birthday Paradox

Consider a hash table with M slots and n random insertions. The probability that *no* collision occurs is:

$$P(\text{no collision}) = \frac{M}{M} \times \frac{M-1}{M} \times \dots \times \frac{M-(n-1)}{M} \approx e^{-\frac{n(n-1)}{2M}} \approx e^{-\frac{n^2}{2M}}$$

For a collision probability of 50%, we solve $0.5 \approx e^{-n^2/2M}$, which yields:

$$n \approx \sqrt{2M \ln 2} \approx 1.17\sqrt{M}$$

Implication: In a table of size $M = 10^6$, you only need ≈ 1170 elements to have a 50% chance of a collision. This is why collision handling is very important.

10.12.3 Collision Resolution Strategies

1. Separate Chaining

In this method, each bucket in the main array points to a secondary data structure (typically a linked list). Colliding elements are simply appended to the chain.

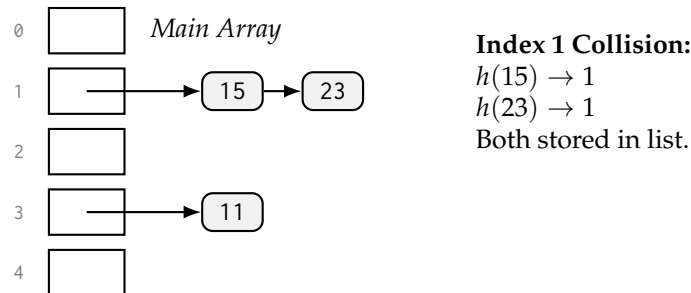


Figure 10.3: Separate chaining. Each bucket anchors a linked list of colliding elements. Performance degrades to $O(n)$ only if one chain becomes disproportionately long.

2. Open Addressing

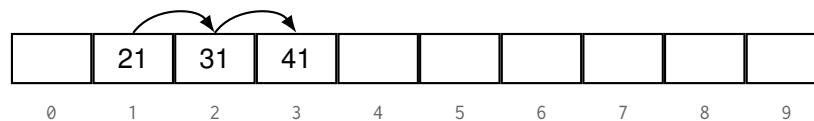
In open addressing, all elements are stored directly in the main array. If the calculated slot $h(x)$ is occupied, the algorithm strictly follows a deterministic **probing sequence** until an empty slot is found.

Common probing sequences include:

- **Linear Probing:** $i, (i + 1), (i + 2), \dots$
- **Quadratic Probing:** $i, (i + 1^2), (i + 2^2), \dots$ (reduces primary clustering).

Linear Probing Example:

$M = 10, h(x) = x \bmod 10$. Insert 21, 31, 41 (all hash to 1).



Insert(41):

1. $h(41) = 1$. Slot 1 occupied (21).
2. Probe Slot 2. Occupied (31).
3. Probe Slot 3. Empty \Rightarrow **Place 41**.

Figure 10.4: Open addressing. Collisions trigger a "probing" search for the next available slot. Deletion requires replacing items with a special TOMBSTONE marker to preserve the search path.

10.12.4 Load Factor (α)

The performance of any hash table is fundamentally governed by its **load factor**, defined as:

$$\alpha = \frac{n}{M}$$

where n is the number of stored elements and M is the number of buckets (array size).

- **In separate chaining:** α can technically exceed 1 (more items than buckets). However, as α grows, the average chain length grows, degrading lookup time to $O(\alpha)$.
- **In open addressing:** α must strictly be < 1 . As α approaches 1, the probability of collisions skyrockets.

In open addressing, the expected number of probes required to find a key is roughly $\frac{1}{1-\alpha}$.

- At $\alpha = 0.5$, we expect 2 probes.
- At $\alpha = 0.9$, we expect 10 probes.
- At $\alpha = 0.99$, we expect 100 probes.

Because of this exponential degradation, hash tables impose a strict **resizing threshold** (typically when $\alpha \approx 0.66$ or 0.75). When this limit is reached, the table doubles in size ($M \rightarrow 2M$) and every element is **rehashed** to a new position.

Mathematical Derivation: Expected Probes

Let $\alpha = n/M$ be the load factor. We assume **uniform hashing**, where the probe sequence for any key is effectively random.

The probability that any random slot is occupied is α . The probability that a slot is **empty** is $1 - \alpha$.

Finding an empty slot is analogous to flipping a biased coin until we get "heads" (where "heads" = empty slot). This follows a **geometric distribution** with success probability $p = 1 - \alpha$.

The expected number of trials E for a geometric distribution is $\frac{1}{p}$. Therefore:

$$E_{\text{probes}} \approx \frac{1}{1 - \alpha}$$

Intuitively,

- If the table is 90% full ($\alpha = 0.9$), only 10% of slots are empty ($p = 0.1$).
- On average, you must check $1/0.1 = 10$ slots to find one empty space.
- As $\alpha \rightarrow 1$, the denominator $\rightarrow 0$, causing $E \rightarrow \infty$.

10.12.5 Deletion and Tombstones

Deleting keys in **open addressing** is significantly more complex than insertion. We cannot simply mark a slot as EMPTY because doing so might break the probe chain for *other* elements that collided previously.

Broken Chains

Imagine inserting keys **A** and **B**, which both hash to index 5.

1. **Insert A:** Goes into index 5.
2. **Insert B:** Hashes to 5 (collision) → Probes to index 6.
3. **Delete A:** If we strictly empty index 5,
4. **Search B:** Hashes to 5. Sees EMPTY. IT **stops and returns "not found"**, which is incorrect!

The Solution: Tombstones. Instead of emptying the slot, we replace the deleted item with a special sentinel value often called a TOMBSTONE (or DUMMY).

The probe logic is then modified to distinguish between "new" empty slots and "recycled" tombstone slots:

1. **Search:** Treat a TOMBSTONE as **occupied**. Do not stop; continue probing.
2. **Insert:** Treat a TOMBSTONE as **empty**. You can overwrite it with the new value.

Note: Over time, a table can become cluttered with tombstones, which slows down searching (since searches step over them but do not stop). Re-hashing the table clears these tombstones.

10.12.6 Complexity Summary of Hash Tables

Operation	Average	Worst	Notes
Search	$O(1)$	$O(n)$	Degrades if many tombstones exist.
Insert	$O(1)$	$O(n)$	Amortized $O(1)$; occasional $O(n)$ spikes during resize.
Delete	$O(1)$	$O(n)$	Fast, but increases load of the table.

Python's set and dict use **open addressing** with the following specific tunings:

1. **Probing Schema:** To avoid the "clustering" issues of simple linear probing, Python uses a pseudo-random permutation based on the hash bits:

$$j = (5j + 1 + \text{perturb}) \pmod{M}$$

The perturb variable shifts bits from the hash, ensuring that two keys with similar hashes probe vastly different sequences.

2. **Growth Rate:** CPython typically resizes the table when $\alpha > 2/3$.
3. **Power of Two:** The array size M is always a power of 2 (e.g. 8, 16, 32...). This allows the expensive modulo operation ($x \pmod{M}$) to be replaced by a fast bitwise AND operation ($x \& (M-1)$).