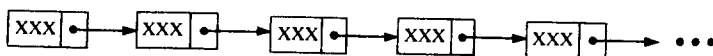


می‌شود. در اینجا ما فقط نگاهی اجمالی به چگونگی ظاهر شدن رشته‌ها در این ساختمان داده خواهیم داشت.

منظور از یک لیست پیوندی (یک طرفه)، یک دنباله از خانه‌های حافظه به نام گره است که به صورت خطی مرتب شده‌اند و هر گره شامل یک فیلد موسوم به پیوند یا اتصال است که به گره بعدی لیست اشاره می‌کند (یعنی دارای آدرس گره بعدی است).
شکل ۶-۳ نمودار فرضی یک لیست پیوندی را نشان می‌دهد.



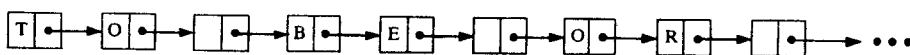
شکل ۶-۳. لیست پیوندی

رشته‌ها را می‌توان به صورت زیر در لیستهای پیوندی ذخیره کرد. در هر خانه حافظه یک کاراکتر یا تعداد معینی کاراکتر جایگزین می‌شود و یک پیوند که مقداری از حافظه را اشغال می‌کند آدرس خانه حافظه‌ای را که شامل کاراکتر بعدی یا تعدادی از کاراکترهای درون رشته است به دست می‌دهد. برای مثال، جمله مشهور زیر را در نظر بگیرید!

بودن یا نبودن، مسأله این است.

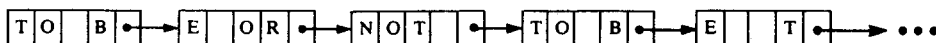
To be or not to be, that is the question.

شکل ۷-۳ (الف) چگونگی ذخیره شدن این رشته را در حافظه نشان می‌دهد که هر گره دارای یک کاراکتر است.



شکل ۷-۳ (الف). هر کاراکتر در یک گره

و شکل ۷-۳ (ب) چگونگی ذخیره شدن این حافظه را نشان می‌دهد که هر گره دارای چهار کاراکتر است.



شکل ۷-۳ (ب). چهار کاراکتر در یک گره

۴-۳ نوع داده کاراکتری

در این بخش مرور کلی بر روشهایی داریم که زبانهای برنامه‌نویسی مختلف داده نوع کاراکتری را مورد استفاده قرار می‌دهند. همانگونه که در فصل قبل (در بخش ۷-۲) متذکر شدیم، هر نوع داده‌ای دارای فرمول خاص خود جهت از کد در آوردن یک دنباله از بیتها، در حافظه است.

ثابت‌ها

اکثر زبانهای برنامه‌نویسی ثابتهای رشته‌ای را با قراردادن آن، در داخل یک یا دو کوتیشن نمایش می‌دهند. برای مثال:

'TO BE OR NOT TO BE' و 'THE END'

به ترتیب ثابت‌های رشته‌ای با طول ۷ و ۱۸ کاراکتر هستند. الگوریتم‌های این کتاب نیز ثابت‌های کاراکتری را به همین صورت تعریف می‌کند.

متغیرها

هر زبان برنامه‌نویسی دارای قانون و قاعده خاص خود برای ساختن متغیرهای کاراکتری است. با وجود این، متغیرهای کاراکتری به یکی از سه دسته زیر تقسیم می‌شوند: ایستا، نیمه‌ایستا و پویا. منظور از یک متغیر کاراکتری ایستا، متغیری است که طول آن قبل از اجرای برنامه تعریف شده و در سراسر برنامه ثابت است. منظور از یک متغیر کاراکتری نیمه‌ایستا، متغیری است که طول آن ممکن است در خلال اجرای برنامه تغییر کند طوری که طول آن نمی‌تواند از یک مقدار ماگزیممی که قبل از اجرای برنامه توسط برنامه معین می‌شود بزرگتر شود. منظور از یک متغیر کاراکتری پویا، متغیری است که طول آن می‌تواند در خلال اجرای برنامه تغییر کند. این سه دسته متغیر کاراکتری، به ترتیب متناظر با روشهایی هستند که رشته‌ها در حافظه کامپیوتر، به صورتی که در بخش قبل توضیح داده شد، ذخیره می‌شوند.

مثال ۲-۳

(الف) بسیاری از نسخه‌های زبان FORTRAN از متغیرهای کاراکتری ایستا CHARACTER استفاده می‌کنند. برای مثال قطعه برنامه FORTRAN زیر را در نظر بگیرید:

```
CHARACTER ST1*10, ST2*14
ST1 = 'THE END'
ST2 = 'TO BE OR NOT TO BE'
```

دستور اول اعلام می‌کند که ST1 و ST2 به ترتیب متغیرهای کاراکتری CHARACTER به طول 10 و 14 هستند. پس از اجرای هر یک از دو دستور جایگزینی بالا، ST1 و ST2 در حافظه به صورت زیر ذخیره می‌شوند:

```
ST1  T H E   E N D   ST2  T O B E   O R N O T   T
```

یعنی یک رشته در حافظه به صورت چیده شده از چپ ذخیره می‌شود. هرگاه طول رشته بزرگتر از طول خانه‌های حافظه باشد فضاهای خالی در سمت راست رشته جمع می‌شوند و اگر طول رشته کوچکتر از طول خانه‌های حافظه باشد، رشته از سمت راست بریده می‌شود.

(ب) زبان BASIC متغیرهای کاراکتری را به صورتی تعریف می‌کند که در انتهای نام متغیر یک علامت \$ قرار می‌گیرد. به بیان کلی‌تر، متغیرها، از نوع متغیرهای نیمه‌ایستا هستند که طول آنها از یک حد معینی نمی‌تواند بزرگتر باشد. برای مثال، قطعه برنامه BASIC زیر:

```
A$ = 'THE BEGINNING'  
B$ = 'THE END'
```

متغیرهای A\$ و B\$ را به صورت کاراکتری تعریف کرده است. هنگامی که این قطعه برنامه اجرا می‌شود، طول A\$ و B\$ به ترتیب 13 و 7 خواهد بود.

علاوه بر این در زبان BASIC برای نمایش ثابت‌های رشته‌ای، از دو کوتیشن استفاده می‌شود.

(ج) زبان SNOBOL از متغیرهای کاراکتری پویا استفاده می‌کند. برای مثال، قطعه برنامه SNOBOL

زیر:

```
WORD = 'COMPUTER'  
TEXT = 'IN THE BEGINNING'
```

متغیرهای WORD و TEXT را به صورت کاراکتری تعریف می‌کند. هنگامی که این قطعه برنامه اجرا می‌شود طول WORD و TEXT به ترتیب 18 و 16 خواهد بود. با وجود این، طول این متغیرها می‌تواند دیرتر در برنامه تغییر کند.

(د) زبان PL / I هم از متغیرهای کاراکتری ایستا و هم از متغیرهای کاراکتری نیمه‌ایستا استفاده

می‌کند. برای مثال، دستور PL / I زیر :

DECLARE NAME CHARACTER(20),
WORD CHARACTER(15) VARYING;

متغیر NAME را به صورت کاراکتری ایستا به طول 20 و متغیر WORD را به صورت کاراکتری نیمه‌ایستا تعریف می‌کند که طول آن می‌تواند تغییر کند ولی نمی‌تواند بیشتر از 15 باشد.
(۵) در زبان PASCAL، یک متغیر کاراکتری (که به اختصار به صورت CHAR نوشته می‌شود) می‌تواند تنها یک کاراکتر را نمایش دهد و از این رو است که یک رشته به صورت آرایه خطی از کاراکترها تعریف می‌شود. برای مثال :

VAR WORD: ARRAY[1..20] OF CHAR

WORD را به صورت یک رشته 20 کاراکتری تعریف می‌کند. علاوه بر این، WORD[1] اولین کاراکتر این رشته است و WORD[2] دومین کاراکتر رشته و غیره می‌باشد. به‌ویژه این که، آرایه‌های کاراکتری در پاسکال طول ثابت دارند و از این رو متغیرهای ایستا هستند.

۵-۳ عملیات بر روی رشته‌ها

هرچند یک رشته را می‌توان به صورت ساده، یک آرایه خطی از کاراکترها در نظر گرفت. با وجود این در کاربردها بین رشته‌ها و انواع دیگر آرایه‌ها یک تفاوت اساسی وجود دارد. به‌ویژه این که، گروههایی از عناصر متوالی، در یک رشته (نظیر کلمه‌ها، عبارتها و جمله‌ها) موسوم به زیررشته‌ها، می‌توانند واحدهایی مستقل روی این رشته‌ها باشند. علاوه بر این، واحدهای اصلی قابل دسترس در یک رشته معمولاً، زیررشته‌ها هستند نه تک تک کاراکترها .
برای مثال رشته زیر را در نظر بگیرید :

'TO BE OR NOT TO BE'

می‌توانیم این رشته را به صورت دنباله‌ای از 18 کاراکتر T, O, □, B, ..., E در نظر بگیریم. با وجود این، زیر رشته‌های TO, BE, OR, ... دارای معنی خاص خود هستند.

از سوی دیگر، یک آرایه خطی 18 عنصری با 18 عدد صحیح زیر را در نظر بگیرید :

4, 8, 6, 15, 9, 5, 4, 13, 8, 5, 11, 9, 9, 13, 7, 10, 6, 11

واحد اصلی قابل دسترس در چنین آرایه‌ای، معمولاً تک تک عناصر آن آرایه است. گروههایی متوالی

از عناصر این آرایه، معمولاً معنی خاصی ندارند. بنا به دلایل بالا، عملیات متعددی روی رشته‌ها تعریف شده است که معمولاً با انواع دیگر آرایه‌ها بکار برده نمی‌شود. در این بخش، این عملیات روی رشته‌ها توضیح داده می‌شود. بخش بعد، چگونگی استفاده از این عملیات را در پردازش کلمات یا رشته‌ها نشان می‌دهد. بجز در حالتی که به صورت صریح بیان می‌شود یا به طور ضمنی توضیح داده می‌شود، فرض می‌کنیم متغیرهای نوع کاراکتری پویا هستند و طول متغیر دارند و طول آن، در متنی که از متغیر استفاده می‌کند، تعیین می‌شود.

زیر رشته

دسترسی به یک زیررشته از یک رشته داده شده، مستلزم داشتن اطلاعات زیر است: (۱) نام رشته یا خود رشته (۲) مکان اولین کاراکتر زیررشته در رشته داده شده و (۳) طول زیر رشته یا مکان آخرین کاراکتر زیر رشته. ما این عمل را SUBSTRING می‌نامیم. به ویژه این که، می‌نویسیم:

SUBSTRING(string, initial, length)

که زیر رشته، یک رشته S (String) را نشان می‌دهد، اولین کاراکتر آن از مکان K (Initial) شروع شده و طول (Length) آن L است.

مثال ۳-۳

(الف) با استفاده از تابع بالا داریم:

SUBSTRING('TO BE OR NOT TO BE', 4, 7) = 'BE OR N'
SUBSTRING('THE END', 4, 4) = '□END'

(ب) تابع SUBSTRING(S, 4, 7) در برخی از زبانهای برنامه‌نویسی به صورت زیر نمایش داده می‌شود:

PL/1:	SUBSTR(S, 4, 7)
FORTRAN 77:	\$(4:10)
UCSD Pascal:	COPY(S, 4, 7)
BASIC:	MID\$(S, 4, 7)

شاخص گذاری

شاخص گذاری که تطبیق الگو نیز نامیده می‌شود، عبارت است از تعیین مکان الگوی رشته‌ای P که برای نخستین بار، در رشته داده شده T ظاهر شده است. ما این عمل را شاخص‌گذاری INDEX می‌نامیم

و می‌نویسیم:

$INDEX(\text{text}, \text{pattern})$

اگر الگوی P در متن T ظاهر نشود، آنگاه در $INDEX$ مقدار ۰ جایگزین می‌شود. آرگومان‌های text و pattern در عمل بالا، می‌تواند ثابت‌های رشته‌ای یا متغیرهای رشته‌ای باشند.

