

step	$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
	7	2	3	2	3	2	3
Init	TH						
1		T	H				
2			<u>T</u>		<u>H</u>		

Table 5.2: Part 1: Finding  $k\lambda$ ,  $f(x) = (3 \times x + 1)\%4$ ,  $x_0 = 7$ 

The Floyd's cycle-finding algorithm maintains two pointers called 'tortoise' (the slower one) at  $x_i$  and 'hare' (the faster one that keeps jumping around) at  $x_{2i}$ . Initially, both are at  $x_0$ . At each step of the algorithm, tortoise is moved *one step* to the right and the hare is moved *two steps* to the right<sup>16</sup> in the sequence. Then, the algorithm compares the sequence values at these two pointers. The smallest value of  $i > 0$  for which both tortoise and hare point to equal values is the value of  $k\lambda$  (multiple of  $\lambda$ ). We will determine the actual  $\lambda$  from  $k\lambda$  using the next two steps. In Table 5.2, when  $i = 2$ , we have  $x_2 = x_4 = x_{2+\underline{2}} = x_{2+k\lambda} = 3$ . So,  $k\lambda = 2$ . In this example, we will see below that  $k$  is eventually 1, so  $\lambda = 2$  too.

### Finding $\mu$

Next, we reset hare back to  $x_0$  and keep tortoise at its current position. Now, we advance *both* pointers to the right one step at a time, thus maintaining the  $k\lambda$  gap between the two pointers. When tortoise and hare points to the same value, we have just found the *first* repetition of length  $k\lambda$ . Since  $k\lambda$  is a multiple of  $\lambda$ , it must be true that  $x_\mu = x_{\mu+k\lambda}$ . The first time we encounter the first repetition of length  $k\lambda$  is the value of the  $\mu$ . In Table 5.3, we find that  $\mu = 1$ .

step	$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
	7	2	3	2	3	2	3
1	H		T				
2		<u>H</u>		<u>T</u>			

Table 5.3: Part 2: Finding  $\mu$ 

### Finding $\lambda$

Once we get  $\mu$ , we let tortoise stays in its current position and set hare next to it. Now, we move hare iteratively to the right one by one. Hare will point to a value that is the same as tortoise for the *first* time after  $\lambda$  steps. In Table 5.4, after hare moves once,  $x_3 = x_{3+2} = x_5 = 2$ . So,  $\lambda = 2$ .

step	$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
	7	2	3	2	3	2	3
1				T	H		
2				<u>T</u>		<u>H</u>	

Table 5.4: Part 3: Finding  $\lambda$ 

Therefore, we report  $\mu = 1$  and  $\lambda = 2$  for  $f(x) = (3 \times x + 1)\%4$  and  $x_0 = 7$ .

In overall, this algorithm runs in  $O(\mu + \lambda)$ .

<sup>16</sup>To move right one step from  $x_i$ , we use  $x_i = f(x_i)$ . To move right two steps from  $x_i$ , we use  $x_i = f(f(x_i))$ .

## The Implementation

The working C/C++ implementation of this algorithm (with comments) is shown below:

```
ii floydCycleFinding(int x0) { // function int f(int x) is defined earlier
    // 1st part: finding k*mu, hare's speed is 2x tortoise's
    int tortoise = f(x0), hare = f(f(x0)); // f(x0) is the node next to x0
    while (tortoise != hare) { tortoise = f(tortoise); hare = f(f(hare)); }
    // 2nd part: finding mu, hare and tortoise move at the same speed
    int mu = 0; hare = x0;
    while (tortoise != hare) { tortoise = f(tortoise); hare = f(hare); mu++; }
    // 3rd part: finding lambda, hare moves, tortoise stays
    int lambda = 1; hare = f(tortoise);
    while (tortoise != hare) { hare = f(hare); lambda++; }
    return ii(mu, lambda);
}
```

Source code: `ch5_07_UVa350.cpp/java`

---

**Exercise 5.7.1\*:** Richard P. Brent invented an improved version of Floyd's cycle-finding algorithm shown above. Study and implement Brent's algorithm [3]!

---



---

Programming Exercises related to Cycle-Finding:

1. UVa 00202 - Repeating Decimals (do expansion digit by digit until it cycles)
  2. UVa 00275 - Expanding Fractions (same as UVa 202 except the output format)
  3. **UVa 00350 - Pseudo-Random Numbers \*** (discussed in this section)
  4. UVa 00944 - Happy Numbers (similar to UVa 10591)
  5. UVa 10162 - Last Digit (cycle after 100 steps, use Java BigInteger to read the input, precalculate)
  6. UVa 10515 - Power et al (concentrate on the last digit)
  7. UVa 10591 - Happy Number (this sequence is 'eventually periodic')
  8. UVa 11036 - Eventually periodic ... (cycle-finding, evaluate Reverse Polish `f` with a `stack`—also see Section 9.27)
  9. **UVa 11053 - Flavius Josephus ... \*** (cycle-finding, the answer is  $N - \lambda$ )
  10. UVa 11549 - Calculator Conundrum (repeat squaring with limited digits until it cycles; that is, the Floyd's cycle-finding algorithm is only used to detect the cycle, we do not use the value of  $\mu$  or  $\lambda$ ; instead, we keep track the largest iterated function value found before any cycle is encountered)
  11. **UVa 11634 - Generate random ... \*** (use DAT of size 10K, extract digits; the programming trick to square 4 digits 'a' and get the resulting middle 4 digits is `a = (a * a / 100) % 10000`)
  12. *UVa 12464 - Professor Lazy, Ph.D.* (although  $n$  can be very huge, the pattern is actually cyclic; find the length of the cycle  $l$  and modulo  $n$  with  $l$ )
-

## 5.8 Game Theory

**Game Theory** is a mathematical model of strategic situations (not necessarily *games* as in the common meaning of ‘games’) in which a player’s success in making choices depends on the choices of *others*. Many programming problems involving game theory are classified as **Zero-Sum Game**—a mathematical way of saying that if one player wins, then the other player loses. For example, a game of Tic-Tac-Toe (e.g. UVa 10111), Chess, various number/integer games (e.g. UVa 847, 10368, 10578, 10891, 11489), and others (UVa 10165, 10404, 11311) are games with two players playing alternately (usually perfectly) and there can only be one winner.

The common question asked in programming contest problems related to game theory is whether the starting player of a two player competitive game has a winning move assuming that both players are doing **Perfect Play**. That is, each player always choose the most optimal choice available to him.

### 5.8.1 Decision Tree

One solution is to write a recursive code to explore the **Decision Tree** of the game (a.k.a. the Game Tree). If there is no overlapping subproblem, pure recursive backtracking is suitable. Otherwise, Dynamic Programming is needed. Each vertex describes the current player and the current state of the game. Each vertex is connected to all other vertices legally reachable from that vertex according to the game rules. The root vertex describes the starting player and the initial game state. If the game state at a leaf vertex is a winning state, it is a win for the current player (and a lose for the other player). At an internal vertex, the current player chooses a vertex that guarantees a win with the largest margin (or if a win is not possible, choose a vertex with the least loss). This is called the **Minimax** strategy.

For example, in UVa 10368 - Euclid’s Game, there are two players: Stan (player 0) and Ollie (player 1). The state of the game is a triple of integers  $(id, a, b)$ . The current player  $id$  can subtracts any positive multiple of the lesser of the two numbers, integer  $b$ , from the greater of the two numbers, integer  $a$ , provided that the resulting number must be nonnegative. We always maintain that  $a \geq b$ . Stan and Ollie plays alternately, until one player is able to subtract a multiple of the lesser number from the greater to reach 0, and thereby wins. The first player is Stan. The decision tree for a game with initial state  $id = 0$ ,  $a = 34$ , and  $b = 12$  is shown below in Figure 5.2 .

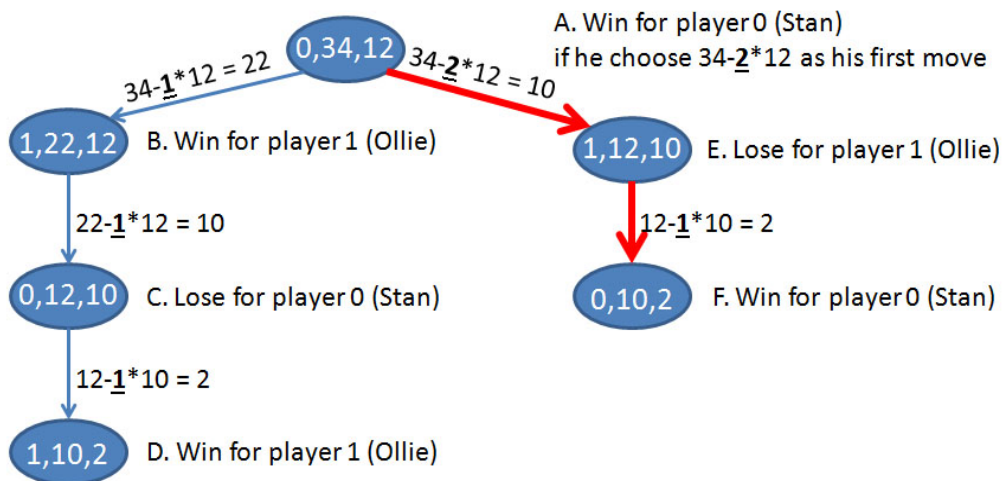


Figure 5.2: Decision Tree for an instance of ‘Euclid’s Game’

Let's trace what happens in Figure 5.2. At the root (initial state), we have triple  $(0, 34, 12)$ . At this point, player 0 (Stan) has two choices: Either to subtract  $a - b = 34 - 12 = 22$  and move to vertex  $(1, 22, 12)$  (the left branch) or to subtract  $a - 2 \times b = 34 - 2 \times 12 = 10$  and move to vertex  $(1, 12, 10)$  (the right branch). We try both choices recursively.

Let's start with the left branch. At vertex  $(1, 22, 12)$ —(Figure 5.2.B), the current player 1 (Ollie) has no choice but to subtract  $a - b = 22 - 12 = 10$ . We are now at vertex  $(0, 12, 10)$ —(Figure 5.2.C). Again, Stan only has one choice which is to subtract  $a - b = 12 - 10 = 2$ . We are now at leaf vertex  $(1, 10, 2)$ —(Figure 5.2.D). Ollie has several choices but Ollie can definitely win as  $a - 5 \times b = 10 - 5 \times 2 = 0$  and it implies that vertex  $(0, 12, 10)$  is a losing state for Stan and vertex  $(1, 22, 12)$  is a winning state for Ollie.

Now we explore the right branch. At vertex  $(1, 12, 10)$ —(Figure 5.2.E), the current player 1 (Ollie) has no choice but to subtract  $a - b = 12 - 10 = 2$ . We are now at leaf vertex  $(0, 10, 2)$ —(Figure 5.2.F). Stan has several choices but Stan can definitely win as  $a - 5 \times b = 10 - 5 \times 2 = 0$  and it implies that vertex  $(1, 12, 10)$  is a losing state for Ollie.