مثال ۳-۴

(الف) فرض کنید T شامل متن

'HIS FATHER IS THE PROFESSOR'

باشد آنگاه: $INDEX(T, 'THE')$ ، $INDEX(T, 'THEN')$ و $INDEX(T, '□THE□')$ به ترتیب مقادیر ۷، ۰ و ۱۴ دارند.

(ب) تابع $INDEX(\text{text}, \text{pattern})$ در برخی از زبانهای برنامه‌نویسی به صورت زیر نمایش داده می‌شود.

PL/1:	$INDEX(\text{text}, \text{pattern})$
UCSD Pascal:	$POS(\text{pattern}, \text{text})$

توجه دارید که در پاسکال UCSD آرگومانها به ترتیب عکس قرار گرفته‌اند.

اتصال دو رشته

فرض کنید S_1 و S_2 دو رشته باشند. از بخش ۲-۳ یادآور می‌شویم که اتصال S_1 و S_2 که ما آن را با $S_1 // S_2$ نمایش می‌دهیم، رشته‌ای است که شامل کاراکترهای S_1 و بدنبال آن کاراکترهای S_2 است.

مثال ۳-۵

(الف) فرض کنید $S_1 = 'MARK'$ و $S_2 = 'TWIN'$ باشند. آنگاه

$S_1 // S_2 = 'MARK TWAIN'$ اما $S_1 // S_2 = 'MARKTWIN'$

(ب) اتصال در برخی از زبانهای برنامه‌نویسی به صورت زیر نمایش داده می‌شود:

PL/1:	$S_1 // S_2$
FORTAN 77:	$S_1 // S_2$
BASIC:	$S_1 + S_2$
SNOBOL:	$S_1 S_2$

(بین S_1 و S_2 یک فضای خالی وجود دارد)

طول رشته

تعداد کاراکترهای داخل یک رشته، طول آن رشته نامیده می‌شود. برای تعیین طول یک رشته معلوم

می‌نویسیم:

LENGTH(string)

بدین ترتیب:

LENGTH('COMPUTER') = 8 LENGTH('□') = 0 LENGTH('') = 0

برخی از زبانهای برنامه‌نویسی این تابع را به صورت زیر نمایش می‌دهند:

PL/1:	LENGTH(string)
BASIC:	LEN(string)
UCSD Pascal:	LENGTH(string)
SNOBOL:	SIZE(string)

FORTRAN و **PASCAL** استاندارد که از متغیرهای رشته‌ای با طول ثابت استفاده می‌کنند و هیچ تابع کتابخانه‌ای **LENGTH** برای تعیین طول رشته‌ها ندارند. با وجود این، هرگاه از تمام فضای خالی انتهایی رشته صرف‌نظر کنیم، می‌توان چنین متغیرهایی را متغیرهایی با طول متغیر در نظر گرفت. بنابراین می‌توان

در این زبانها، یک زیربرنامه **LENGTH** نوشت طوری که: $LENGTH('MARC ') = 4$

در واقع، زبان **SNOBOL** دارای تابع کتابخانه‌ای برای رشته‌ها به نام **TRIM** است که فضاهای خالی

انتهای رشته را حذف می‌کند:

$TRIM('ERIK ') = 'ERIK'$

ما در الگوریتم‌های خود، هراز چندگاهی از این تابع **TRIM** استفاده خواهیم کرد.

۳-۶ پردازش کلمه یا رشته

در گذشته، داده‌های کاراکتری پردازش شده توسط کامپیوتر، اساساً از عناصر داده‌ای نظیر نام و آدرس تشکیل می‌شد. امروزه کامپیوتر، کارهای تاپی و نشریاتی از قبیل نامه‌ها، مقالات و گزارشها را پردازش می‌کند. بخاطر متنها و مطالب بعدی است که از اصطلاح "پردازش کلمه" یا رشته استفاده می‌کنیم.

متن تاپ شده‌ای داده شده است. معمولاً عملیات مرتبط با پردازش کلمه یا رشته به قرار زیر است:

(الف) جانشین سازی: یک رشته را جانشین رشته دیگری در داخل متن کنیم.

(ب) اضافه کردن: یک رشته را در وسط متن اضافه کنیم.

(ج) حذف کردن: یک رشته را از متن حذف کنیم.

عملیات بالا را می‌توان با استفاده از عملیات روی رشته‌ها، که در بخش قبل توضیح دادیم، اجرا نمود. در زیر، این عملیات را همراه با بررسی آنها به طور مستقل، شرح می‌دهیم. بسیاری از این عملیات، در هر یک از زبانهای برنامه‌نویسی شرح داده شده در بالا، یا به صورت کتابخانه‌ای وجود دارند یا می‌توان آنها را به‌سادگی تعریف کرد.

اضافه کردن

متن داده شده T را در نظر بگیرید، می‌خواهیم رشته S را طوری به آن اضافه کنیم که S از مکان K شروع شود. ما این عمل را به صورت

`INSERT(text, position, string)`

نمایش می‌دهیم. برای مثال :

`INSERT('ABCDEFGF', 3, 'XYZ') = 'ABXYZCDEFG'`
`INSERT('ABCDEFGF', 6, 'XYZ') = 'ABCDEXYZFG'`

این تابع `INSERT` را می‌توان با استفاده از عملیات روی رشته‌ها، که در بخش قبل تعریف شده است پیاده‌سازی کرد :

$INSERT(T, K, S) = SUBSTRING(T, 1, K - 1) // S // SUBSTRING(T, K, LENGTH(T) - K + 1)$
 یعنی زیررشته ابتدایی T که قبل از مکان K است و طول $K - 1$ دارد، به رشته S متصل شده است و نتیجه به بقیه قسمت T متصل می‌شود که از مکان K شروع می‌شود و طول $LENGTH(T) - (K - 1) = LENGTH(T) - K + 1$ دارد. ما به صورت ضمنی فرض کرده‌ایم که T یک متغیر پویا است و اندازه T خیلی بزرگ نمی‌شود.

حذف کردن

متن داده شده T را در نظر بگیرید، می‌خواهیم زیررشته‌ای را از آن حذف کنیم که از مکان K شروع می‌شود و طول L دارد. این عمل را به صورت زیر نمایش می‌دهیم :

`DELETE(text, position, length)`

برای مثال

`DELETE('ABCDEFGF', 4, 2) = 'ABCFGF'`
`DELETE('ABCDEFGF', 2, 4) = 'AFG'`

فرض می‌کنیم که عمل حذف زمانی که مکان $K = 0$ است انجام نمی‌شود. بنابراین اهمیت این "حالت صفر" دیرتر مشاهده می‌شود.

تابع DELETE را می‌توان با استفاده از عملیات بر روی رشته‌ها، که در بخش قبل به صورت زیر داده شده است پیاده‌سازی کرد:

$$\text{DELETE}(T, K, L) = \text{SUBSTRING}(T, 1, K - 1) // \text{SUBSTRING}(T, K + L, \text{LENGTH}(T) - K - L + 1)$$

یعنی زیررشته ابتدایی T که قبل از مکان K است به زیررشته انتهایی T که از مکان $K + L$ شروع می‌شود متصل شده است. طول زیررشته ابتدایی $K - 1$ است و طول زیررشته انتهایی برابر است با:

$$\text{LENGTH}(T) - (K + L - 1) = \text{LENGTH}(T) - K - L + 1$$

علاوه بر این فرض می‌کنیم به ازای $K = 0$ ، $\text{DELETE}(T, K, L) = T$. حال فرض کنید متن T و الگوی P داده شده است و بخواهیم از متن T اولین وقوع الگوی P را حذف کنیم. این عمل با استفاده از تابع DELETE بالا به صورت زیر انجام می‌شود:

$$\text{DELETE}(T, \text{INDEX}(T, P), \text{LENGTH}(P))$$

یعنی در متن T، نخست $\text{INDEX}(T, P)$ را محاسبه می‌کنیم یعنی مکانی را که الگوی P برای اولین بار در T ظاهر شده است پیدا می‌کنیم و بدنبال آن $\text{LENGTH}(P)$ یعنی تعداد کاراکترهای P را محاسبه می‌کنیم. یادآور می‌شویم که وقتی $\text{INDEX}(T, P) = 0$ (یعنی هنگامی که P در T ظاهر نشده باشد) متن T تغییر نمی‌کند.

مثال ۳-۶

(الف) فرض کنید $T = \text{'ABCDEFGG'}$ و $P = \text{'CD'}$ آنگاه $\text{INDEX}(T, P) = 3$ و $\text{LENGTH}(P) = 2$ از این رو

$$\text{DELETE}(\text{'ABCDEFGG'}, 3, 2) = \text{'ABEFG'}$$

(ب) فرض کنید $T = \text{'ABCDEFGG'}$ و $P = \text{'CD'}$ آنگاه $\text{INDEX}(T, P) = 0$ و $\text{LENGTH}(P) = 2$ از این رو، بنا به "حالت صفر":

$$\text{DELETE}(\text{'ABCDEFGG'}, 0, 2) = \text{'ABCDEFGG'}$$

که همان نتیجه موردانتظار است.

فرض کنید پس از خواندن متن T و الگوی P در کامپیوتر، بخواهیم هر وقوع الگوی P در متن T را حذف کنیم. این عمل را می‌توان با تکرار چند بار دستور

$$\text{DELETE}(T, \text{INDEX}(T, P), \text{LENGTH}(P))$$

انجام داد تا این که $\text{INDEX}(T, P) = 0$ (یعنی به محض این که P در T ظاهر نشده است). الگوریتمی که این عمل را انجام می‌دهد به صورت زیر است:

Algorithm 3.1: A text T and a pattern P are in memory. This algorithm deletes every occurrence of P in T .

1. [Find index of P .] Set $K := \text{INDEX}(T, P)$.
2. Repeat while $K \neq 0$:
 - (a) [Delete P from T .]
Set $T := \text{DELETE}(T, \text{INDEX}(T, P), \text{LENGTH}(P))$
 - (b) [Update index.] Set $K := \text{INDEX}(T, P)$.
- [End of loop.]
3. Write: T .
4. Exit.

تأکید می‌کنیم که پس از هر عمل حذف طول T کاهش می‌یابد و از این رو الگوریتم باید متوقف شود. با وجود این، تعداد دفعاتی که حلقه اجرا می‌شود می‌تواند تعداد دفعات ظاهر شدن P در متن اولیه T را افزایش دهد، که در مثال زیر نشان داده شده است.

مثال ۷-۳

(الف) فرض کنید الگوریتم 3.1 با داده‌های زیر اجرا شده است:

$$T = \text{XABYABZ}, \quad P = \text{AB}$$

آنگاه حلقه داخلی الگوریتم دوبار اجرا خواهد شد. در خلال اجرای اول الگوریتم، اولین وقوع AB از T حذف خواهد شد که در نتیجه آن T به صورت $T = \text{XYABZ}$ در خواهد آمد. در خلال اجرای دوم، وقوع باقیمانده AB از T حذف خواهد شد. از این رو $T = \text{XYZ}$. بنابراین خروجی الگوریتم است:

(ب) فرض کنید الگوریتم 3.1 با داده‌های زیر اجرا شده است:

$$T = \text{XAAABBBY}, \quad P = \text{AB}$$

ملاحظه می‌کنید که الگوی AB تنها یکبار در T ظاهر شده است اما حلقه الگوریتم سه بار اجرا می‌شود. به‌ویژه این‌که، پس از حذف AB در مرحله اول از T داریم. $T=XAABBY$ و از این رو AB مجدداً در T ظاهر شده است. پس از حذف AB در مرحله دوم از T ملاحظه می‌کنید که $T=XABY$ و AB همچنان در T وجود دارد. بالاخره پس از حذف AB در مرحله سوم از T داریم $T=XY$ و AB در T ظاهر نشده است و بدین ترتیب $INDEX(T,P) = 0$. از این رو XY خروجی الگوریتم است.

مثال بالا نشان می‌دهد که با تغییر متن T در اثر حذف الگو، الگو که در ابتدا در متن وجود داشت در انتهای این عمل در متن ظاهر نخواهد شد.

جانشین سازی

فرض کنید متن T داده شده است، می‌خواهیم الگوی P_2 را جانشین اولین وقوع الگوی P_1 کنیم. این عمل را به صورت زیر نمایش می‌دهیم:

$REPLACE(text, pattern_1, pattern_2)$

برای مثال:

$REPLACE('XABYABZ', 'AB', 'C') = 'XCZYABZ'$
 $REPLACE('XABYABZ', 'BA', 'C') = 'XABYABZ'$

در حالت دوم، الگو BA در متن ظاهر نشده است و از این رو هیچ جانشینی صورت نمی‌گیرد. توجه داشته باشید که تابع REPLACE را می‌توان به صورت ترکیب یک عمل حذف و یک عمل اضافه کردن بیان کرد، مشروط بر این که از همان تابع‌های قبلی DELETE و INSERT استفاده کنید. به‌ویژه این که، تابع REPLACE را می‌توان با استفاده از سه مرحله زیر اجرا کرد:

$K := INDEX(T, P_1)$
 $T := DELETE(T, K, LENGTH(P_1))$
 $INSERT(T, K, P_2)$

دو مرحله اول P_1 را از T حذف می‌کند و مرحله سوم P_2 را در مکان X ای که P_1 حذف شد به T اضافه می‌کند.

فرض کنید متن T و الگوهای P و Q در حافظه یک کامپیوتر قرار دارند و بخواهیم هر وقوع P در متن T توسط الگوی Q جانشین شود. این عمل را با استفاده مکرر از دستور

$REPLACE(T, P, Q)$

تا وقتی که $INDEX(T, P) = 0$ می‌توان انجام داد (یعنی تا وقتی که الگوی P در T وجود نداشته باشد).

الگوریتمی که این کار را انجام می‌دهد به شرح زیر است:

Algorithm 3.2: A text T and patterns P and Q are in memory. This algorithm replaces every occurrence of P in T by Q .

1. [Find index of P .] Set $K := \text{INDEX}(T, P)$.
2. Repeat while $K \neq 0$:
 - (a) [Replace P by Q .] Set $T := \text{REPLACE}(T, P, Q)$.
 - (b) [Update index.] Set $K := \text{INDEX}(T, P)$.
- [End of loop.]
3. Write: T .
4. Exit.

توجه کنید: هرچند این الگوریتم شباهت زیادی با الگوریتم 3.1 دارد اما پایان یافتن الگوریتم را تضمین نمی‌کند. این واقعیت در مثال ۸-۳ (ب) توضیح داده شده است. از طرف دیگر فرض کنید که طول Q کوچکتر از طول P باشد. در این صورت طول T پس از جانشانی کاهش می‌یابد. این مطلب تضمین می‌کند که در این حالت خاص Q کوچکتر از P است و باید الگوریتم پایان یابد.

مثال ۸-۳

(الف) فرض کنید الگوریتم 3.2 با داده‌های زیر اجرا شده است:

$$T = XABYABZ, \quad P = AB, \quad Q = C$$

آنگاه حلقه داخلی الگوریتم دوبار اجرا خواهد شد. در خلال اجرای اول الگوریتم، C جانشین اولین وقوع AB در T می‌شود و نتیجه می‌دهد $T = XCYABZ$. در خلال اجرای دوم، C جانشین وقوع باقیمانده AB

در T می‌شود و نتیجه می‌دهد $T = XCYCZ$. از این رو خروجی الگوریتم است.

(ب) فرض کنید الگوریتم 3.2 با داده‌های زیر اجرا شده است:

$$T = XAY, \quad P = A, \quad Q = AB$$

در آن صورت این الگوریتم هیچوقت به پایان نمی‌رسد. به این دلیل که صرفنظر از تعداد دفعاتی که حلقه اجرا می‌شود P همواره در متن T ظاهر خواهد شد. به‌ویژه این که:

$$T = XABY \text{ در پایان اجرای اول حلقه است.}$$

$$T = XAB^2Y \text{ در پایان اجرای دوم حلقه است.}$$

.....

$$T = XAB^nY \text{ در پایان اجرای } n \text{ ام حلقه است.}$$

چون P یک زیررشته Q است از این رو در اینجا حلقه بی پایان Infinite Loop یا LOOP خواهیم داشت.

۷-۳ الگوریتم‌های تطبیق الگو

تطبیق الگو، مسأله‌ای است که تعیین می‌کند یک الگوی رشته‌ای داده شده P در متن رشته‌ای T وجود دارد یا خیر. فرض می‌کنیم که طول P کوچکتر از طول T است. این بخش دو الگوریتم در مورد تطبیق الگوها را مورد بحث و بررسی کامل قرار می‌دهد. علاوه بر آن، پیچیدگی این الگوریتم‌ها مورد بررسی قرار می‌گیرد تا بتوان کارایی آنها را مورد مقایسه قرار داد.

توجه کنید: در خلال بررسی الگوریتم‌های تطبیق الگو، کاراکترهای گاهی اوقات با حروف کوچک a, b, c, ... نمایش داده می‌شوند و از توانها برای نمایش تعداد دفعات تکرار آنها استفاده می‌شود، مانند:

$$a^2b^3ab^2 \quad \text{به جای } aabbabb \quad \text{و} \quad (cd)^3 \quad \text{به جای } cdc dcd$$

علاوه بر این، رشته تھی با حرف یونانی لاندا Λ نمایش داده می‌شود و اتصال رشته‌های X و Y به صورت X.Y یا فقط XY نمایش داده می‌شود.

الگوریتم اول تطبیق الگو

الگوریتم اول تطبیق الگو، الگوریتم روشنی است که در آن الگوی داده شده P با هر یک از زیررشته‌های T مقایسه می‌شود. عمل مقایسه با حرکت از چپ به راست متن T انجام می‌شود تا به یک تطبیق با P برسیم. به طور مشروح این که، فرض کنید:

$$W_k = \text{SUBSTRING}(T, K, \text{LENGTH}(P))$$

باشد. به عبارت دیگر فرض کنید W_k زیررشته‌ای از T با همان طول P و با شروع از K امین کاراکتر T باشد. نخست P را کاراکتر به کاراکتر با اولین زیررشته، W_1 مقایسه می‌کنیم. اگر تمام کاراکترها مساوی باشند در آن صورت $P = W_1$ و در نتیجه P در T ظاهر شده است و $\text{INDEX}(T, P) = 1$. از سوی دیگر، فرض کنید به این نتیجه رسیدیم که یک کاراکتر P همان کاراکتر متناظر در W_1 نیست. در آن صورت $P \neq W_1$ و بلافاصله می‌توانیم به زیررشته بعدی W_2 برویم به عبارت دیگر، در مرحله بعد P را با W_2 مقایسه می‌کنیم. اگر $P = W_2$ ، در آن صورت P را با W_3 مقایسه می‌کنیم و الی آخر. عمل مقایسه متوقف می‌شود:

(الف) هرگاه به یک تطبیق P با زیررشته W_k برسیم و از این رو P در متن T وجود دارد و $\text{INDEX}(T, P) = K$ یا (ب) هرگاه به هیچ تطبیق P با زیررشته W_k نرسیم و از این رو P در متن T وجود نخواهد داشت.

MAX مقدار ماکزیمم اندیس K برابر است با $LENGTH(T) - LENGTH(P) + 1$. حال به عنوان مثال فرض کنید که P یک رشته 4 کاراکتری و T یک رشته 20 کاراکتری باشد و فرض کنید P و T در حافظه به صورت آرایه‌های خطی نمایش داده شده‌اند که در هر عنصر آرایه یک کاراکتر ذخیره شده است. یعنی

$$T = T[1]T[2]T[3]...T[19]T[20] \quad \text{و} \quad P = P[1]P[2]P[3]P[4]$$

آنگاه P با هر یک از زیررشته‌های 4 کاراکتری بعدی T مقایسه می‌شود:

$$W_1 = T[1]T[2]T[3]T[4], \quad W_2 = T[2]T[3]T[4]T[5], \quad \dots, \quad W_{17} = T[17]T[18]T[19]T[20]$$

توجه داشته باشید که تعداد $MAX = 20 - 4 + 1 = 17$ زیررشته 4 کاراکتری در T وجود دارد. نمایش رسمی این الگوریتم که در آن P یک رشته R کاراکتری و T یک رشته S کاراکتری است در الگوریتم 3.3 نشان داده شده است.

ملاحظه می‌کنید الگوریتم 3.3 شامل دو حلقه است، یکی از حلقه‌ها داخل حلقه دیگر قرار دارد. حلقه خارجی برای هر زیررشته R کاراکتری متوالی

$$W_k = T[K]T[K+1] \dots T[K+R-1]$$

رشته T اجرا می‌شود. حلقه داخلی P را با W_k ، کاراکتر به کاراکتر مقایسه می‌کند. اگر هیچ کاراکتری تطبیق نکنند، در آن صورت کنترل به مرحله S داده می‌شود که K را افزایش می‌دهد و پس از آن منتهی به زیررشته بعدی T می‌شود. اگر تمام R کاراکتر P با کاراکترهای W_k تطابق داشته باشد، آنگاه P در T وجود دارد و K اندیس INDEX الگوی P در T است. از سوی دیگر، اگر حلقه خارجی به پایان تمام مراحل خود برسد ولی در T هیچ P ظاهر نشود در آن صورت $INDEX = 0$.

Algorithm 3.3: (Pattern Matching) P and T are strings with lengths R and S, respectively, and are stored as arrays with one character per element. This algorithm finds the INDEX of P in T.

1. [Initialize.] Set $K := 1$ and $MAX := S - R + 1$.
2. Repeat Steps 3 to 5 while $K \leq MAX$:
3. Repeat for $L = 1$ to R : [Tests each character of P.]
If $P[L] \neq T[K + L - 1]$, then: Go to Step 5.
[End of inner loop.]
4. [Success.] Set $INDEX = K$, and Exit.
5. Set $K := K + 1$.
[End of Step 2 outer loop.]
6. [Failure.] Set $INDEX = 0$.
7. Exit.

پیچیدگی این الگوریتم تطبیق متن به وسیلهٔ C تعداد مقایسه بین کاراکترهای الگوی P و کاراکترهای متن T اندازه گرفته می‌شود. برای پیدا کردن C ، فرض می‌کنیم N_k نمایش تعداد مقایسه‌هایی باشد که هنگام مقایسهٔ P با W_k در حلقهٔ داخلی اتفاق می‌افتد. در آن صورت

$$C = N_1 + N_2 + \dots + N_L$$

که در آن L ، مکان L در متن T است که در آن P برای اولین بار در T ظاهر می‌شود یا $L = \text{MAX}$ است اگر P در متن T ظاهر نشده باشد.