Therefore, for player 0 (Stan) to win this game, Stan should choose  $a - 2 \times b = 34 - 2 \times 12$  first, as this is a winning move for Stan—(Figure 5.2.A).

Implementation wise, the first integer  $id$  in the triple can be dropped as we know that depth 0 (root), 2, 4, ... are always Stan's turns and depth 1, 3, 5, ... are always Ollie's turns. This integer  $id$  is used in Figure 5.2 to simplify the explanation.

## 5.8.2 Mathematical Insights to Speed-up the Solution

Not all game theory problems can be solved by exploring the *entire* decision tree of the game, especially if the size of the tree is large. If the problem involves numbers, we may need to come up with some mathematical insights to speed up the computation.

For example, in UVa 847 - A multiplication game, there are two players: Stan (player 0) and Ollie (player 1) again. The state of the game<sup>17</sup> is an integer  $p$ . The current player can multiply  $p$  with any number between 2 to 9. Stan and Ollie also plays alternately, until one player is able to multiply  $p$  with a number between 2 to 9 such that  $p \geq n$  ( $n$  is the target number), thereby wins. The first player is Stan with  $p = 1$ .

Figure 5.3 shows an instance of this multiplication game with  $n = 17$ . Initially, player 0 has up to 8 choices (to multiply  $p = 1$  by  $[2..9]$ ). However, all of these 8 states are winning states of player 1 as player 1 can always multiply the current  $p$  by  $[2..9]$  to make  $p \geq 17$ —(Figure 5.3.B). Therefore player 0 will surely lose—(Figure 5.3.A).

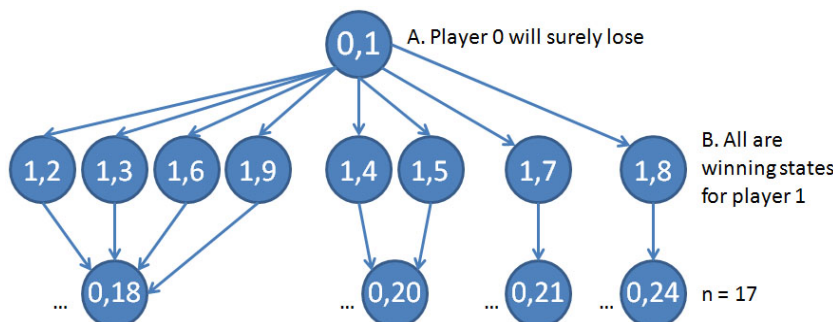


Figure 5.3: Partial Decision Tree for an instance of 'A multiplication game'

As  $1 < n < 4294967295$ , the resulting decision tree on the largest test case can be extremely huge. This is because each vertex in this decision tree has a *huge* branching factor of 8 (as

<sup>17</sup>This time we omit the player  $id$ . However, this parameter  $id$  is still shown in Figure 5.3 for clarity.

there are 8 possible numbers to choose from between 2 to 9). It is not feasible to actually explore the decision tree.

It turns out that the optimal strategy for Stan to win is to *always* multiply  $p$  with 9 (the largest possible) while Ollie will *always* multiply  $p$  with 2 (the smallest possible). Such optimization insights can be obtained by observing the pattern found in the output of smaller instances of this problem. Note that a maths-savvy contestants may want to prove this observation first before coding the solution.

### 5.8.3 Nim Game

There is a special game that is worth mentioning as it may appear in programming contests: The **Nim** game<sup>18</sup>. In Nim game, two players take turns to remove objects from distinct heaps. On each turn, a player must remove *at least one object* and may remove *any number of objects* provided they all come from the same heap. The initial state of the game is the number of objects  $n_i$  at each of the  $k$  heaps:  $\{n_1, n_2, \dots, n_k\}$ . There is a nice solution for this game. For the first (starting) player to win, the value of  $n_1 \wedge n_2 \wedge \dots \wedge n_k$  must be *non zero* where  $\wedge$  is the bit operator xor (exclusive or)—proof omitted.

---

Programming Exercises related to Game Theory:

1. UVa 00847 - A multiplication game (simulate the perfect play, discussed above)
  2. **UVa 10111 - Find the Winning ... \*** (Tic-Tac-Toe, minimax, backtracking)
  3. UVa 10165 - Stone Game (Nim game, application of Sprague-Grundy theorem)
  4. UVa 10368 - Euclid's Game (minimax, backtracking, discussed in this section)
  5. UVa 10404 - Bachet's Game (2 players game, Dynamic Programming)
  6. UVa 10578 - The Game of 31 (backtracking; try all; see who wins the game)
  7. **UVa 11311 - Exclusively Edible \*** (game theory, reducible to Nim game; we can view the game that Handel and Gretel are playing as Nim game, where there are 4 heaps - cakes left/below/right/above the topping; take the Nim sum of these 4 values and if they are equal to 0, Handel loses)
  8. **UVa 11489 - Integer Game \*** (game theory, reducible to simple math)
  9. [UVa 12293 - Box Game](#) (analyze the game tree of smaller instances to get the mathematical insight to solve this problem)
  10. [UVa 12469 - Stones](#) (game playing, Dynamic Programming, pruning)
- 

<sup>18</sup>The general form of two player games is inside the IOI syllabus [20], but Nim game is not.

## 5.9 Solution to Non-Starred Exercises

**Exercise 5.2.1:** The `<cmath>` library in C/C++ has two functions: `log` (base  $e$ ) and `log10` (base 10); `Java.lang.Math` only has `log` (base  $e$ ). To get  $\log_b(a)$  (base  $b$ ), we use the fact that  $\log_b(a) = \log(a) / \log(b)$ .

**Exercise 5.2.2:** `(int)floor(1 + log10((double)a))` returns the number of digits in decimal number  $a$ . To count the number of digits in other base  $b$ , we can use similar formula: `(int)floor(1 + log10((double)a) / log10((double)b))`.

**Exercise 5.2.3:**  $\sqrt[n]{a}$  can be rewritten as  $a^{1/n}$ . We can then use built in formula like `pow((double)a, 1.0 / (double)n)` or `exp(log((double)a) * 1.0 / (double)n)`.

**Exercise 5.3.1.1:** Possible, keep the intermediate computations **modulo**  $10^6$ . Keep chipping away the trailing zeroes (either none or a few zeroes are added after a multiplication from  $n!$  to  $(n+1)!)$ .

**Exercise 5.3.1.2:** Possible.  $9317 = 7 \times 11^3$ . We also list  $25!$  as its prime factors. Then, we check if there are one factor 7 (yes) and three factors 11 (unfortunately no). So  $25!$  is not divisible by 9317. Alternative approach: Use modulo arithmetic (see Section 5.5.8).

**Exercise 5.3.2.1:** For base number conversion of 32-bit integers, use `parseInt(String s, int radix)` and `toString(int i, int radix)` in the faster Java **Integer** class. You can also use `BufferedReader` and `BufferedWriter` for I/O (see Section 3.2.3).

**Exercise 5.4.1.1:** Binet's closed-form formula for Fibonacci:  $fib(n) = (\phi^n - (-\phi)^{-n})/\sqrt{5}$  should be correct for larger  $n$ . But since double precision data type is limited, we have discrepancies for larger  $n$ . This closed form formula is correct up to  $fib(75)$  if implemented using typical double data type in a computer program. This is unfortunately too small to be useful in typical programming contest problems involving Fibonacci numbers.

**Exercise 5.4.2.1:**  $C(n, 2) = \frac{n!}{(n-2)! \times 2!} = \frac{n \times (n-1) \times (n-2)!}{(n-2)! \times 2} = \frac{n \times (n-1)}{2} = 0.5n^2 - 0.5n = O(n^2)$ .

**Exercise 5.4.4.1:** Fundamental counting principle: If there are  $m$  ways to do one thing and  $n$  ways to do another thing, then there are  $m \times n$  ways to do both. The answer for this exercise is therefore:  $6 \times 6 \times 2 \times 2 = 6^2 \times 2^2 = 36 \times 4 = 144$  different possible outcomes.

**Exercise 5.4.4.2:** See above. The answer is:  $9 \times 9 \times 8 = 648$ . Initially there are 9 choices (1-9), then there are still 9 choices (1-9 minus 1, plus 0), then finally there are only 8 choices.

**Exercise 5.4.4.3:** A permutation is an arrangement of objects without repetition and the order is important. The formula is  ${}_nP_r = \frac{n!}{(n-r)!}$ . The answer for this exercise is therefore:  $\frac{6!}{(6-3)!} = 6 \times 5 \times 4 = 120$  3-letters words.

**Exercise 5.4.4.4:** The formula to count different permutations is:  $\frac{n!}{(n_1)! \times (n_2)! \times \dots \times (n_k)!}$  where  $n_i$  is the frequency of each unique letter  $i$  and  $n_1 + n_2 + \dots + n_k = n$ . The answer for this exercise is therefore:  $\frac{5!}{3! \times 1! \times 1!} = \frac{120}{6} = 20$  because there are 3 'B's, 1 'O', and 1 'Y'.

**Exercise 5.4.4.5:** The answers for few small  $n = 3, 4, 5, 6, 7, 8, 9$ , and 10 are 0, 1, 3, 7, 13, 22, 34, and 50, respectively. You can generate these numbers using brute force solution first. Then find the pattern and use it.

**Exercise 5.5.2.1:** Multiplying  $a \times b$  first before dividing the result by  $\gcd(a, b)$  has a higher chance of overflow in programming contest than  $a \times (b/\gcd(a, b))$ . In the example given, we have  $a = 1000000000$  and  $b = 8$ . The LCM is 1000000000—which should fit in 32-bit signed integers—can only be properly computed with  $a \times (b/\gcd(a, b))$ .

**Exercise 5.5.4.1:** Since the largest prime in `vi 'primes'` is 9999991, this code can therefore



handles  $N \leq 9999991^2 = 99999820000081 \approx 9 \times 10^{13}$ . If the smallest prime factor of  $N$  is greater than 9999991, for example,  $N = 1010189899^2 = 1020483632041630201 \approx 1 \times 10^{18}$  (this still within the capacity of 64-bit signed integer), this code will crash or produce wrong result. If we decide to drop the usage of vi ‘primes’ and uses  $PF = 3, 5, 7, \dots$  (with special check for  $PF = 2$ ), then we have a slower code and the newer limit for  $N$  is now  $N$  with smallest prime factor up to  $2^{63} - 1$ . However, if such input is given, we need to use the algorithms mentioned in **Exercise 5.5.4.2\*** and in Section 9.26.

**Exercise 5.5.4.2:** See Section 9.26.

**Exercise 5.5.5.1:** GCD(A, B) can be obtained by taking the lower power of the common prime factors of A and B. LCM(A, B) can be obtained by taking the greater power of all the prime factors of A and B. So,  $\text{GCD}(2^6 \times 3^3 \times 97^1, 2^5 \times 5^2 \times 11^2) = 2^5 = 32$  and  $\text{LCM}(2^6 \times 3^3 \times 97^1, 2^5 \times 5^2 \times 11^2) = 2^6 \times 3^3 \times 5^2 \times 11^2 \times 97^1 = 507038400$ .

**Exercise 5.5.6.1:**

```
ll numDiffPF(ll N) {
    ll PF_idx = 0, PF = primes[PF_idx], ans = 0;
    while (PF * PF <= N) {
        if (N % PF == 0) ans++;
        while (N % PF == 0) N /= PF;
        PF = primes[++PF_idx];
    }
    if (N != 1) ans++;
    return ans;
}
```

```
ll sumPF(ll N) {
    ll PF_idx = 0, PF = primes[PF_idx], ans = 0;
    while (PF * PF <= N) {
        while (N % PF == 0) { N /= PF; ans += PF; }
        PF = primes[++PF_idx];
    }
    if (N != 1) ans += N;
    return ans;
}
```

**Exercise 5.5.7.1:** The modified sieve code to compute the Euler Totient function up to  $10^6$  is shown below:

```
for (int i = 1; i <= 1000000; i++) EulerPhi[i] = i;
for (int i = 2; i <= 1000000; i++)
    if (EulerPhi[i] == i)
        for (int j = i; j <= 1000000; j += i)
            EulerPhi[j] = (EulerPhi[j] / i) * (i - 1);
```

**Exercise 5.5.8.1:** Statement 2 and 4 are not valid. The other 3 are valid.

## 5.10 Chapter Notes

This chapter has grown significantly since the first edition of this book. However, even until the third edition, we become more aware that there are still many more mathematics problems and algorithms that have not been discussed in this chapter, e.g.

- There are many more but rare **combinatorics** problems and formulas that are not yet discussed: **Burnside’s lemma**, **Stirling Numbers**, etc.
- There are other theorems, hypothesis, and conjectures that cannot be discussed one by one, e.g. **Carmichael’s function**, **Riemann’s hypothesis**, **Fermat’s Little Test**, **Chinese Remainder Theorem**, **Sprague-Grundy Theorem**, etc.
- We only briefly mention Brent’s cycle-finding algorithm (that is slightly faster than Floyd’s version) in **Exercise 5.7.1\***.
- (Computational) Geometry is also part of Mathematics, but since we have a special chapter for that, we reserve the discussions about geometry problems in Chapter 7.
- Later in Chapter 9, we briefly discuss a few more mathematics-related algorithm, e.g. Gaussian Elimination for solving system of linear equations (Section 9.9), Matrix Power and its usages (Section 9.21), Pollard’s rho algorithm (Section 9.26), Postfix Calculator and (Infix to Postfix) Conversion and (Section 9.27), Roman Numerals (Section 9.28).

There are really *many* topics about mathematics. This is not surprising since various mathematics problems have been investigated by people since hundreds years ago. Some of them are discussed in this chapter, many others are not, and yet only 1 or 2 will actually appear in a problem set. To do well in ICPC, it is a good idea to have at least *one strong mathematician* in your ICPC team in order to have those 1 or 2 mathematics problems solved. Mathematical prowess is also important for IOI contestants. Although the amount of problem-specific topics to be mastered is smaller, many IOI tasks require some form of ‘mathematical insights’.

We end this chapter by listing some pointers that may be of interest to some readers: Read number theory books, e.g. [56], investigate mathematical topics in [mathworld.wolfram.com](http://mathworld.wolfram.com) or Wikipedia, and attempt many more programming exercises related to mathematics problems like the ones in <http://projecteuler.net> [17] and <https://brilliant.org> [4].

Statistics	First Edition	Second Edition	Third Edition
Number of Pages	17	29 (+71%)	41 (+41%)
Written Exercises	-	19	20+10*=30 (+58%)
Programming Exercises	175	296 (+69%)	369 (+25%)

The breakdown of the number of programming exercises from each section is shown below:

Section	Title	Appearance	% in Chapter	% in Book
5.2	<b>Ad Hoc Mathematics ...</b>	144	39%	9%
5.3	Java BigInteger Class	45	12%	3%
5.4	Combinatorics	54	15%	3%
5.5	<b>Number Theory</b>	86	23%	5%
5.6	Probability Theory	18	5%	1%
5.7	Cycle-Finding	13	3%	1%
5.8	Game Theory	10	3%	1%





# Chapter 6

## String Processing

*The Human Genome has approximately 3.2 Giga base pairs*  
— **Human Genome Project**

### 6.1 Overview and Motivation

In this chapter, we present one more topic that is tested in ICPC—although not as frequent<sup>1</sup> as graph and mathematics problems—namely: String processing. String processing is common in the research field of *bioinformatics*. As the strings (e.g. DNA strings) that researchers deal with are usually (very) long, efficient string-specific data structures and algorithms are necessary. Some of these problems are presented as contest problems in ICPCs. By mastering the content of this chapter, ICPC contestants will have a better chance at tackling those string processing problems.

String processing tasks also appear in IOI, but usually they do not require advanced string data structures or algorithms due to syllabus [20] restriction. Additionally, the input and output format of IOI tasks are usually simple<sup>2</sup>. This eliminates the need to code tedious input parsing or output formatting commonly found in ICPC problems. IOI tasks that require string processing are usually still solvable using the problem solving paradigms mentioned in Chapter 3. It is sufficient for IOI contestants to skim through all sections in this chapter except Section 6.5 about string processing with DP. However, we believe that it may be advantageous for IOI contestants to learn some of the more advanced materials outside of their syllabus ahead of time.

This chapter is structured as follows: It starts with an overview of basic string processing skills and a *long* list of Ad Hoc string problems solvable with that basic string processing skills. Although the Ad Hoc string problems constitute the majority of the problems listed in this chapter, we have to make a remark that recent contest problems in ACM ICPC (and also IOI) usually do not ask for basic string processing solutions *except* for the ‘giveaway’ problem that most teams (contestants) should be able to solve. The more important sections are the string matching problems (Section 6.4), string processing problems solvable with Dynamic Programming (DP) (Section 6.5), and finally an extensive discussion on string processing problems where we have to deal with reasonably **long** strings (Section 6.6). The last section involves a discussion on an efficient data structure for strings like Suffix **Trie**, Suffix **Tree**, and Suffix **Array**.

---

<sup>1</sup>One potential reason: String input is harder to parse correctly and string output is harder to format correctly, making such string-based I/O less preferred over the more precise integer-based I/O.

<sup>2</sup>IOI 2010-2012 require contestants to implement functions instead of coding I/O routines.

## 6.2 Basic String Processing Skills

We begin this chapter by listing several *basic* string processing skills that every competitive programmer must have. In this section, we give a series of mini tasks that you should solve one after another without skipping. You can use any of the three programming languages: C, C++, and/or Java. Try your best to come up with the shortest, most efficient implementation that you can think of. Then, compare your implementations with ours (see the answers at the back of this chapter). If you are not surprised with any of our implementations (or can even give simpler implementations), then you are already in a good shape for tackling various string processing problems. Go ahead and read the next sections. Otherwise, please spend some time studying our implementations.

1. Given a text file that contains only alphabet characters [A-Za-z], digits [0-9], space, and period ('.'), write a program to read this text file line by line until we encounter a line that *starts with* seven periods ('.....'). Concatenate (combine) each line into one long string T. When two lines are combined, give one space between them so that the last word of the previous line is separated from the first word of the current line. There can be up to 30 characters per line and no more than 10 lines for this input block. There is no trailing space at the end of each line and each line ends with a newline character. Note: The sample input text file 'ch6.txt' is shown inside a box after question 1.(d) and before task 2.
  - (a) Do you know how to store a string in your favorite programming language?
  - (b) How to read a given text input line by line?
  - (c) How to concatenate (combine) two strings into a larger one?
  - (d) How to check if a line starts with a string '.....' to stop reading input?

```
I love CS3233 Competitive
Programming. i also love
AlGoRiThM
.....you must stop after reading this line as it starts with 7 dots
after the first input block, there will be one loooooooooooooong line...
```

2. Suppose that we have one long string T. We want to check if another string P can be found in T. Report all the indices where P appears in T or report -1 if P cannot be found in T. For example, if T = 'I love CS3233 Competitive Programming. i also love AlGoRiThM' and P = 'I', then the output is only {0} (0-based indexing). If uppercase 'I' and lowercase 'i' are considered different, then the character 'i' at index {39} is not part of the output. If P = 'love', then the output is {2, 46}. If P = 'book', then the output is {-1}.
  - (a) How to find the first occurrence of a substring in a string (if any)?  
Do we need to implement a string matching algorithm (e.g. Knuth-Morris-Pratt algorithm discussed in Section 6.4, etc) or can we just use library functions?
  - (b) How to find the next occurrence(s) of a substring in a string (if any)?
3. Suppose we want to do some simple analysis of the characters in T and also to transform each character in T into lowercase. The required analysis are: How many digits, vowels [aeiouAEIOU], and consonants (other alphabets that are not vowels) are there in T? Can you do all these in  $O(n)$  where  $n$  is the length of the string T?

4. Next, we want to break this one long string `T` into *tokens* (substrings) and store them into an array of strings called `tokens`. For this mini task, the *delimiters* of these tokens are spaces and periods (thus breaking sentences into words). For example, if we *tokenize* the string `T` (in lowercase), we will have these `tokens` = `{'i', 'love', 'cs3233', 'competitive', 'programming', 'i', 'also', 'love', 'algorithm'}`. Then, we want to sort this array of strings lexicographically<sup>3</sup> and then find the lexicographically smallest string. That is, we have sorted `tokens`: `{'algorithm', 'also', 'competitive', 'cs3233', 'i', 'i', 'love', 'love', 'programming'}`. Thus, the lexicographically smallest string for this example is `'algorithm'`.
  - (a) How to tokenize a string?
  - (b) How to store the tokens (the shorter strings) in an *array* of strings?
  - (c) How to sort an array of strings lexicographically?
5. Now, identify which word appears the most in `T`. In order to answer this query, we need to count the frequency of each word. For `T`, the output is either `'i'` or `'love'`, as both appear twice. Which data structure should be used for this mini task?
6. The given text file has one more line after a line that starts with `'.....'` but the length of this last line is not constrained. Your task is to count how many characters there are in the last line. How to read a string if its length is not known in advance?

Tasks and Source code: [ch6\\_01\\_basic\\_string.html/cpp/java](http://ch6_01_basic_string.html/cpp/java)

## Profile of Algorithm Inventors

**Donald Ervin Knuth** (born 1938) is a computer scientist and Professor Emeritus at Stanford University. He is the author of the popular Computer Science book: “*The Art of Computer Programming*”. Knuth has been called the ‘father’ of the analysis of algorithms. Knuth is also the creator of the `TeX`, the computer typesetting system used in this book.