مثال بعدی C را برای یک P مشخص و T محاسبه می‌کند که در آن $\text{LENGTH}(P) = 4$ و

$$\text{LENGTH}(T) = 20 \text{ و در نتیجه } \text{MAX} = 20 - 4 + 1 = 17$$

مثال ۹-۳

(الف) فرض کنید $T = \text{cdcd...cd} = (\text{cd})^{10}$. واضح است که P در T ظاهر نشده است. علاوه بر این، برای هر 17 بار اجرای حلقه، $N_k = 1$ ، چون اولین کاراکتر P با W_k مطابقت نمی‌کند. از این رو

$$C = 1 + 1 + 1 + \dots + 1 = 17$$

(ب) فرض کنید $P = \text{aaba}$ و $T = \text{ababaaba...}$ ملاحظه می‌کنید که P یک زیررشتهٔ T است. در واقع $P = W_5$ و در نتیجه $N_5 = 4$. علاوه بر این از مقایسهٔ P با $W_1 = \text{abab}$ نتیجه می‌گیریم که $N_1 = 2$ ، چون اولین حرف این دو رشته با هم مطابقت دارد اما از مقایسهٔ P با $W_2 = \text{baba}$ ، ملاحظه می‌کنید که $N_2 = 1$ چون اولین حرفهای آنها با هم مطابقت ندارد. به‌طور مشابه، $N_3 = 2$ و $N_4 = 1$. بنابراین

$$C = 2 + 1 + 2 + 1 + 4 = 10$$

(ج) فرض کنید $P = \text{aaab}$ و $T = \text{aa...a} = \text{a}^{20}$. در اینجا P در T ظاهر نشده است. همچنین هر $W_k = \text{aaaa}$ ، از این رو هر $N_k = 4$ ، چون سه حرف اول P با سه حرف اول T مطابقت می‌کند. بنابراین

$$C = 4 + 4 + \dots + 4 = 17 \cdot 4 = 68$$

در حالت کلی، هنگامی که P یک رشتهٔ r کاراکتری و T یک رشتهٔ s کاراکتر است، اندازهٔ داده‌های این الگوریتم برابر است با:

$$n = r + s$$

بدترین حالت وقتی اتفاق می‌افتد که همانند مثال ۹-۳ (ج) تمام کاراکترهای P بجز آخرین کاراکتر با هر زیررشتهٔ W_k مطابقت داشته باشد. در این حالت، $C(n) = r(s - r + 1)$ برای مقدار ثابت n داریم

$$s = n - r \text{ طوری که}$$

$$C(n) = r(n - 2r + 1)$$

مقدار ماگزیمم $C(n)$ وقتی اتفاق می‌افتد که $r = (n + 1)/4$ (مسأله ۱۹-۳ را ببینید) باشد. بنابراین با قراردادن این مقدار به جای r در فرمول مربوط به $C(n)$ نتیجه می‌گیریم:

$$C(n) = \frac{(n+1)^2}{8} = O(n^2)$$

پیچیدگی حالت میانگین در هر وضعیت واقعی، بستگی به حالت‌های احتمالی دارد که معمولاً ناشناخته هستند. هرگاه کاراکترهای P و T به تصادف از روی یکی از حروف الفبا انتخاب شوند، تجزیه و تحلیل پیچیدگی حالت میانگین همچنان مشکل است، اما پیچیدگی حالت میانگین همچنان به صورت ضربی از بدترین حالت است. بنابراین، پیچیدگی این الگوریتم را چنین بیان می‌کنیم: پیچیدگی الگوریتم تطبیق الگو برابر $O(n^2)$ است. به بیان دیگر، زمان موردنیاز برای اجرای این الگوریتم متناسب با n^2 است. (این نتیجه را با نتیجه صفحه ۸۹ مقایسه کنید.)

الگوریتم دوم تطبیق الگو

الگوریتم دوم تطبیق الگو، از جدولی استفاده می‌کند که از یک الگوی خاص P مشتق شده است اما مستقل از متن T است. برای روشن شدن مطلب فرض کنید:

$$P = aaba$$

نخست دلیلی برای درایه‌های جدول اقامه می‌کنیم و چگونگی استفاده از آنها را شرح می‌دهیم. فرض کنید $T = T_1 T_2 T_3 \dots$ که در آن T_i کاراکتر i ام T را نمایش می‌دهد و فرض کنید کاراکتر اول T با دو کاراکتر اول P مطابقت می‌کند یعنی فرض کنید $T = aa..$. آنگاه T دارای یکی از سه صورت زیر است:

$$(i) T = aab\dots, \quad (ii) T = aaa\dots, \quad (iii) T = aax$$

که در آن x می‌تواند هر کاراکتر دلخواهی به غیر از a یا b باشد. فرض کنید T_3 را خواندیم و دریافتیم که $T_3 = b$. در آن صورت، بدنبال آن T_4 را می‌خوانیم تا ببینیم آیا $T_4 = a$ است یا خیر. P با W_1 تطابق خواهد داشت. از طرف دیگر، فرض کنید $T_3 = a$. در آن صورت، می‌دانیم که $P \neq W_1$ ، اما همچنین می‌دانیم که $W_2 = aa\dots$ یعنی دو کاراکتر اول زیررشته W_2 با دو کاراکتر اول P مطابقت می‌کند. از این رو بعد از آن T_4 را می‌خوانیم تا ببینیم آیا $T_4 = b$ است. بالاخره، فرض کنید $T_3 = x$. در آن صورت می‌دانیم $P \neq W_1$ ، اما همچنین می‌دانیم که $P \neq W_2$ و $P \neq W_3$ ، چون x در P ظاهر نمی‌شود. از این رو بعد از آن، T_4 را می‌خوانیم تا ببینیم آیا $T_4 = a$ است یا خیر. یعنی ببینیم که آیا کاراکتر اول W_4 با کاراکتر اول P مطابقت می‌کند یا خیر.

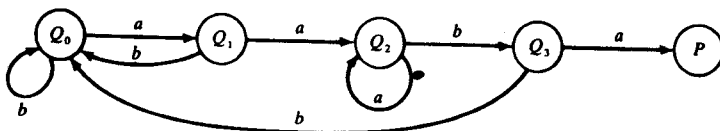
دو نکته قابل توجه و مهم در روش بالا وجود دارد. نخست، هنگام خواندن T_3 ، تنها لازم است T_3 با

آن تعداد از کاراکترها، که در P وجود دارد مقایسه شود. اگر هیچ یک از این کاراکترها مطابقت نداشت، در آن صورت ما در حالت آخری هستیم که کاراکتری مانند x در P وجود ندارد. دوم، پس از خواندن و بررسی T_3 ، T_4 را می‌خوانیم. لازم نیست مجدداً به متن T برگردیم.

شکل ۳-۸ (الف) شامل جدولی است که برای الگوی $P = aaba$ از الگوریتم دوم تطبیق الگو استفاده می‌کند. در هر دو مورد جدول و گراف همراه آن، یعنی الگوی P و زیررشته‌های Q با حروف کج لاتین نمایش داده شده است.

	a	b	x
Q_0	Q_1	Q_0	Q_0
Q_1	Q_2	Q_0	Q_0
Q_2	Q_2	Q_3	Q_0
Q_3	P	Q_0	Q_0

(الف) جدول تطبیق الگو



(ب) گراف تطبیق الگو

شکل ۳-۸

طرز تشکیل جدول به صورت زیر است. قبل از هر کاری، فرض می‌کنیم Q_i نمایش زیررشته اولیه P با طول i باشد. از این رو

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = a^2, \quad Q_3 = a^2b, \quad Q_4 = a^2ba = P$$

در اینجا $Q_i = \Lambda$ رشته تهی است. سطرهای این جدول را با زیررشته‌های اولیه P به جزء خود P ، برچسب‌گذاری می‌کنیم. ستونهای این جدول با a ، b ، و x برچسب‌گذاری می‌شوند که در آن x هر کاراکتر دلخواهی می‌تواند باشد که در الگوی P ظاهر نشده است. فرض کنید f تابعی باشد که این جدول معین می‌کند، یعنی فرض کنید تابع

$$f(Q_i, t)$$

درایه سطر Q_i و ستون t را نمایش دهد که در آن t می‌تواند هر کاراکتر دلخواهی باشد. درایه $f(Q_i, t)$ بنابه تعریف بزرگترین Q ای است که در رشته $Q_i t$ ، یعنی در نتیجه اتصال Q_i و t به صورت زیررشته‌نهایی ظاهر می‌شود. برای مثال :

$$a^2 \text{ بزرگترین } Q \text{ ای است که یک زیررشته پایانی } Q_2 a = a^3 \text{ است از این رو } Q_2 a = a^3$$

$$a \text{ بزرگترین } Q \text{ ای است که یک زیررشته پایانی } Q_1 b = ab \text{ است از این رو } Q_1 b = ab$$

$$a \text{ بزرگترین } Q \text{ ای است که یک زیررشته پایانی } Q_0 a = a \text{ است از این رو } Q_0 a = a$$

$$a \text{ بزرگترین } Q \text{ ای است که یک زیررشته پایانی } Q_3 x = a^2 b x \text{ است از این رو } Q_3 x = a^2 b x$$

و الی آخر، اگرچه $Q_1 = a$ زیررشته پایانی $Q_2 a = a^3$ است، داریم $f(Q_2, a) = Q_2$ زیرا Q_2 نیز یک زیررشته پایانی $Q_2 a = a^3$ است و Q_2 بزرگتر از Q_1 است. توجه دارید که به ازای هر Q ، $f(Q_i, x) = Q_0$ چون x در الگوی P ظاهر نشده است. بنابراین، ستون متناظر با x معمولاً از جدول حذف می‌شود.

جدول را می‌توان با استفاده از گراف جهت‌دار برچسب‌گذاری شده شکل ۸-۳ (ب) رسم کرد. گراف به صورت زیر به دست می‌آید. نخست، متناظر با هر زیررشته اولیه Q_i از P یک گره در گراف وجود دارد. Q ها حالت‌های گراف نامیده می‌شوند و Q حالت اولیه نام دارد. دوم، متناظر با هر درایه در جدول، یک پیکان (یک یال جهت‌دار) در گراف وجود دارد. به ویژه این که، اگر

$$f(Q_i, t) = Q_j$$

آنگاه یک پیکان با برچسب t از Q_i به Q_j وجود دارد. برای مثال، $f(Q_2, b) = Q_3$ از این رو یک پیکان با برچسب b از Q_2 به Q_3 وجود دارد. جهت سهولت در نمادگذاری، تمام پیکان‌های با برچسب x را که اجباراً منتهی به حالت اولیه Q_0 می‌شوند حذف کرده‌ایم. اکنون آماده‌ایم که برای الگوی $P = aaba$ الگوریتم دوم تطبیق الگو را ارائه دهیم. توجه دارید که در بحث زیر، برای نام تمام متغیرهای یک حرفی داخل الگوریتم، از حروف بزرگ استفاده می‌کنیم. فرض کنید $T = T_1 T_2 T_3 \dots T_N$ نمایش متن رشته n کاراکتری باشد که در آن برای یافتن الگوی P عمل جستجو انجام می‌شود. با شروع از حالت اولیه Q_0 و استفاده از متن T به یک دنباله از حالت‌های S_1, S_2, S_3, \dots به صورت زیر دست می‌یابیم. فرض کنید $S_1 = Q_0$ و اولین کاراکتر T_1 را می‌خوانیم. چه از طریق جدول عمل کنیم و چه از گراف شکل ۸-۳ استفاده کنیم، زوج (S_1, T_1) حالت دوم S_2 را نتیجه می‌دهد یعنی $F(S_1, T_1) = S_2$ ، کاراکتر بعدی T_2 را می‌خوانیم، زوج (S_2, T_2) حالت S_3 را نتیجه می‌دهد و الی آخر. دو حالت ممکن وجود دارد :

(۱) حالت $S_k = P$ ، الگوی موردنظر باشد. در این حالت، در P وجود دارد و اندیس آن

$K - \text{LENGTH}(P)$ است.

(۲) هیچ حالتی از S_1, S_2, \dots, S_{N+1} برابر P نیست. در این حالت، P در T وجود ندارد. الگوریتم را با استفاده از الگوی $P = aaba$ و دو متن مختلف توضیح می‌دهیم.