**James Hiram Morris** (born 1941) is a Professor of Computer Science. He is a co-discoverer of the Knuth-Morris-Pratt algorithm for string-search.

**Vaughan Ronald Pratt** (born 1944) is a Professor Emeritus at Stanford University. He was one of the earliest pioneers in the field of computer science. He has made several contributions to foundational areas such as search algorithms, sorting algorithms, and primality testing. He is also a co-discoverer of the Knuth-Morris-Pratt algorithm for string-search.

**Saul B. Needleman** and **Christian D. Wunsch** jointly published the string alignment Dynamic Programming algorithm in 1970. Their DP algorithm is discussed in this book.

**Temple F. Smith** is a Professor in biomedical engineering who helped to develop the Smith-Waterman algorithm developed with Michael Waterman in 1981. The Smith-Waterman algorithm serves as the basis for multi sequence comparisons, identifying the segment with the maximum *local* sequence similarity for identifying similar DNA, RNA, and protein segments.

**Michael S. Waterman** is a Professor at the University of Southern California. Waterman is one of the founders and current leaders in the area of computational biology. His work has contributed to some of the most widely-used tools in the field. In particular, the Smith-Waterman algorithm (developed with Temple F. Smith) is the basis for many sequence comparison programs.

<sup>3</sup>Basically, this is a sort order like the one used in our common dictionary.

## 6.3 Ad Hoc String Processing Problems

Next, we continue our discussion with something light: The Ad Hoc string processing problems. They are programming contest problems involving strings that require no more than basic programming skills and perhaps some basic string processing skills discussed in Section 6.2 earlier. We only need to read the requirements in the problem description carefully and code the usually short solution. Below, we give a list of such Ad Hoc string processing problems with hints. These programming exercises have been further divided into sub-categories.

- Cipher/Encode/Encrypt/Decode/Decrypt

It is everyone's wish that their private digital communications are secure. That is, their (string) messages can only be read by the intended recipient(s). Many ciphers have been invented for this purpose and many (of the simpler ones) end up as Ad Hoc programming contest problems, each with its own encoding/decoding rules. There are many such problems in UVa online judge [47]. Thus, we have further split this category into two: the easier versus the harder ones. Try solving some of them, especially those that we classify as **must try** \*. It is interesting to learn a bit about *Computer Security/Cryptography* by solving these problems.

- Frequency Counting

In this group of problems, the contestants are asked to count the frequency of a letter (easy, use Direct Addressing Table) or a word (harder, the solution is either using a balanced Binary Search Tree—like C++ STL `map`/Java `TreeMap`—or Hash table). Some of these problems are actually related to Cryptography (the previous sub-category).

- Input Parsing

This group of problems is not for IOI contestants as the IOI syllabus enforces the input of IOI tasks to be formatted as simple as possible. However, there is no such restriction in ICPC. Parsing problems range from the simpler ones that can be dealt with an iterative parser and the more complex ones involving some grammars that requires recursive descent parser or Java `String/Pattern` class.

- Solvable with Java `String/Pattern` class (Regular Expression)

Some (but rare) string processing problems are solvable with one liner<sup>4</sup> code that use `matches(String regex)`, `replaceAll(String regex, String replacement)`, and/or other useful functions of Java `String` class. To be able to do this, one has to master the concept of Regular Expression (Regex). We will not discuss Regex in detail but we will show two usage examples:

1. In UVa 325 - Identifying Legal Pascal Real Constants, we are asked to decide if the given line of input is a legal Pascal Real constant. Suppose the line is stored in `String s`, then the following one-liner Java code is the required solution:

```
s.matches("[+-]?\\d+(\\.\\d+([eE][+-]?\\d+)?|[eE][+-]?\\d+)")
```

2. In UVa 494 - Kindergarten Counting Game, we are asked to count how many words are there in a given line. Here, a word is defined as a consecutive sequence of letters (upper and/or lower case). Suppose the line is stored in `String s`, then the following one-liner Java code is the required solution:

```
s.replaceAll("[^a-zA-Z]+", " ").trim().split(" ").length
```

---

<sup>4</sup>We can solve these problems without Regular Expression, but the code may be longer.



- Output Formatting

This is another group of problems that is also not for IOI contestants. This time, the output is the problematic one. In an ICPC problem set, such problems are used as ‘coding warm up’ or the ‘time-waster problem’ for the contestants. Practice your coding skills by solving these problems *as fast as possible* as such problems can differentiate the penalty time for each team.

- String Comparison

In this group of problems, the contestants are asked to compare strings with various criteria. This sub-category is similar to the string matching problems in the next section, but these problems mostly use `strcmp`-related functions.

- Just Ad Hoc

These are other Ad Hoc string related problems that cannot be classified as one of the other sub categories above.

Programming Exercises related to Ad Hoc String Processing:

- Cipher/Encode/Encrypt/Decode/Decrypt, Easier

1. UVa 00245 - Uncompress (use the given algorithm)
2. UVa 00306 - Cipher (can be made faster by avoiding cycle)
3. UVa 00444 - Encoder and Decoder (each char is mapped to 2 or 3 digits)
4. UVa 00458 - The Decoder (shift each character’s ASCII value by -7)
5. UVa 00483 - Word Scramble (read char by char from left to right)
6. UVa 00492 - Pig Latin (ad hoc, similar to UVa 483)
7. UVa 00641 - Do the Untwist (reverse the given formula and simulate)
8. UVa 00739 - Soundex Indexing (straightforward conversion problem)
9. UVa 00795 - Sandorf’s Cipher (prepare an ‘inverse mapper’)
10. UVa 00865 - Substitution Cypher (simple character substitution mapping)
11. UVa 10019 - Funny Encryption Method (not hard, find the pattern)
12. UVa 10222 - Decode the Mad Man (simple decoding mechanism)
13. **UVa 10851 - 2D Hieroglyphs ... \*** (ignore border; treat ‘\’ as 1/0; read from bottom)
14. **UVa 10878 - Decode the Tape \*** (treat space/‘o’ as 0/1, then it is binary to decimal conversion)
15. UVa 10896 - Known Plaintext Attack (try all possible keys; use tokenizer)
16. UVa 10921 - Find the Telephone (simple conversion problem)
17. UVa 11220 - Decoding the message (follow instruction in the problem)
18. **UVa 11278 - One-Handed Typist \*** (map QWERTY keys to DVORAK)
19. UVa 11541 - Decoding (read char by char and simulate)
20. UVa 11716 - Digital Fortress (simple cipher)
21. UVa 11787 - Numeral Hieroglyphs (follow the description)
22. UVa 11946 - Code Number (ad hoc)

- Cipher/Encode/Encrypt/Decode/Decrypt, Harder

1. UVa 00213 - Message Decoding (decrypt the message)
2. UVa 00468 - Key to Success (letter frequency mapping)
3. **UVa 00554 - Caesar Cypher \*** (try all shifts; output formatting)



4. [UVa 00632 - Compression \(II\)](#) (simulate the process, use sorting)
  5. [UVa 00726 - Decode](#) (frequency cypher)
  6. UVa 00740 - Baudot Data ... (just simulate the process)
  7. UVa 00741 - Burrows Wheeler Decoder (simulate the process)
  8. UVa 00850 - Crypt Kicker II (plaintext attack, tricky test cases)
  9. UVa 00856 - The Vigenère Cipher (3 nested loops; one for each digit)
  10. **UVa 11385 - Da Vinci Code \*** (string manipulation + Fibonacci)
  11. **UVa 11697 - Playfair Cipher \*** (follow the description, a bit tedious)
- Frequency Counting
    1. UVa 00499 - What's The Frequency ... (use 1D array for frequency counting)
    2. UVa 00895 - Word Problem (get the letter frequency of each word, compare with puzzle line)
    3. **UVa 00902 - Password Search \*** (read char by char; count word freq)
    4. UVa 10008 - What's Cryptanalysis? (character frequency count)
    5. UVa 10062 - Tell me the frequencies (ASCII character frequency count)
    6. **UVa 10252 - Common Permutation \*** (count freq of each alphabet)
    7. UVa 10293 - Word Length and Frequency (straightforward)
    8. UVa 10374 - Election (use `map` for frequency counting)
    9. UVa 10420 - List of Conquests (word frequency counting, use `map`)
    10. UVa 10625 - GNU = GNU'sNotUnix (frequency addition  $n$  times)
    11. UVa 10789 - Prime Frequency (check if a letter's frequency is prime)
    12. **UVa 11203 - Can you decide it ... \*** (problem description is convoluted, but this problem is actually easy)
    13. UVa 11577 - Letter Frequency (straightforward problem)
  - Input Parsing (Non Recursive)
    1. UVa 00271 - Simply Syntax (grammar check, linear scan)
    2. UVa 00327 - Evaluating Simple C ... (implementation can be tricky)
    3. UVa 00391 - Mark-up (use flags, tedious parsing)
    4. UVa 00397 - Equation Elation (iteratively perform the next operation)
    5. UVa 00442 - Matrix Chain Multiplication (properties of matrix chain mult)
    6. UVa 00486 - English-Number Translator (parsing)
    7. UVa 00537 - Artificial Intelligence? (simple formula; parsing is difficult)
    8. UVa 01200 - A DP Problem (LA 2972, Tehran03, tokenize linear equation)
    9. [UVa 10906 - Strange Integration \\*](#) (BNF parsing, iterative solution)
    10. UVa 11148 - Moliu Fractions (extract integers, simple/mixed fractions from a line; a bit of gcd—see Section 5.5.2)
    11. **UVa 11357 - Ensuring Truth \*** (the problem description looks scary—a SAT (satisfiability) problem; the presence of BNF grammar makes one think of recursive descent parser; however, only one clause needs to be satisfied to get TRUE; a clause can be satisfied if for all variables in the clause, its inverse is not in the clause too; now, we have a much simpler problem)
    12. **UVa 11878 - Homework Checker \*** (mathematical expression parsing)
    13. [UVa 12543 - Longest Word](#) (LA6150, HatYai12, iterative parser)

- Input Parsing (Recursive)
  1. UVa 00384 - Slurpys (recursive grammar check)
  2. UVa 00464 - Sentence/Phrase Generator (generate output based on the given BNF grammar)
  3. UVa 00620 - Cellular Structure (recursive grammar check)
  4. **UVa 00622 - Grammar Evaluation \*** (recursive BNF grammar check/evaluation)
  5. UVa 00743 - The MTM Machine (recursive grammar check)
  6. **UVa 10854 - Number of Paths \*** (recursive parsing plus counting)
  7. [UVa 11070 - The Good Old Times](#) (recursive grammar evaluation)
  8. [UVa 11291 - Smeech \\*](#) (recursive descent parser)
- Solvable with Java String/Pattern class (Regular Expression)
  1. **UVa 00325 - Identifying Legal ... \*** (see the Java solution above)
  2. **UVa 00494 - Kindergarten Counting ... \*** (see the Java solution above)
  3. UVa 00576 - Haiku Review (parsing, grammar)
  4. **UVa 10058 - Jimmi's Riddles \*** (solvable with Java regular expression)
- Output Formatting
  1. UVa 00110 - Meta-loopless sort (actually an ad hoc sorting problem)
  2. [UVa 00159 - Word Crosses](#) (tedious output formatting problem)
  3. UVa 00320 - Border (requires flood fill technique)
  4. [UVa 00330 - Inventory Maintenance](#) (use `map` to help)
  5. [UVa 00338 - Long Multiplication](#) (tedious)
  6. [UVa 00373 - Romulan Spelling](#) (check 'g' versus 'p', ad hoc)
  7. [UVa 00426 - Fifth Bank of ...](#) (tokenize; sort; reformat output)
  8. UVa 00445 - Marvelous Mazes (simulation, output formatting)
  9. **UVa 00488 - Triangle Wave \*** (use several loops)
  10. UVa 00490 - Rotating Sentences (2d array manipulation, output formatting)
  11. [UVa 00570 - Stats](#) (use `map` to help)
  12. [UVa 00645 - File Mapping](#) (use recursion to simulate directory structure, it helps the output formatting)
  13. [UVa 00890 - Maze \(II\)](#) (simulation, follow the steps, tedious)
  14. UVa 01219 - Team Arrangement (LA 3791, Tehran06)
  15. [UVa 10333 - The Tower of ASCII](#) (a real time waster problem)
  16. UVa 10500 - Robot maps (simulate, output formatting)
  17. UVa 10761 - Broken Keyboard (tricky with output formatting; note that 'END' is part of input!)
  18. **UVa 10800 - Not That Kind of Graph \*** (tedious problem)
  19. [UVa 10875 - Big Math](#) (simple but tedious problem)
  20. UVa 10894 - Save Hridoy (how fast can you can solve this problem?)
  21. UVa 11074 - Draw Grid (output formatting)
  22. [UVa 11482 - Building a Triangular ...](#) (tedious...)
  23. UVa 11965 - Extra Spaces (replace consecutive spaces with only one space)
  24. **UVa 12155 - ASCII Diamondi \*** (use proper index manipulation)
  25. [UVa 12364 - In Braille](#) (2D array check, check all possible digits [0..9])

- String Comparison
    1. UVa 00409 - Excuses, Excuses (tokenize and compare with list of excuses)
    2. **UVa 00644 - Immediate Decodability \*** (use brute force)
    3. UVa 00671 - Spell Checker (string comparison)
    4. *UVa 00912 - Live From Mars* (simulation, find and replace)
    5. **UVa 11048 - Automatic Correction ... \*** (flexible string comparison with respect to a dictionary)
    6. **UVa 11056 - Formula 1 \*** (sorting, case-insensitive string comparison)
    7. UVa 11233 - Deli Deli (string comparison)
    8. UVa 11713 - Abstract Names (modified string comparison)
    9. UVa 11734 - Big Number of Teams ... (modified string comparison)
  - Just Ad Hoc
    1. UVa 00153 - Permalex (find formula for this, similar to UVa 941)
    2. UVa 00263 - Number Chains (sort digits, convert to integers, check cycle)
    3. UVa 00892 - Finding words (basic string processing problem)
    4. **UVa 00941 - Permutations \*** (formula to get the  $n$ -th permutation)
    5. UVa 01215 - String Cutting (LA 3669, Hanoi06)
    6. UVa 01239 - Greatest K-Palindrome ... (LA 4144, Jakarta08, brute-force)
    7. UVa 10115 - Automatic Editing (simply do what they want, uses string)
    8. *UVa 10126 - Zipf's Law* (sort the words to simplify this problem)
    9. UVa 10197 - Learning Portuguese (must follow the description very closely)
    10. UVa 10361 - Automatic Poetry (read, tokenize, process as requested)
    11. UVa 10391 - Compound Words (more like data structure problem)
    12. **UVa 10393 - The One-Handed Typist \*** (follow problem description)
    13. UVa 10508 - Word Morphing (number of words = number of letters + 1)
    14. UVa 10679 - I Love Strings (the test data weak; just checking if  $T$  is a prefix of  $S$  is AC when it should not)
    15. **UVa 11452 - Dancing the Cheeky ... \*** (string period, small input, BF)
    16. UVa 11483 - Code Creator (straightforward, use 'escape character')
    17. UVa 11839 - Optical Reader (illegal if mark 0 or > 1 alternatives)
    18. UVa 11962 - DNA II (find formula; similar to UVa 941; base 4)
    19. *UVa 12243 - Flowers Flourish ...* (simple string tokenizer problem)
    20. *UVa 12414 - Calculating Yuan Fen* (brute force problem involving string)
-

## 6.4 String Matching

String *Matching* (a.k.a String *Searching*<sup>5</sup>) is a problem of finding the starting index (or indices) of a (sub)string (called *pattern*  $P$ ) in a longer string (called *text*  $T$ ). Example: Let's assume that we have  $T = \text{'STEVEN EVENT'}$ . If  $P = \text{'EVE'}$ , then the answers are index 2 and 7 (0-based indexing). If  $P = \text{'EVENT'}$ , then the answer is index 7 only. If  $P = \text{'EVENING'}$ , then there is no answer (no matching found and usually we return either -1 or NULL).

### 6.4.1 Library Solutions

For most *pure* String Matching problems on reasonably short strings, we can just use string library in our programming language. It is `strstr` in C `<string.h>`, `find` in C++ `<string>`, `indexOf` in Java `String` class. Please revisit Section 6.2, mini task 2 that discusses these string library solutions.

### 6.4.2 Knuth-Morris-Pratt's (KMP) Algorithm

In Section 1.2.3, Question 7, we have an exercise of finding all the occurrences of a substring  $P$  (of length  $m$ ) in a (long) string  $T$  (of length  $n$ ), if any. The code snippet, reproduced below with comments, is actually the *naïve* implementation of String Matching algorithm.

```
void naiveMatching() {
    for (int i = 0; i < n; i++) {           // try all potential starting indices
        bool found = true;
        for (int j = 0; j < m && found; j++) // use boolean flag 'found'
            if (i + j >= n || P[j] != T[i + j]) // if mismatch found
                found = false;                // abort this, shift the starting index i by +1
        if (found)                          // if P[0..m-1] == T[i..i+m-1]
            printf("P is found at index %d in T\n", i);
    } }
```

This naïve algorithm can run in  $O(n)$  *on average* if applied to natural text like the paragraphs of this book, but it can run in  $O(nm)$  with the worst case programming contest input like this:  $T = \text{'AAAAAAAAAAB'}$  ('A' ten times and then one 'B') and  $P = \text{'AAAAB'}$ . The naive algorithm will keep failing at the last character of pattern  $P$  and then try the next starting index which is just +1 than the previous attempt. This is not efficient. Unfortunately, a good problem author will include such test case in their secret test data.

In 1977, Knuth, Morris, and Pratt—thus the name of KMP—invented a better String Matching algorithm that makes use of the information gained by previous character comparisons, especially those that matches. KMP algorithm *never* re-compares a character in  $T$  that has matched a character in  $P$ . However, it works similar to the naïve algorithm if the *first* character of pattern  $P$  and the current character in  $T$  is a mismatch. In the example below<sup>6</sup>, comparing  $P[j]$  and  $T[i]$  and from  $i = 0$  to 13 with  $j = 0$  (the first character of  $P$ ) is no different than the naïve algorithm.

<sup>5</sup>We deal with this String Matching problem almost every time we read/edit text using computer. How many times have you pressed the well-known 'CTRL + F' button (standard Windows shortcut for the 'find feature') in typical word processing softwares, web browsers, etc?

<sup>6</sup>The sentence in string  $T$  below is just for illustration. It is not grammatically correct.

```

      1      2      3      4      5
01234567890123456789012345678901234567890
T = I DO NOT LIKE SEVENTY SEV BUT SEVENTY SEVENTY SEVEN
P = SEVENTY SEVEN
0123456789012
      1
^ the first character of P mismatch with T[i] from index i = 0 to 13
KMP has to shift the starting index i by +1, as with naive matching.
... at i = 14 and j = 0 ...

```

```

      1      2      3      4      5
01234567890123456789012345678901234567890
T = I DO NOT LIKE SEVENTY SEV BUT SEVENTY SEVENTY SEVEN
P =          SEVENTY SEVEN
          0123456789012
                1
                ^ then mismatch at index i = 25 and j = 11

```

There are 11 matches from index  $i = 14$  to  $24$ , but one mismatch at  $i = 25$  ( $j = 11$ ). The naïve matching algorithm will inefficiently restart from index  $i = 15$  but KMP can resume from  $i = 25$ . This is because the matched characters before the mismatch is ‘SEVENTY SEV’. ‘SEV’ (of length 3) appears as BOTH proper suffix and prefix of ‘SEVENTY SEV’. This ‘SEV’ is also called as the **border** of ‘SEVENTY SEV’. We can safely skip index  $i = 14$  to  $21$ : ‘SEVENTY ’ in ‘SEVENTY SEV’ as it will not match again, but we cannot rule out the possibility that the next match starts from the second ‘SEV’. So, KMP resets  $j$  back to 3, skipping  $11 - 3 = 8$  characters of ‘SEVENTY ’ (notice the trailing space), while  $i$  remains at index 25. This is the major difference between KMP and the naïve matching algorithm.

```

... at i = 25 and j = 3 (This makes KMP efficient) ...
      1      2      3      4      5
01234567890123456789012345678901234567890
T = I DO NOT LIKE SEVENTY SEV BUT SEVENTY SEVENTY SEVEN
P =          SEVENTY SEVEN
          0123456789012
                1
                ^ then immediate mismatch at index i = 25, j = 3

```

This time the prefix of  $P$  before mismatch is ‘SEV’, but it does not have a border, so KMP resets  $j$  back to 0 (or in another word, restart matching pattern  $P$  from the front again).

```

... mismatches from i = 25 to i = 29... then matches from i = 30 to i = 42 ...
      1      2      3      4      5
01234567890123456789012345678901234567890
T = I DO NOT LIKE SEVENTY SEV BUT SEVENTY SEVENTY SEVEN
P =          SEVENTY SEVEN
          0123456789012
                1

```

This is a match, so  $P = \text{‘SEVENTY SEVEN’}$  is found at index  $i = 30$ . After this, KMP knows that ‘SEVENTY SEVEN’ has ‘SEVEN’ (of length 5) as border, so KMP resets  $j$  back to 5, effectively skipping  $13 - 5 = 8$  characters of ‘SEVENTY ’ (notice the trailing space), immediately resumes the search from  $i = 43$ , and gets another match. This is efficient.