مثال ۱۰-۳

(الف) فرض کنید $T = aabcaba$ با شروع از Q_0 و استفاده از کاراکترهای T و گراف (یا جدول) شکل ۳-۸، دنباله حالت‌های زیر به دست می‌آید:

$$Q_0 \xrightarrow{ca} Q_1 \xrightarrow{ca} Q_2 \xrightarrow{cb} Q_3 \xrightarrow{cc} Q_0 \xrightarrow{ca} Q_1 \xrightarrow{cb} Q_0 \xrightarrow{ca} Q_1$$

به حالت P دست نمی‌یابیم، از این رو P در T وجود ندارد.

(ب) فرض کنید $T = abcaabaca$. آنگاه به دنباله حالت‌های زیر دست می‌یابیم:

$$Q_0 \xrightarrow{ca} Q_1 \xrightarrow{cb} Q_0 \xrightarrow{cc} Q_0 \xrightarrow{ca} Q_1 \xrightarrow{ca} Q_2 \xrightarrow{cb} Q_3 \xrightarrow{ca} P$$

در اینجا، در حالت S_8 به الگوی P دست می‌یابیم. از این رو P در T ظاهر شده است و اندیس آن $\text{LENGTH}(P) = 8$ است.

بیان رسمی الگوریتم دوم تطبیق الگو به صورت زیر است:

Algorithm 3.4: (Pattern Matching). The pattern matching table $F(Q_i, T)$ of a pattern P is in memory, and the input is an N -character string $T = T_1T_2 \dots T_N$. This algorithm finds the INDEX of P in T .

1. [Initialize.] Set $K := 1$ and $S_1 = Q_0$.
2. Repeat Steps 3 to 5 while $S_K \neq P$ and $K \leq N$.
3. Read T_K .
4. Set $S_{K+1} := F(S_K, T_K)$. [Finds next state.]
5. Set $K := K + 1$. [Updates counter.]
- [End of Step 2 loop.]
6. [Successful?]
 - If $S_K = P$, then:
 - INDEX = $K - \text{LENGTH}(P)$.
 - Else:
 - INDEX = 0.
 - [End of If structure.]
7. Exit.

زمان اجرای الگوریتم بالا با تعداد دفعات اجرای حلقه مرحله ۲ متناسب است. بدترین حالت وقتی

اتفاق می‌افتد که تمام متن T خوانده شود یعنی هنگامی که حلقه $n = \text{LENGTH}(T)$ بار اجرا شود. بنابراین پیچیدگی الگوریتم بالا را می‌توان به صورت زیر بیان کرد:

پیچیدگی این الگوریتم تطبیق متن برابر $O(n)$ است.

توجه کنید: یک مسأله ترکیباتی را با زمان چندجمله‌ای، حل پذیر گویند اگر به ازای بعضی از مقادیر m یک جواب الگوریتمی با پیچیدگی ای برابر $O(n^m)$ ، وجود داشته باشد و با زمان خطی، حل پذیر گویند اگر یک جواب الگوریتمی با پیچیدگی ای برابر $O(n)$ وجود داشته باشد که در آن n تعداد داده‌ها است. بدین ترتیب حالت دوم دو الگوریتم تطبیق الگو که در این بخش تشریح شده است با زمان خطی حل پذیر هستند. الگوریتم اول تطبیق الگو با زمان چندجمله‌ای حل پذیر است.

مسأله‌های حل شده

اصطلاحات، ذخیره رشته‌ها

مسأله ۱-۳: فرض کنید W رشته $ABCD$ است. (الف) طول W را پیدا کنید. (ب) تمام زیررشته‌های W را بنویسید. (ج) تمام زیررشته‌های اولیه W را بنویسید.

حل: (الف) تعداد کاراکترهای رشته W طول آن است، بنابراین طول W برابر ۴ است.

(ب) هر زیر دنباله‌ای از کاراکترهای W یک زیررشته W است. در W ، ۱۱ زیررشته وجود دارد:

زیررشته‌ها: $ABCD$, ABC , BCD , AB , BC , CD , A , B , C , D , Λ

طول آنها: 4 3 2 1 0

در اینجا Λ معرف رشته تهی است.

(ج) زیررشته‌های اولیه عبارتند از $ABCD$ ، ABC ، AB ، A ، یعنی هم رشته تهی و هم زیررشته‌های اولیه با A شروع می‌شوند.

مسأله ۲-۳: فرض کنید یک زبان برنامه‌نویسی حداقل از ۴۸ کاراکتر شامل ۲۶ حرف، ۱۰ رقم و حداقل ۱۲ کاراکتر مخصوص استفاده می‌کند. حداقل تعداد بیتها و تعداد بیتهای متداول برای نمایش یک کاراکتر در حافظه کامپیوتر را به دست آورید.

حل: از آنجا که $2^6 < 48 < 2^7$ ، حداقل یک کد ۶ بیتی برای نمایش این ۴۸ کاراکتر مورد نیاز است. معمولاً یک کامپیوتر از یک کد ۷ بیتی نظیر کد $ASCII$ یا یک کد ۸ بیتی نظیر کد $EBCDIC$ برای نمایش کاراکترها استفاده می‌کند. این کدگذاری‌ها به کامپیوتر اجازه می‌دهند کاراکترهای مخصوص بسیار زیادتری را در

حافظه ایجاد کرده به نمایش و پردازش آنها پردازند.

مسئله ۳-۳: به اختصار سه نوع ساختار مختلف را که برای ذخیره رشته‌ها بکار می‌رود توضیح دهید.
 حل: (الف) ساختارهای حافظه با طول ثابت. در این حالت رشته‌هایی که در خانه‌های حافظه ذخیره می‌شود همگی دارای طول مساوی هستند، معمولاً طول حافظه برای رشته‌ها 80 کاراکتر است.
 (ب) حافظه با طول متغیر که ماگزیم طول آن ثابت است. در این حالت، رشته‌ها نیز در خانه‌های حافظه ذخیره می‌شوند همگی دارای طول مساوی هستند، بنابراین باید طول واقعی رشته داخل حافظه معلوم باشد.

(ج) حافظه به صورت لیست پیوندی. در اینجا هر خانه حافظه به دو قسمت تقسیم می‌شود، در قسمت اول یک کاراکتر (یا تعداد ثابت کوچکی از کاراکتر) ذخیره می‌شود و قسمت دوم شامل آدرس حافظه کاراکتر بعدی است.

مسئله ۳-۴: رشته ذخیره شده در شکل ۹-۳ را تعیین کنید. فرض می‌شود مقدار پیوند 0 علامت پایان لیست پیوندی است.

	CHAR	LINK
1	OY F	10
2	ING	7
3		
4	A TH	2
5		
6	ER.	0
7	OF B	11
8	A J	1
9		
10	OREV	6
11	EAUT	12
12	Y IS	8

START → 4

شکل ۹-۳

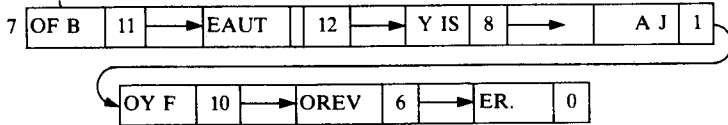
حل: در اینجا رشته در یک ساختمان لیست پیوندی ذخیره می‌شود که هر گره آن شامل 4 کاراکتر است. مقدار START، مکان گره اول لیست را به دست می‌دهد:

4 A TH 2

مقدار پیوند در این گره، مکان گره بعدی لیست را به دست می‌دهد:

2 ING 7

اگر این روش را ادامه بدهیم، دنباله گره‌های زیر به دست می‌آید:



بنابراین رشته ذخیره شده عبارت است از:

A THING OF BEAUTY IS A JOY FOREVER.

مسئله ۵-۳: برخی از (الف) مزایا (ب) معایب استفاده از حافظه پیوندی برای ذخیره رشته‌ها را بیان کنید.

حل: (الف) هنگام استفاده از حافظه پیوندی به راحتی می‌توان عملیات اضافه کردن، حذف کردن، اتصال و جابجایی زیررشته‌ها را انجام داد.

(ب) برای ذخیره پیوندها به حافظه اضافی نیاز است. علاوه بر این به کاراکتر وسط لیست مستقیماً نمی‌توان دسترسی پیدا کرد.

مسئله ۶-۳: به اختصار درباره معنی (الف) متغیرهای کاراکتری ایستا (ب) نیمه ایستا (ج) پویا توضیح دهید.

حل: (الف) طول این نوع متغیرها قبل از اجرای برنامه تعریف می‌شود و نمی‌توان طی اجرای برنامه طول آنها را تغییر داد.

(ب) طول این نوع متغیرها می‌تواند طی اجرای برنامه تغییر کند، اما طول متغیر نمی‌تواند از یک مقدار ماگزیمی که قبل از اجرای برنامه تعریف می‌شود بیشتر شود.

(ج) طول این نوع متغیرها می‌تواند طی اجرا برنامه، تغییر کند.

مسئله ۷-۳: فرض کنید MEMBER یک متغیر کاراکتری با طول ثابت 20 است. فرض کنید یک رشته به صورت چیده شده از چپ در خانه حافظه ذخیره می‌شود که فضاها خالی، طرف راست حافظه را پر

مسئله ۱۰-۳: یادآور می‌شویم که برای نمایش مکانی که یک الگو برای نخستین بار در متن T ظاهر می‌شود، از INDEX(TEXT, PATTERN) استفاده می‌کنیم. در این تابع مقدار 0 جایگزین می‌شود اگر الگو در متن ظاهر نشده باشد. مطلوب است تعیین

(الف) INDEX(S, 'JO'), (ب) INDEX(S, 'JOY'), (ج) INDEX(S, 'JO'),
 (د) INDEX(T, 'A'), (ه) INDEX(T, 'A□'), (و) INDEX(T, 'THE')
 حل: (الف) INDEX(S, 'JO') = 1, (ب) INDEX(S, 'JOY') = 0, (ج) INDEX(S, 'JO') = 10, INDEX(T, 'THE') = 0

(د) INDEX(T, 'A') = 1, (ه) INDEX(T, 'A□') = 21 و (و) INDEX(T, 'THE') = 0
 یادآوری می‌کنیم که □ برای نمایش یک فضای خالی مورد استفاده قرار می‌گیرد.

مسئله ۱۱-۳: یادآور می‌شویم که برای نمایش اتصال دو رشته S₁ و S₂ از S₁ // S₂ استفاده می‌کنیم.
 (الف) این تابع در زبانهای (i) PL/I (ii) FORTRAN (iii) BASIC (iv) SNOBOL (v) UCSD PASCAL چگونه نمایش داده می‌شود؟

(ب) مطلوب است تعیین (i) 'THE' // 'END' و (ii) 'THE' // 'END' □.
 (ج) مطلوب است تعیین (i) SUBSTRING(S, 1, 9) // □, ' و SUBSTRING(S, 11, 5) و
 (ii) 'GIVEN' // SUBSTRING(T, 28, 3)

حل: (الف) (i) S₁ // S₂, (ii) S₁ // S₂, (iii) S₁ + S₂, (iv) S₁ S₂ (بین S₁ و S₂ یک فضای خالی وجود دارد) (v) CONCAT(S₁, S₂)

(ب) منظور از S₁ // S₂ رشته‌ای است که از کاراکترهای S₁ و بدنبال آن کاراکترهای S₂ تشکیل شده باشد.
 از این رو (i) THEEND (ii) THE END
 (ج) (i) JONES, JOHN PAUL (ii) FORGIVEN

مسئله ۱۲-۳: یادآور می‌شویم که برای نمایش اضافه کردن یک رشته S در یک متن داده شده T با شروع از مکان K از INSERT(text, position, string) استفاده می‌کنیم.