... at  $i = 43$  and  $j = 5$ , we have matches from  $i = 43$  to  $i = 50$  ...

So  $P = \text{'SEVENTY SEVEN'}$  is found again at index  $i = 38$ .

```

          1         2         3         4         5
012345678901234567890123456789012345678901234567890
T = I DO NOT LIKE SEVENTY SEV BUT SEVENTY SEVENTY SEVEN
P =                               SEVENTY SEVEN
                                0123456789012
                                  1

```

To achieve such speed up, KMP has to preprocess the pattern string and get the 'reset table'  $b$  (back). If given pattern string  $P = \text{'SEVENTY SEVEN'}$ , then table  $b$  will look like this:

```

          1
    0 1 2 3 4 5 6 7 8 9 0 1 2 3
P =  S E V E N T Y   S E V E N
b = -1 0 0 0 0 0 0 0 0 1 2 3 4 5

```

This means, if mismatch happens in  $j = 11$  (see the example above), i.e. after finding matches for 'SEVENTY SEV', then we know that we have to re-try matching  $P$  from index  $j = b[11] = 3$ , i.e. KMP now assumes that it has matched only the first three characters of 'SEVENTY SEV', which is 'SEV', because the next match can start with that prefix 'SEV'. The relatively short implementation of the KMP algorithm with comments is shown below. This implementation has a time complexity of  $O(n + m)$ .

```

#define MAX_N 100010
char T[MAX_N], P[MAX_N];           // T = text, P = pattern
int b[MAX_N], n, m;                // b = back table, n = length of T, m = length of P

void kmpPreprocess() {              // call this before calling kmpSearch()
    int i = 0, j = -1; b[0] = -1;   // starting values
    while (i < m) {                 // pre-process the pattern string P
        while (j >= 0 && P[i] != P[j]) j = b[j]; // different, reset j using b
        i++; j++;                   // if same, advance both pointers
        b[i] = j; // observe i = 8, 9, 10, 11, 12, 13 with j = 0, 1, 2, 3, 4, 5
    } }                             // in the example of P = "SEVENTY SEVEN" above

void kmpSearch() {                 // this is similar as kmpPreprocess(), but on string T
    int i = 0, j = 0;               // starting values
    while (i < n) {                 // search through string T
        while (j >= 0 && T[i] != P[j]) j = b[j]; // different, reset j using b
        i++; j++;                   // if same, advance both pointers
        if (j == m) {               // a match found when j == m
            printf("P is found at index %d in T\n", i - j);
            j = b[j];               // prepare j for the next possible match
        }
    } }

```

Source code: `ch6_02_kmp.cpp/java`

**Exercise 6.4.1\*:** Run `kmpPreprocess()` on  $P = \text{'ABABA'}$  and show the reset table  $b$ !

**Exercise 6.4.2\*:** Run `kmpSearch()` with  $P = \text{'ABABA'}$  and  $T = \text{'ACABAABABDABABA'}$ . Explain how the KMP search looks like?



### 6.4.3 String Matching in a 2D Grid

The string matching problem can also be posed in 2D. Given a 2D grid/array of characters (instead of the well-known 1D array of characters), find the occurrence(s) of pattern P in the grid. Depending on the problem requirement, the search direction can be to 4 or 8 cardinal directions, and either the pattern must be found in a straight line or it can bend. See the following example below.

```

abcdefghigg      // From UVa 10010 - Where's Waldorf?
hebkWaldork      // We can go to 8 directions, but must be straight
ftyawAldorm      // 'WALDORF' is highlighted as capital letters in the grid
ftsimrLqsrc
byoarbeDeyv      // Can you find 'BAMBI' and 'BETTY'?
klcbqwikOmk
strebghadhRb     // Can you find 'DAGBERT' in this row?
yuiqlxcnbjF

```

The solution for such string matching in a 2D grid is usually a *recursive backtracking* (see Section 3.2.2). This is because unlike the 1D counterpart where we always go to the right, at every coordinate (row, col) of the 2D grid, we have *more than one choice* to explore.

To speed up the backtracking process, usually we employ this simple pruning strategy: Once the recursion depth exceeds the length of pattern P, we can immediately prune that recursive branch. This is also called as *depth-limited search* (see Section 8.2.5).

---

#### Programming Exercises related to String Matching

- Standard
    1. UVa 00455 - Periodic String (find s in s + s)
    2. [UVa 00886 - Named Extension Dialing](#) (convert first letter of given name and all the letters of the surname into digits; then do a kind of special string matching where we want the matching to start at the *prefix* of a string)
    3. [UVa 10298 - Power Strings \\*](#) (find s in s + s, similar to UVa 455)
    4. UVa 11362 - Phone List (string sort, matching)
    5. [UVa 11475 - Extend to Palindromes \\*](#) ('border' of KMP)
    6. [UVa 11576 - Scrolling Sign \\*](#) (modified string matching; complete search)
    7. UVa 11888 - Abnormal 89's (to check 'alindrome', find reverse of s in s + s)
    8. [UVa 12467 - Secret word](#) (similar idea with UVa 11475, if you can solve that problem, you should be able to solve this problem)
  - In 2D Grid
    1. [UVa 00422 - Word Search Wonder \\*](#) (2D grid, backtracking)
    2. [UVa 00604 - The Boggle Game](#) (2D matrix, backtracking, sort, and compare)
    3. [UVa 00736 - Lost in Space](#) (2D grid, a bit modified)
    4. [UVa 10010 - Where's Waldorf? \\*](#) (discussed in this section)
    5. [UVa 11283 - Playing Boggle \\*](#) (2D grid, backtracking, do not count twice)
-

## 6.5 String Processing with Dynamic Programming

In this section, we discuss several string processing problems that are solvable with DP technique discussed in Section 3.5. The first two (String Alignment and Longest Common Subsequence) are *classical* problems and should be known by all competitive programmers. Additionally, we have added a collection of some known twists of these problems.

An important note: For various DP problems on string, we usually manipulate the *integer indices* of the strings and not the actual strings (or substrings) themselves. Passing substrings as parameters of recursive functions is strongly not recommended as it is very slow and hard to memoize.

### 6.5.1 String Alignment (Edit Distance)

The String Alignment (or Edit Distance<sup>7</sup>) problem is defined as follows: Align<sup>8</sup> two strings A with B with the maximum alignment score (or minimum number of edit operations):

After aligning A with B, there are a few possibilities between character A[i] and B[i]:

1. Character A[i] and B[i] **match** and we do nothing (assume this worth '+2' score),
2. Character A[i] and B[i] **mismatch** and we replace A[i] with B[i] (assume '-1' score),
3. We insert a space in A[i] (also '-1' score),
4. We delete a letter from A[i] (also '-1' score).

For example: (note that we use a special symbol '\_' to denote a space)

```
A = 'ACAATCC' -> 'A_CAATCC'           // Example of a non optimal alignment
B = 'AGCATGC' -> 'AGCATGC_'           // Check the optimal one below
                        2-22--2-         // Alignment Score = 4*2 + 4*-1 = 4
```

A brute force solution that tries all possible alignments will get TLE even for medium-length strings A and/or B. The solution for this problem is the Needleman-Wunsch's (bottom-up) DP algorithm [62]. Consider two strings A[1..n] and B[1..m]. We define  $V(i, j)$  to be the score of the optimal alignment between prefix A[1..i] and B[1..j] and  $score(C1, C2)$  is a function that returns the score if character C1 is aligned with character C2.

Base cases:

$V(0, 0) = 0$  // no score for matching two empty strings

$V(i, 0) = i \times score(A[i], \_)$  // delete substring A[1..i] to make the alignment,  $i > 0$

$V(0, j) = j \times score(\_, B[j])$  // insert spaces in B[1..j] to make the alignment,  $j > 0$

Recurrences: For  $i > 0$  and  $j > 0$ :

$V(i, j) = \max(option1, option2, option3)$ , where

$option1 = V(i - 1, j - 1) + score(A[i], B[j])$  // score of match or mismatch

$option2 = V(i - 1, j) + score(A[i], \_)$  // delete  $A_i$

$option3 = V(i, j - 1) + score(\_, B[j])$  // insert  $B_j$

In short, this DP algorithm concentrates on the three possibilities for the last pair of characters, which must be either a match/mismatch, a deletion, or an insertion. Although we do not know which one is the best, we can try all possibilities while avoiding the re-computation of overlapping subproblems (i.e. basically a DP technique).

<sup>7</sup>Another name for 'edit distance' is 'Levenshtein Distance'. One notable application of this algorithm is the spelling checker feature commonly found in popular text editors. If a user misspells a word, like 'probelm', then a clever text editor that realizes that this word has a very close edit distance to the correct word 'problem' can do the correction automatically.

<sup>8</sup>Align is a process of inserting spaces to strings A or B such that they have the same number of characters. You can view 'inserting spaces to B' as 'deleting the corresponding aligned characters of A'.

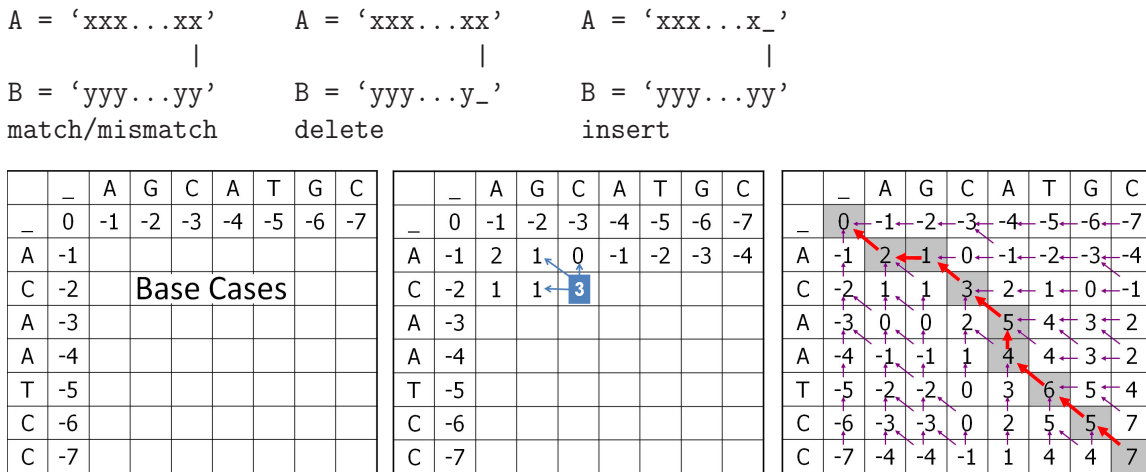


Figure 6.1: Example: A = 'ACAATCC' and B = 'AGCATGC' (alignment score = 7)

With a simple scoring function where a match gets a +2 points and mismatch, insert, delete all get a -1 point, the detail of string alignment score of A = 'ACAATCC' and B = 'AGCATGC' is shown in Figure 6.1. Initially, only the base cases are known. Then, we can fill the values row by row, left to right. To fill in  $V(i, j)$  for  $i, j > 0$ , we just need three other values:  $V(i - 1, j - 1)$ ,  $V(i - 1, j)$ , and  $V(i, j - 1)$ —see Figure 6.1, middle, row 2, column 3. The maximum alignment score is stored at the bottom right cell (7 in this example).