(الف) مطلوب است تعیین (i) INSERT('AAAAA', 1, 'BBB'),
 (ii) INSERT('AAAAA', 3, 'BBB') و (iii) INSERT('AAAAA', 6, 'BBB')

(ب) فرض کنید T متن 'THE STUDENT IS ILL' باشد. با استفاده از INSERT، T را طوری تغییر

دهید که (i) The student is very ill (ii) The student is ill today
 و (iii) The student is very ill today را بخواند.

حل: (الف) (i) BBBAAAAA, (ii) AABBBAAA و (iii) AAAAABBB

(ب) دقت کنید در جاهای لازم فضاهای خالی بگنجانید. (i) INSERT(T, 15, '□ VERY')

(ii) INSERT(T, 19, '□ TODAY') (iii) INSERT(INSERT(T, 19, '□ TODAY'), 15, '□ VERY')

. INSERT(INSERT(T, 15, '□ VERY'), 24, '□ TODAY')

مسأله ۱۳-۳: مطلوب است تعیین

(الف) DELETE('JOHN PAUL JONES', 6, 5) و DELETE('AAABBB', 2, 2)

(ب) REPLACE('AAABBB', 'AA', 'BB')

و REPLACE('JOHN PAUL JONES', 'PAUL', 'DAVID')

حل: (الف) DELETE(T, K, L) از متن T زیررشته‌ای را حذف می‌کند که از مکان K شروع شده طول L

دارد. از این رو جوابها عبارتند از ABBB و JOHN JONES .

(ب) REPLACE(T, P₁, P₂) در متن T اولین وقوع الگوی P₁ توسط الگوی P₂ جانشین می‌شود. از

این رو جوابها عبارتند از BBABBB و JOHN DAVID JONES .

پردازش کلمه

در مسأله‌های ۱۴-۳ تا ۱۷-۳، S یک داستان کوتاه است که در یک آرایه خطی LINE با n عنصر

به گونه‌ای ذخیره شده است که در آن هر LINE[K] یک متغیر کاراکتری ایستا با قدرت ذخیره‌سازی 80

کاراکتر است و یک خط از داستان را نمایش می‌دهد. همچنین LINE[1] یعنی اولین خط داستان تنها

شامل عنوان داستان است و LINE[N] آخرین خط داستان تنها شامل نام نویسنده است. علاوه بر این، هر

پاراگراف با 5 فضای خالی شروع می‌شود و در هیچ کجای داستان بجز احتمالاً در عنوان LINE[1] یا نام

نویسنده LINE[N] هیچ تورفتگی مانند ابتدای پاراگراف وجود ندارد.

مسأله ۱۴-۳: یک زیربرنامه Procedure بنویسید تا NUM تعداد پاراگرافهای داستان کوتاه S را بشمارد.

حل: از خط LINE[2] شروع می‌کنیم و با خط LINE [N - 1] شمارش تعداد خطوط را که با 5 فضای

خالی شروع می‌شود به پایان می‌بریم. زیربرنامه به صورت زیر است:

Procedure P3.14: PAR(LINE, N, NUM)

1. Set NUM := 0 and BLANK := '□□□□□'.
2. [Initialize counter.] Set K := 2.
3. Repeat Steps 4 and 5 while K ≤ N - 1.
4. [Compare first 5 characters of each line with BLANK.]
If SUBSTRING(LINE[K], 1, 5) = BLANK, then:
Set NUM := NUM + 1.
[End of If structure.] \
5. Set K := K + 1. [Increments counter.]
[End of Step 3 loop.]
6. Return.

مسئله ۱۵-۳: یک زیربرنامه `procedure` بنویسید تا `NUM` تعداد دفعاتی را که کلمه `the` در این داستان کوتاه `S` ظاهر می‌شود بشمارد. `the` در کلمه `mother` را در شمارش شرکت ندهید و فرض می‌شود که هیچ جمله‌ای با کلمه `the` پایان نمی‌یابد.

حل: توجه دارید که کلمه `the` می‌تواند به صورت `THE` در آغاز یک خط، به صورت `THE` در پایان یک خط یا به صورت `THE` در هر جای دیگری از یک خط می‌تواند ظاهر شود. از این‌رو برای هر خط باید این سه حالت را مورد توجه قرار دهیم. زیربرنامه به صورت زیر است:

Procedure P3.15: COUNT(LINE, N, NUM)

1. Set WORD := 'THE' and NUM := 0.
2. [Prepare for the three cases.]
Set BEG := WORD // '□', END := '□' // WORD and
MID := '□' // WORD // '□'.
3. Repeat Steps 4 through 6 for K = 1 to N:
4. [First case.] If SUBSTRING(LINE[K], 1, 4) = BEG, then:
Set NUM := NUM + 1.
5. [Second case.] If SUBSTRING(LINE[K], 77, 4) = END, then:
Set NUM := NUM + 1.
6. [General case.] Repeat for J = 2 to 76.
If SUBSTRING(LINE[K], J, 5) = MID, then:
Set NUM := NUM + 1.
[End of If structure.]
[End of Step 6 loop.]
7. Return.

مسئله ۱۶-۳: چنانچه بخواهیم تعداد دفعات وقوع کلمه دلخواهی مانند `W` را با طول `r` بشماریم توضیح دهید چه تغییراتی باید در `Procedure P3.15` بدهیم.

حل: لازم است سه نوع تغییر اصلی در این زیربرنامه داده شود.

(الف) واضح است که در مرحله ۱ باید `THE` به `W` تغییر داده شود.

(ب) از آنجا که طول `W` برابر `r` است و نه `3` از این‌رو باید تغییرات مناسب در مراحل `3` تا `6` داده شود.

(ج) علاوه بر این باید این امکان در نظر گرفته شود که بعد از `W` می‌تواند علامت نقطه‌گذاری نظیر

`W, W; W. W?`

باشد. از این‌رو بیشتر از سه حالت باید مورد بررسی قرار گیرد.

مسئله ۱۷-۳: الگوریتمی را شرح دهید که پاراگراف‌های `K` ام و `L` ام داستان کوتاه `S` را جابجا می‌کند.

حل: الگوریتم را به دو زیربرنامه `Procedure` تجزیه می‌کنیم.

زیربرنامه `A`، `Procedure A`: مقادیر آرایه‌های `BEG` و `END` را پیدا می‌کند که در آن

`LINE[END[K]]` و `LINE[BEG[K]]`

به ترتیب شامل اولین و آخرین خط پاراگراف K ی داستان S است.
 زیربرنامه B ، Procedure B : با استفاده از مقادیر END[K] و BEG[K] و مقادیر END[L] و BEG[L] بلاک خطوط پاراگراف K با بلاک خطوط پاراگراف L جابجا می‌شود.

تطبيق الگو

مسأله ۱۸-۳: برای هر یک از الگوهای P و متنهای T زیر، C تعداد مقایسه‌هایی را پیدا کنید که با آن INDEX اندیس P در T با استفاده از الگوریتم "آرام"، الگوریتم ۳-۳ به دست می‌آید.

$$P = aaa, T = (aabb)^3 = aabbaabbaabb \quad (\text{ج}) \quad P = abc, T = (ab)^5 = ababababab \quad (\text{الف})$$

$$P = aaa, T = abaabbaabbbbaaaabbbb \quad (\text{د}) \quad P = abc, T = (ab)^{2n} \quad (\text{ب})$$

حل: خاطر نشان می‌کنیم که $C = N_1 + N_2 + \dots + \dots + N_L$ که در آن N_k تعداد مقایسه‌هایی را نشان می‌دهد که هنگام مقایسه P با W_k در حلقه داخلی انجام می‌شود.

(الف) نخست توجه داشته باشید که تعداد

$$\text{LENGTH}(T) - \text{LENGTH}(P) + 1 = 10 - 3 + 1 = 8$$

زیررشته W_k وجود دارد. داریم:

$$C = 2 + 1 + 2 + 1 + 2 + 1 + 2 + 1 = 4(3) = 12$$

و $\text{INDEX}(T, P) = 0$ ، چون P در T ظاهر نشده است.

(ب) تعداد $2n - 3 + 1 = 2(n - 1) + 1$ زیرکلمه W_k وجود دارد. داریم:

$$C = 2 + 1 + 2 + 1 + \dots + 2 + 1 = (n + 1)(3) = 3n + 3$$

و $\text{INDEX}(T, P) = 0$

(ج) تعداد $12 - 3 + 1 = 10$ زیرکلمه W_k وجود دارد. داریم:

$$C = 3 + 2 + 1 + 1 + 3 + 2 + 1 + 1 + 3 + 2 = 19$$

و $\text{INDEX}(T, P) = 0$

(د) داریم:

$$C = 2 + 1 + 3 + 2 + 1 + 1 + 3 = 13$$

و $\text{INDEX}(T, P) = 7$

مسأله ۱۹-۳: فرض کنید P یک رشته r کاراکتری و T یک رشته s کاراکتری باشد و فرض کنید هنگامی که الگوریتم ۳-۳ بر P و T اعمال می‌شود C(n) نمایش تعداد مقایسه‌های صورت گرفته باشد. در اینجا

$$. n = r + s$$

(الف) $C(n)$ پیچیدگی بهترین حالت را پیدا کنید.

(ب) ثابت کنید مقدار ماگزیمم $C(n)$ وقتی اتفاق می‌افتد که $r = (n+1) / 4$.

حل: (الف) بهترین حالت وقتی اتفاق می‌افتد که P یک زیررشته اولیه T باشد یا به عبارت دیگر وقتی

که $INDEX(T, P) = 1$ در این حالت $C(n) = r$. (فرض می‌کنیم $r \leq s$).

(ب) بنا به بحث بخش ۷-۳،

$$C = C(n) = r(n - 2r + 1) = nr - 2r^2 + r$$

در اینجا n ثابت است از این رو $C = C(n)$ را می‌توان به صورت تابعی از r در نظر گرفت. حساب

دیفرانسیل و انتگرال می‌گوید ماگزیمم مقدار C وقتی اتفاق می‌افتد که $C' = dC/dr = 0$ در اینجا C'

مشتق C نسبت به r است. با استفاده از حساب دیفرانسیل و انتگرال به دست می‌آید:

$$C' = n - 4r + 1$$

با قراردادن $C' = 0$ و حل آن نسبت به r نتیجه موردنظر به دست می‌آید.

مسئله ۲۰-۳: الگوی $P = aaabb$ را در نظر بگیرید. جدول و گراف جهت‌دار برجسب‌گذاری شده متناظر

با آن را که در الگوریتم "سریع" یا الگوریتم دوم تطبیق الگو مورد استفاده قرار گرفت، رسم کنید.

حل: نخست لیست قطعه‌های اولیه P را می‌نویسیم:

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = a^2, \quad Q_3 = a^3, \quad Q_4 = a^3b, \quad Q_5 = a^3b^2$$

برای هر کاراکتر t ، درایه $f(Q_i, t)$ جدول بزرگترین Q ای است که به صورت یک زیررشته انتهایی در

رشته $Q_i t$ ظاهر می‌شود. محاسبه می‌کنیم:

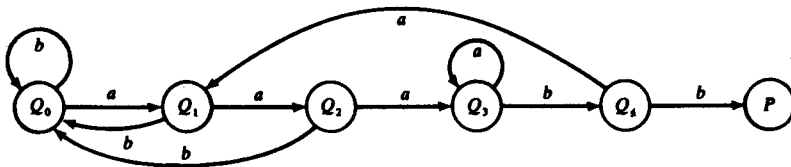
$$\begin{array}{lllll} f(\Lambda, a) = a, & f(a, a) = a^2, & f(a^2, a) = a^3, & f(a^3, a) = a^3, & f(a^3b, a) = a \\ f(\Lambda, b) = \Lambda, & f(a, b) = \Lambda, & f(a^2, b) = \Lambda, & f(a^3, b) = a^3b, & f(a^3b, b) = P \end{array}$$

از این رو جدول موردنظر در شکل ۱۰-۳ (الف) ارائه شده است.

	a	b
Q_0	Q_1	Q_0
Q_1	Q_2	Q_0
Q_2	Q_3	Q_0
Q_3	Q_3	Q_4
Q_4	Q_4	P

شکل ۱۰-۳ (الف)

گراف متناظر با آن در شکل ۳-۱۰ (ب) رسم شده است که در آن برای هر Q_i یک گره وجود دارد و پیکان از Q_i و Q_j برای هر درایه $f(Q_i, t) = Q_j$ داخل جدول با کاراکتر t برچسب‌گذاری شده است.