To reconstruct the solution, we follow the darker cells from the bottom right cell. The solution for the given strings A and B is shown below. Diagonal arrow means a match or a mismatch (e.g. the last character ..C). Vertical arrow means a deletion (e.g. ..CAA.. to ..C.A..). Horizontal arrow means an insertion (e.g. A.C.. to AGC..).

```
A = 'A_CAAAT[C]C' // Optimal alignment
B = 'AGC_AT[G]C'  // Alignment score = 5*2 + 3*-1 = 7
```

The space complexity of this (bottom-up) DP algorithm is  $O(nm)$ —the size of the DP table. We need to fill in all cells in the table in  $O(1)$  per cell. Thus, the time complexity is  $O(nm)$ .

Source code: `ch6_03_str_align.cpp/java`

**Exercise 6.5.1.1:** Why is the cost of a match +2 and the costs of replace, insert, delete are all -1? Are they magic numbers? Will +1 for match work? Can the costs for replace, insert, delete be different? Restudy the algorithm and discover the answer.

**Exercise 6.5.1.2:** The example source code: `ch6_03_str_align.cpp/java` only show the optimal alignment *score*. Modify the given code to actually show the *actual alignment*!

**Exercise 6.5.1.3:** Show how to use the 'space saving trick' shown in Section 3.5 to improve this Needleman-Wunsch's (bottom-up) DP algorithm! What will be the new space and time complexity of your solution? What is the drawback of using such a formulation?

**Exercise 6.5.1.4:** The String Alignment problem in this section is called the **global** alignment problem and runs in  $O(nm)$ . If the given contest problem is limited to  $d$  insertions or deletions only, we can have a faster algorithm. Find a simple tweak to the Needleman-Wunsch's algorithm so that it performs at most  $d$  insertions or deletions and runs faster!

**Exercise 6.5.1.5:** Investigate the improvement of Needleman-Wunsch's algorithm (the **Smith-Waterman's** algorithm [62]) to solve the **local** alignment problem!

### 6.5.2 Longest Common Subsequence

The Longest Common Subsequence (LCS) problem is defined as follows: Given two strings A and B, what is the longest common subsequence between them. For example, A = 'ACAATCC' and B = 'AGCATGC' have LCS of length 5, i.e. 'ACATC'.

This LCS problem can be reduced to the String Alignment problem presented earlier, so we can use the same DP algorithm. We set the cost for mismatch as negative infinity (e.g. -1 Billion), cost for insertion and deletion as 0, and the cost for match as 1. This makes the Needleman-Wunsch's algorithm for String Alignment to never consider mismatches.

---

**Exercise 6.5.2.1:** What is the LCS of A = 'apple' and B = 'people'?

**Exercise 6.5.2.2:** The Hamming distance problem, i.e. finding the number of different characters between two equal-length strings, can be reduced to String Alignment problem. Assign an appropriate cost to match, mismatch, insert, and delete so that we can compute the Hamming distance between two strings using Needleman-Wunsch's algorithm!

**Exercise 6.5.2.3:** The LCS problem can be solved in  $O(n \log k)$  when all characters are distinct, e.g. if you are given two permutations as in UVa 10635. Solve this variant!

---

### 6.5.3 Non Classical String Processing with DP

#### UVa 11151 - Longest Palindrome

A palindrome is a string that can be read the same way in either direction. Some variants of palindrome finding problems are solvable with DP technique, e.g. UVa 11151 - Longest Palindrome: Given a string of up to  $n = 1000$  characters, determine the length of the longest palindrome that you can make from it by deleting zero or more characters. Examples:

'ADAM' → 'ADA' (of length 3, delete 'M')

'MADAM' → 'MADAM' (of length 5, delete nothing)

'NEVERODDOREVENING' → 'NEVERODDOREVEN' (of length 14, delete 'ING')

'RACEF1CARFAST' → 'RACECAR' (of length 7, delete 'F1' and 'FAST')

The DP solution: let  $len(l, r)$  be the length of the longest palindrome from string  $A[1..r]$ .

Base cases:

If  $(l = r)$ , then  $len(l, r) = 1$ . // odd-length palindrome

If  $(l + 1 = r)$ , then  $len(l, r) = 2$  if  $(A[l] = A[r])$ , or 1 otherwise. // even-length palindrome

Recurrences:

If  $(A[l] = A[r])$ , then  $len(l, r) = 2 + len(l + 1, r - 1)$ . // both corner characters are the same  
 else  $len(l, r) = \max(len(l, r - 1), len(l + 1, r))$ . // increase left side or decrease right side

This DP solution has time complexity of  $O(n^2)$ .

---

**Exercise 6.5.3.1\*:** Can we use the Longest Common Subsequence solution shown in Section 6.5.2 to solve UVa 11151? If we can, how? What is the time complexity?

**Exercise 6.5.3.2\*:** Suppose that we are now interested to find the longest palindrome in a given string with length up to  $n = 10000$  characters. This time, we are not allowed to delete any character. What should be the solution?

---

---

Programming Exercises related to String Processing with DP:

- Classic
    1. UVa 00164 - String Computer (String Alignment/Edit Distance)
    2. **UVa 00526 - Edit Distance \*** (String Alignment/Edit Distance)
    3. UVa 00531 - Compromise (Longest Common Subsequence; print the solution)
    4. UVa 01207 - AGTC (LA 3170, Manila06, classical String Edit problem)
    5. UVa 10066 - The Twin Towers (Longest Common Subsequence problem, but not on 'string')
    6. UVa 10100 - Longest Match (Longest Common Subsequence)
    7. **UVa 10192 - Vacation \*** (Longest Common Subsequence)
    8. UVa 10405 - Longest Common ... (Longest Common Subsequence)
    9. **UVa 10635 - Prince and Princess \*** (find LCS of two permutations)
    10. UVa 10739 - String to Palindrome (variation of edit distance)
  - Non Classic
    1. *UVa 00257 - Palindromes* (standard DP palindrome plus brute force checks)
    2. *UVa 10453 - Make Palindrome* (s: (L, R); t: (L+1, R-1) if  $S[L] == S[R]$ ; or one plus min of (L + 1, R) or (L, R - 1); also print the required solution)
    3. UVa 10617 - Again Palindrome (manipulate indices, not the actual string)
    4. *UVa 11022 - String Factoring \** (s: the min weight of substring [i..j])
    5. **UVa 11151 - Longest Palindrome \*** (discussed in this section)
    6. **UVa 11258 - String Partition \*** (discussed in this section)
    7. *UVa 11552 - Fewest Flops* ( $dp(i, c)$  = minimum number of chunks after considering the first i segments ending with character c)
- 

## Profile of Algorithm Inventors

**Udi Manber** is an Israeli computer scientist. He works in Google as one of their vice presidents of engineering. Along with Gene Myers, Manber invented Suffix Array data structure in 1991.

**Eugene “Gene” Wimberly Myers, Jr.** is an American computer scientist and bioinformatician, who is best known for his development of the BLAST (Basic Local Alignment Search Tool) tool for sequence analysis. His 1990 paper that describes BLAST has received over 24000 citations making it among the most highly cited paper ever. He also invented Suffix Array with Udi Manber.

## 6.6 Suffix Trie/Tree/Array

Suffix Trie, Suffix Tree, and Suffix Array are efficient and related data structures for strings. We do not discuss this topic in Section 2.4 as these data structures are unique to strings.

### 6.6.1 Suffix Trie and Applications

The **suffix**  $i$  (or the  $i$ -th suffix) of a string is a ‘special case’ of substring that goes from the  $i$ -th character of the string up to the *last* character of the string. For example, the 2-th suffix of ‘STEVEN’ is ‘EVEN’, the 4-th suffix of ‘STEVEN’ is ‘EN’ (0-based indexing).

A **Suffix Trie**<sup>9</sup> of a set of strings  $S$  is a tree of all possible suffixes of strings in  $S$ . Each edge label represents a character. Each vertex represents a suffix indicated by its path label: A sequence of edge labels from root to that vertex. Each vertex is connected to (some of) the other 26 vertices (assuming that we only use upper-case Latin letters) according to the suffixes of strings in  $S$ . The common prefix of two suffixes is shared. Each vertex has two boolean flags. The first/second one is to indicate that there exists a suffix/word in  $S$  *terminating* in that vertex, respectively. Example: If we have  $S = \{\text{‘CAR’}, \text{‘CAT’}, \text{‘RAT’}\}$ , we have the following suffixes  $\{\text{‘CAR’}, \text{‘AR’}, \text{‘R’}, \text{‘CAT’}, \text{‘AT’}, \text{‘T’}, \text{‘RAT’}, \text{‘AT’}, \text{‘T’}\}$ . After sorting and removing duplicates, we have:  $\{\text{‘AR’}, \text{‘AT’}, \text{‘CAR’}, \text{‘CAT’}, \text{‘R’}, \text{‘RAT’}, \text{‘T’}\}$ . Figure 6.2 shows the Suffix Trie with 7 suffix terminating vertices (filled circles) and 3 word terminating vertices (filled circles indicated with label ‘In Dictionary’).

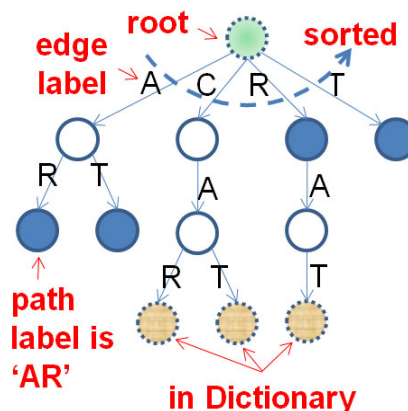


Figure 6.2: Suffix Trie

Suffix Trie is typically used as an efficient data structure for *dictionary*. Assuming that the Suffix Trie of a set of strings in the dictionary has been built, we can determine if a query/pattern string  $P$  exists in this dictionary (Suffix Trie) in  $O(m)$  where  $m$  is the length of string  $P$ —this is efficient<sup>10</sup>. We do this by traversing the Suffix Trie from the root. For example, if we want to find if the word  $P = \text{‘CAT’}$  exists in the Suffix Trie shown in Figure 6.2, we can start from the root node, follow the edge with label ‘C’, then ‘A’, then ‘T’. Since the vertex at this point has the word-terminating flag set to true, then we know that there is a word ‘CAT’ in the dictionary. Whereas, if we search for  $P = \text{‘CAD’}$ , we go through this path: root  $\rightarrow$  ‘C’  $\rightarrow$  ‘A’ but then we do not have edge with edge label ‘D’, so we conclude that ‘CAD’ is not in the dictionary.

---

**Exercise 6.6.1.1\*:** Implement this Suffix Trie data structure using the ideas outlined above, i.e. create a vertex object with up to 26 ordered edges that represent ‘A’ to ‘Z’ and suffix/word terminating flags. Insert each suffix of each string in  $S$  into the Suffix Trie one by one. Analyze the time complexity of such Suffix Trie construction strategy and compare with Suffix Array construction strategy in Section 6.6.4! Also perform  $O(m)$  queries for various pattern strings  $P$  by starting from the root and follow the corresponding edge labels.

---

<sup>9</sup>This is not a typo. The word ‘TRIE’ comes from the word ‘information reTRIEval’.

<sup>10</sup>Another data structure for dictionary is balanced BST—see Section 2.3. It has  $O(\log n \times m)$  performance for each dictionary query where  $n$  is the number of words in the dictionary. This is because one string comparison already costs  $O(m)$ .





### 6.6.3 Applications of Suffix Tree

Assuming that the Suffix Tree of a string  $T$  is *already built*, we can use it for these applications (not exhaustive):

#### String Matching in $O(m + occ)$

With Suffix Tree, we can find all (exact) occurrences of a pattern string  $P$  in  $T$  in  $O(m + occ)$  where  $m$  is the length of the pattern string  $P$  itself and  $occ$  is the total number of occurrences of  $P$  in  $T$ —*no matter how long* the string  $T$  is. When the Suffix Tree is *already built*, this approach is *much faster* than string matching algorithms discussed earlier in Section 6.4.

Given the Suffix Tree of  $T$ , our task is to search for the vertex  $x$  in the Suffix Tree whose path label represents the pattern string  $P$ . Remember, a matching is after all a *common prefix* between pattern string  $P$  and some suffixes of string  $T$ . This is done by just one root to leaf traversal of Suffix Tree of  $T$  following the edge labels. Vertex with path label equals to  $P$  is the desired vertex  $x$ . Then, the suffix indices stored in the terminating vertices (leaves) of the subtree rooted at  $x$  are the occurrences of  $P$  in  $T$ .

Example: In the Suffix Tree of  $T = \text{'GATAGACA\$'}$  shown in Figure 6.4 and  $P = \text{'A'}$ , we can simply traverse from root, go along the edge with edge label 'A' to find vertex  $x$  with the path label 'A'. There are 4 occurrences<sup>12</sup> of 'A' in the subtree rooted at  $x$ . They are suffix 7: 'A\$', suffix 5: 'ACA\$', suffix 3: 'AGACA\$', and suffix 1: 'ATAGACA\$'.

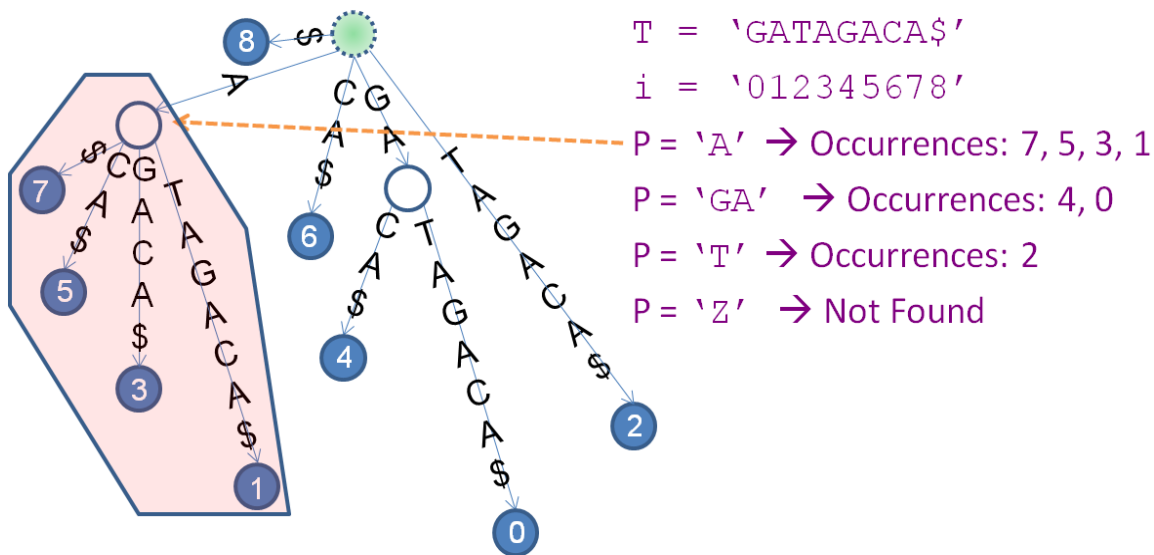


Figure 6.4: String Matching of  $T = \text{'GATAGACA\$'}$  with Various Pattern Strings

#### Finding the Longest Repeated Substring in $O(n)$

Given the Suffix Tree of  $T$ , we can also find the Longest Repeated Substring<sup>13</sup> (LRS) in  $T$  efficiently. The LRS problem is the problem of finding the longest substring of a string that occurs *at least twice*. The path label of the *deepest internal vertex*  $x$  in the Suffix Tree of  $T$  is the answer. Vertex  $x$  can be found with an  $O(n)$  tree traversal. The fact that  $x$  is

<sup>12</sup>To be precise,  $occ$  is the *size* of subtree rooted at  $x$ , which can be larger—but not more than double—than the actual number ( $occ$ ) of terminating vertices (leaves) in the subtree rooted at  $x$ .

<sup>13</sup>This problem has several interesting applications: Finding the chorus section of a song (that is repeated several times); Finding the (longest) repeated sentences in a (long) political speech, etc.

an internal vertex implies that it represents more than one suffixes of  $T$  (there will be  $> 1$  terminating vertices in the subtree rooted at  $x$ ) and these suffixes share a common prefix (which implies a repeated substring). The fact that  $x$  is the *deepest* internal vertex (from root) implies that its path label is the *longest* repeated substring.

Example: In the Suffix Tree of  $T = \text{'GATAGACA\$'}$  in Figure 6.5, the LRS is 'GA' as it is the path label of the deepest internal vertex  $x$ —'GA' is repeated twice in 'GATAGACA\\$'.

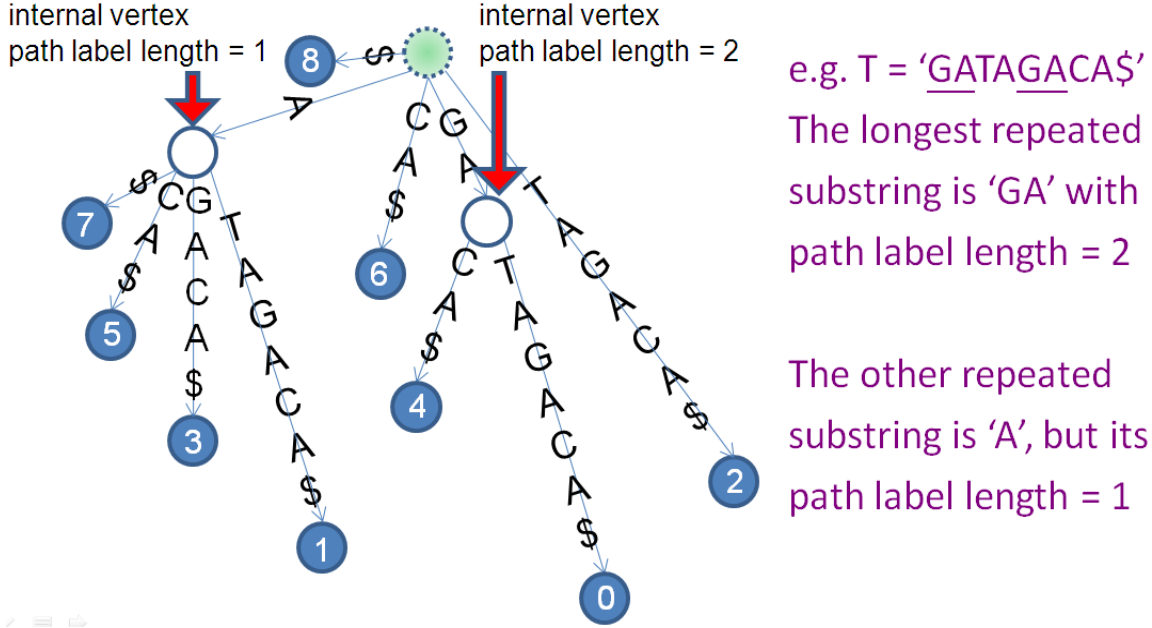


Figure 6.5: Longest Repeated Substring of  $T = \text{'GATAGACA\$'}$

### Finding the Longest Common Substring in $O(n)$

The problem of finding the Longest Common **Substring** (LCS<sup>14</sup>) of two **or more** strings can be solved in linear time<sup>15</sup> with Suffix Tree. Without loss of generality, let's consider the case with *two* strings only:  $T_1$  and  $T_2$ . We can build a **generalized Suffix Tree** that combines the Suffix Tree of  $T_1$  and  $T_2$ . To differentiate the source of each suffix, we use two different terminating vertex symbols, one for each string. Then, we mark *internal vertices* which have vertices in their subtrees with *different* terminating symbols. The suffixes represented by these marked internal vertices share a common prefix and come from *both*  $T_1$  and  $T_2$ . That is, these marked internal vertices represent the common substrings between  $T_1$  and  $T_2$ . As we are interested with the *longest* common substring, we report the path label of the *deepest* marked vertex as the answer.

For example, with  $T_1 = \text{'GATAGACA\$'}$  and  $T_2 = \text{'CATA\#'}$ , The Longest Common Substring is 'ATA' of length 3. In Figure 6.6, we see the vertices with path labels 'A', 'ATA', 'CA', and 'TA' have two different terminating symbols (notice that vertex with path label 'GA' is *not* considered as both suffix 'GACA\\$' and 'GATAGACA\\$' come from  $T_1$ ). These are the common substrings between  $T_1$  and  $T_2$ . The deepest marked vertex is 'ATA' and this is the longest common substring between  $T_1$  and  $T_2$ .

<sup>14</sup>Note that 'Substring' is different from 'Subsequence'. For example, "BCE" is a subsequence but not a substring of "ABCDEF" whereas "BCD" (contiguous) is both a subsequence and a substring of "ABCDEF".

<sup>15</sup>Only if we use the linear time Suffix Tree construction algorithm (not discussed in this book, see [65]).