شکل ۳-۱۰ (ب)

مسئله ۳-۲۱: جدول و گراف متناظر با آن را برای الگوریتم دوم تطبیق الگو که در آن $P = ababab$ الگو است رسم کنید.

حل: زیررشته‌های اولیه P عبارتند از:

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = ab, \quad Q_3 = aba, \quad Q_4 = abab, \quad Q_5 = ababa, \quad Q_6 = ababab = P$$

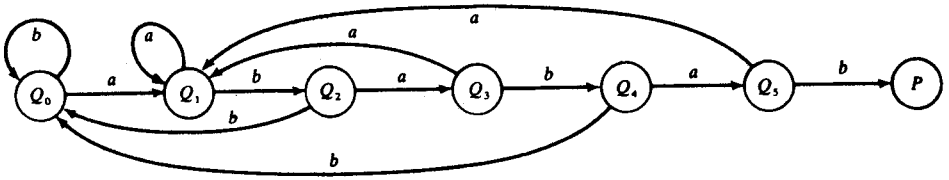
تابع f درایه‌های جدول را به صورت زیر به دست می‌دهد:

$f(\Lambda, a) = a$	$f(\Lambda, b) = \Lambda$
$f(a, a) = a$	$f(a, b) = ab$
$f(ab, a) = aba$	$f(ab, b) = \Lambda$
$f(aba, a) = a$	$f(aba, b) = abab$
$f(abab, a) = ababa$	$f(abab, b) = \Lambda$
$f(ababa, a) = a$	$f(ababa, b) = P$

جدول در شکل ۳-۱۱ (الف) و گراف متناظر با آن در شکل ۳-۱۱ (ب) آمده است.

	a	b
Q_0	Q_1	Q_0
Q_1	Q_1	Q_2
Q_2	Q_3	Q_0
Q_3	Q_1	Q_4
Q_4	Q_5	Q_0
Q_5	Q_1	P

(الف)



(ب)

شکل ۳-۱۱

مسئله‌های تکمیلی

رشته‌ها

مسئله ۳-۲۲: رشته ذخیره شده در شکل ۳-۱۲ را به دست آورید.

START	CHAR	LINK
5	1 UNIT	11
	2 HE P	8
	3	
	4 S	0
	5 WE T	2
	6 THE	1
	7	
	8 EQPL	12
	9 TATE	4
	10	
	11 ED S	9
	12 E OF	6

شکل ۳-۱۲

مسئله ۳-۲۳: رشته $W = 'XYZST'$ را در نظر بگیرید. (الف) تمام زیررشته‌های W و

(ب) تمام زیررشته‌های ابتدایی W را بنویسید.

مسئله ۲۴-۳: فرض کنید W یک رشته به طول n است. مطلوب است تعیین تعداد (الف) زیررشته‌های W و (ب) زیررشته‌های ابتدایی W .

مسئله ۲۵-۳: فرض کنید STATE یک متغیر کاراکتری با طول ثابت ۱۲ است. محتوای STATE را پس از جایگزینی‌های زیر توضیح دهید. (الف) STATE := 'NEW YORK' و (ج) STATE := 'PENNSYLVANIA'

(ب) STATE := 'SOUTH CAROLINA' و (ج) STATE := 'PENNSYLVANIA'

عملیات بر روی رشته‌ها

در مسئله‌های ۲۶-۳ تا ۳۱-۳ فرض کنید S و T متغیرهای کاراکتری هستند به طوری که

$S = 'WE THE PEOPLE'$ و $T = 'OF THE UNITED STATES'$

مسئله ۲۶-۳: طول S و T را تعیین کنید.

مسئله ۲۷-۳: مطلوب است تعیین (الف) SUBSTRING(S , 4, 8) و (ب) SUBSTRING(T , 10, 5)

مسئله ۲۸-۳: مطلوب است تعیین (الف) INDEX(S , 'P') و (ب) INDEX(S , 'E')

(ج) INDEX(S , 'THE') (د) INDEX(T , 'THE') (ه) INDEX(T , 'THEN')

(و) INDEX(T , 'TE')

مسئله ۲۹-۳: با استفاده از $S_1 // S_2$ به عنوان اتصال دو رشته S_1 و S_2 ، مطلوب است تعیین

(الف) 'NO' // 'EXIT' (ب) 'NO' // 'EXIT' // 'NO' // 'EXIT' و

(ج) SUBSTRING(S , 4, 10) // 'ARE' // SUBSTRING(T , 8, 6)

مسئله ۳۰-۳: مطلوب است تعیین (الف) DELETE('AAABBB', 3, 3)

(ب) DELETE('AAABBB', 1, 4) (ج) DELETE(S , 1, 3) و (د) DELETE(T , 1, 7)

مسئله ۳۱-۳: مطلوب است تعیین (الف) REPLACE('ABABAB', 'B', 'BAB')

(ب) REPLACE(S , 'WE', 'ALL') و (ج) REPLACE(T , 'THE', 'THESE')

مسئله ۳۲-۳: مطلوب است تعیین (الف) INSERT('AAA', 2, 'BBB')

(ب) INSERT('ABCDE', 3, 'XYZ') و (ج) INSERT('THE BOY', 5, 'BIG')

مسئله ۳۳-۳: فرض کنید U متن 'MARC STUDIES MATHEMATICS' باشد. با استفاده از INSERT،

متن U را طوری تغییر دهید که (الف) MARC STUDIES ONLY MATHEMATICS

(ب) MARC STUDIES MATHEMATICS AND PHYSICS (ج) MARC STUDIES APPLIED

MATHEMATICS را بخواند.

تطبیق الگو

مسئله ۳-۳۴: الگوی $P = abc$ را در نظر بگیرید. با استفاده از الگوریتم تطبیق الگو "آرام" یعنی الگوریتم ۳-۳، تعداد مقایسه‌هایی را پیدا کنید که با آنها INDEX اندیس P در هر یک از متنهای T زیر به دست می‌آید:

(الف) a^{10} ، (ب) $(aba)^{10}$ (ج) $(cbab)^{10}$ (د) d^{10} و (ه) d^n که در آن $n > 3$ است.

مسئله ۳-۳۵: الگوی $P = a^5b$ را در نظر بگیرید. مسئله ۳-۳۴ را با هر یک از متنهای T زیر مجدداً حل کنید:

(الف) a^{20} (ب) a^n که در آن $n > 6$ (ج) d^{20} و (د) d^n که در آن $n > 6$ است.

مسئله ۳-۳۶: الگوی $P = a^3ba$ را در نظر بگیرید. جدول و گراف جهت‌دار برجسب‌گذاری شده متناظر با آن را که در الگوریتم تطبیق الگو "سریع" مورد استفاده قرار گرفت رسم کنید.

مسئله ۳-۳۷: مسئله ۳-۳۶ را برای الگوی $P = aba^2b$ مجدداً حل کنید.

برای مسأله‌های زیر برنامه بنویسید

در مسأله‌های ۳-۳۸ تا ۳-۴۰، فرض می‌شود که پیشگفتار این متن در یک آرایه خطی LINE ذخیره شده است به طوری که $LINE[K]$ یک متغیر کارا کتری ایستا با قدرت ذخیره 80 کارا کتر است و یک خط از پیشگفتار را نمایش می‌دهد. فرض کنید هر پاراگراف با 5 فضای خالی شروع می‌شود و در هیچ جای دیگر پیشگفتار تورفتگی وجود ندارد. علاوه بر این فرض می‌شود یک متغیر NUM وجود دارد که تعداد خطوط پیشگفتار را به دست می‌دهد.

مسئله ۳-۳۸: یک برنامه بنویسید که یک آرایه خطی PAR را معین می‌کند به گونه‌ای که $PAR[K]$ حاوی مکان K ام پاراگراف است و علاوه بر این یک متغیر NPAR را معین می‌کند که حاوی تعداد پاراگرافها است.

مسئله ۳-۳۹: یک برنامه بنویسید که یک کلمه WORD داده شده را بخواند و C تعداد دفعاتی را که WORD در LINE ظاهر شده است را بشمارد. برنامه را با استفاده از (الف) $WORD = 'THE'$ و (ب) $WORD = 'HENCE'$ آزمایش کنید.

مسئله ۳-۴۰: یک برنامه بنویسید که پاراگرافهای J ام و K ام را جابجا کند. برنامه را با استفاده از $J = 2$ و $K = 2$ آزمایش کنید.

در مسأله‌های ۳-۴۱ تا ۳-۴۶ فرض می‌شود که پیشگفتار این متن در یک متغیر کارا کتری TEXT ذخیره شده است. فرض کنید 5 فضای خالی مبین یک پاراگراف جدید است.

مسئله ۳-۴۱: یک برنامه بنویسید که یک آرایه خطی PAR را به گونه‌ای ایجاد کند که PAR[K] حاوی مکان پاراگراف K ام در TEXT باشد همچنین مقدار یک متغیر NPAR را پیدا کند که حاوی تعداد پاراگرافها است. این مسئله را با مسئله ۳-۳۸ مقایسه کنید.

مسئله ۳-۴۲: یک برنامه بنویسید که یک کلمه WORD داده شده را بخواند و آنگاه C تعداد دفعاتی را که WORD در TEXT ظاهر شده است را بشمارد. برنامه را با استفاده از (الف) 'THE' = WORD و (ب) 'HENCE' = WORD آزمایش کنید. این مسئله را با مسئله ۳-۳۹ مقایسه کنید.

مسئله ۳-۴۳: یک برنامه بنویسید که پاراگرافهای J ام و K ام متن TEXT را جابجا کند. این برنامه را با استفاده از $J = 2$ و $K = 4$ آزمایش کنید. این مسئله را با مسئله ۳-۴۰ مقایسه کنید.

مسئله ۳-۴۴: یک برنامه بنویسید که کلمه‌های WORD1 و WORD2 را بخوانید و آنگاه هر وقوع WORD1 در TEXT توسط WORD2 جایگزین گردد. برنامه را با استفاده از 'HENCE' = WORD1 و 'THUS' = WORD2 آزمایش کنید.

مسئله ۳-۴۵: یک زیربرنامه INST(TEXT, NEW, K) بنویسید که رشته NEW را در داخل متن TEXT در آغاز [K] TEXT اضافه کند.

مسئله ۳-۴۶: یک زیربرنامه PRINT(TEXT, K) بنویسید که رشته کاراکتری TEXT را در خطوطی با حداکثر K کاراکتر چاپ کند. هیچ کلمه‌ای به دو قسمت تقسیم نشده و در دو خط ظاهر نمی‌شود، از اینرو برخی از خطوط می‌توانند شامل چند فضای خالی در انتهایشان باشند. هر پاراگراف با خط مربوطه‌اش شروع شده و به کمک 5 فضای خالی تورفتگی دارد. برنامه را با استفاده از (الف) $K = 80$ (ب) $K = 70$ و (ج) $K = 60$ آزمایش کنید.

فصل ۴

آرایه‌ها، رکوردها و اشاره‌گرها

۱-۴ مقدمه

ساختمان داده‌ها به دو دسته خطی و غیرخطی تقسیم می‌شود. یک ساختمان داده را خطی گویند، هرگاه عناصر آن تشکیل یک دنباله دهند، به بیان دیگر یک لیست خطی باشد. برای نمایش ساختمان داده خطی در حافظه، دو روش اساسی وجود دارد. یکی از این روشها، عبارت است از داشتن رابطه خطی بین عناصری که به وسیله خانه‌های متوالی حافظه نمایش داده می‌شود. این ساختارهای خطی آرایه‌ها نام دارند که موضوع اصلی این فصل را تشکیل می‌دهد. روش دیگر عبارت است از داشتن رابطه خطی بین عناصری که به وسیله اشاره‌گرها یا پیوندها نمایش داده می‌شود. این ساختارهای خطی لیستهای پیوندی نام دارند که موضوع اصلی مطالب فصل ۵ را تشکیل می‌دهد. ساختارهای غیرخطی نظیر درختها و گرافها در فصل‌های بعد مورد مطالعه قرار می‌گیرد.

عملیاتی که معمولاً بر روی یک ساختار خطی انجام می‌شود خواه این ساختار آرایه باشد یا یک لیست پیوندی، شامل عملیات زیر است:

(الف) پیمایش. پردازش هر عنصر داخل لیست را پیمایش گویند.

(ب) جستجو کردن. پیدا کردن مکان یک عنصر با یک مقدار داده‌شده یا رکورد با یک کلید معین را

جستجو کردن گویند.

(ج) اضافه کردن. افزودن یک عنصر جدید به لیست را اضافه کردن گویند.

(د) حذف کردن. حذف یک عنصر از لیست را حذف کردن گویند.

(ه) مرتب کردن. تجدید آرایش عناصر با یک نظم خاص را مرتب کردن گویند.

(و) ادغام کردن. ترکیب دو لیست در یک لیست را ادغام کردن گویند.

ساختار خطی خاصی که برای یک وضعیت معین انتخاب می‌شود بستگی به تعداد دفعاتی دارد که عملیات مختلف بالا روی ساختار اجرا می‌شود.

این فصل یک ساختار خطی کاملاً متداولی را مورد بررسی قرار می‌دهد که آرایه نام دارد. از آنجا که پیمایش، جستجو و مرتب کردن آرایه معمولاً ساده است از آنها اغلب برای ذخیره مجموعه‌هایی از داده‌های نسبتاً دائمی استفاده می‌شود. از سوی دیگر، اگر اندازه ساختار و داده‌های آن پیوسته در حال تغییر باشد آنگاه آرایه نمی‌تواند به خوبی ساختاری نظیر لیست پیوندی باشد که در فصل ۵ بررسی می‌شود.

۲-۴ آرایه‌های خطی

یک آرایه خطی لیستی از n یا تعداد متناهی عنصر داده‌ای همجنس است (یعنی عناصر داده‌ای از یک نوع هستند) بطوری که:

(الف) به عناصر آرایه به ترتیب به کمک یک مجموعه از اندیسها که شامل n عدد متوالی است رجوع می‌شود.

(ب) عناصر آرایه به ترتیب در خانه‌های متوالی حافظه ذخیره می‌شوند.

n یا تعداد عناصر آرایه طول یا اندازه آرایه نام دارد. اگر به صراحت بیان نشود فرض می‌کنیم که مجموعه اندیسها شامل اعداد صحیح 1، 2، ...، n است. در حالت کلی، طول یا تعداد عناصر داده‌ای آرایه را می‌توان از مجموعه اندیسها به کمک فرمول زیر به دست آورد:

$$(۴-۱) \quad \text{طول} = UB - LB + 1$$

که در آن UB بزرگترین اندیس یا کران بالای آرایه و LB کوچکترین اندیس یا کران پائین آرایه است. توجه دارید که وقتی $LB = 1$ باشد طول آرایه برابر UB است.

عناصر یک آرایه A را می‌توان با نماد اندیس‌گذاری

$$A_1, A_2, A_3, \dots, A_n$$

یا با نماد پراتزگذاری (که در BASIC، PL / I، FORTRAN و بکار می‌رود)

$$A(1), A(2), \dots, A(N)$$

یا با نماد کروش‌های (که در PASCAL بکار می‌رود)

$A[1], A[2], A[3], \dots, A[N]$

نمایش داد. ما معمولاً از نماد اندیس‌گذاری یا نماد کروش‌های استفاده می‌کنیم. عدد K در $A[K]$ زیرنویس یا اندیس و $A[K]$ یک متغیر اندیس‌دار نام دارد. توجه دارید که اندیس‌ها اجازه می‌دهند به مکان نسبی هر عنصر A دسترسی پیدا کنیم.

مثال ۴-۱

(الف) فرض کنید DATA یک آرایه خطی 6 عنصری از اعداد صحیح باشد به طوری که

$DATA[1]=247$ $DATA[2]=56$ $DATA[3]=429$ $DATA[4]=135$ $DATA[5]=87$ $DATA[6]=156$

گاهی اوقات ما این آرایه را فقط با نوشتن

DATA: 247, 56, 429, 135, 87, 156

نمایش می‌دهیم. آرایه DATA غالباً به صورت شکل ۴-۱ (الف) یا شکل ۴-۱ (ب) نمایش داده می‌شود.

DATA					
247	56	429	135	87	156
1	2	3	4	5	6

(ب)

DATA	
1	247
2	56
3	429
4	135
5	87
6	156

(الف)

شکل ۴-۱

(ب) یک شرکت فروشنده اتومبیل از یک آرایه AUTO برای ثبت تعداد اتومبیل‌های فروخته شده هر سال از 1932 تا 1984 استفاده می‌کند. بجای اینکه مجموعه اندیسها از 1 شروع شود، بهتر است از 1932 شروع شود، طوری که :

$\text{AUTO}[K] =$ تعداد اتومبیل‌های فروخته‌شده در سال K

آنگاه $\text{LB} = 1932$ کران پائین و $\text{UB} = 1984$ کران بالای آرایه AUTO است.

با توجه به فرمول (۱-۴)

$$\text{طول} = \text{UB} - \text{LB} + 1 = 1984 - 1930 + 1 = 55$$

یعنی AUTO شامل 55 عنصر است و مجموعه اندیسی آن شامل تمام اعداد صحیح از 1932 تا 1984 است.

هر زبان برنامه‌نویسی قاعده و روش خاصی برای معرفی آرایه‌ها دارد. هر یک از این روشها، به صورت صریح یا ضمنی، سه نوع اطلاعات را در مورد آرایه‌ها به دست می‌دهند که عبارتند از: (۱) نام آرایه (۲) نوع داده آرایه و (۳) مجموعه اندیسهای آرایه.

مثال ۲-۴

(الف) فرض کنید DATA یک آرایه خطی 6 عنصری از اعداد اعشاری باشد. زبانهای برنامه‌نویسی مختلف زیر این آرایه را به صورت زیر معرفی می‌کنند:

FORTRAN:	<code>REAL DATA(6)</code>
PL/1:	<code>DECLARE DATA(6) FLOAT;</code>
Pascal:	<code>VAR DATA: ARRAY[1..6] OF REAL</code>

ما این آرایه را در صورت نیاز به صورت $\text{DATA}(6)$ معرفی خواهیم کرد. متن مورد بررسی نوع داده را مشخص می‌کند از این رو نوع آرایه به صراحت بیان نمی‌شود.

(ب) آرایه صحیح AUTO را با کران پائین $\text{LB} = 1932$ و کران بالای $\text{UB} = 1984$ در نظر بگیرید. زبانهای برنامه‌نویسی مختلف زیر این آرایه را به صورت زیر معرفی می‌کنند:

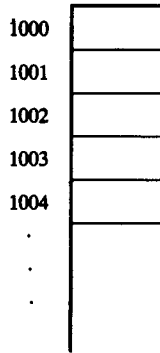
FORTRAN 77	<code>INTEGER AUTO(1932:1984)</code>
PL/1:	<code>DECLARE AUTO(1932:1984) FIXED;</code>
Pascal	<code>VAR AUTO: ARRAY[1932..1984] of INTEGER</code>

ما این آرایه را به صورت $\text{AUTO}(1932 : 1984)$ معرفی خواهیم کرد.

برخی از زبانهای برنامه‌نویسی مانند **FORTRAN** و **PASCAL** برای آرایه‌ها به صورت ایستا یعنی در زمان کامپایل برنامه حافظه اختصاص می‌دهند. از این رو اندازه آرایه در طول اجرای برنامه ثابت است. از طرف دیگر، برخی از زبانهای برنامه‌نویسی اجازه می‌دهند یک عدد صحیح n خوانده شود و آنگاه آرایه n عنصری را معرفی می‌کنند. در این گونه زبانهای برنامه‌نویسی گفته می‌شود حافظه به صورت پویا به برنامه اختصاص می‌یابد.

۳-۴ نمایش آرایه‌های خطی در حافظه

فرض کنید LA یک آرایه خطی در حافظه کامپیوتر باشد. یادآوری می‌کنیم که حافظه کامپیوتر فقط یک دنباله از مکانهای دارای آدرس است که در شکل ۲-۴ به تصویر کشیده شده است.



شکل ۲-۴- حافظه کامپیوتر

اجازه دهید از نماد زیر استفاده کنیم:

$$\text{LOC}(\text{LA}[\text{K}]) = \text{آدرس عنصر LA}[\text{K}] \text{ ی آرایه LA}$$

همانگونه که قبلاً متذکر شدیم، عناصر LA در خانه‌های متوالی ذخیره می‌شوند. بنابراین لازم نیست کامپیوتر آدرس هر عنصر LA را داشته باشد بلکه فقط لازم است آدرس اولین عنصر LA را بداند که ما آن را به صورت زیر نمایش می‌دهیم:

$$\text{Base}(\text{LA})$$

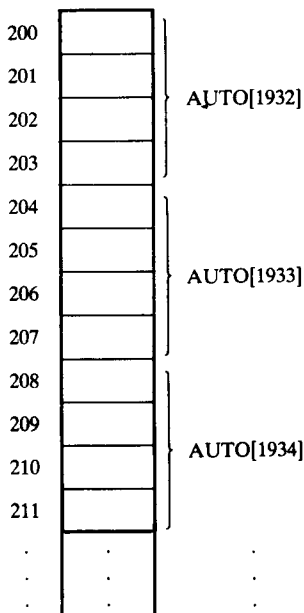
و آن را آدرس پایه یا مبنای LA می‌نامیم. با استفاده از آدرس پایه $\text{Base}(\text{LA})$ ، کامپیوتر آدرس تمام عناصر دیگر آرایه LA را به کمک فرمول زیر محاسبه می‌کند:

$$\text{LOC}(\text{LA}[\text{K}]) = \text{Base}(\text{LA}) + w(\text{K} - \text{lower bound}) \quad (۲-۴)$$

که در آن w تعداد کلمات موجود در خانه حافظه برای آرایه LA است. ملاحظه می‌کنید که زمان موردنیاز برای محاسبه $\text{LOC}(\text{LA}[\text{K}])$ اساساً برابر زمان موردنیاز برای هر مقدار دیگر K است. علاوه بر این با معلوم بودن اندیس K، می‌توان بدون جستجوی هر عنصر دیگر LA مکان $\text{LA}[\text{K}]$ را پیدا کرد و به محتوای آن دسترسی پیدا کرد.

مثال ۳-۴

آرایهٔ AUTO مثال ۱-۴ (ب) را در نظر بگیرید که تعداد اتومبیل‌های فروخته‌شدهٔ هر سال از ۱۹۳۲ تا ۱۹۸۴ را ثبت می‌کند. فرض کنید AUTO در حافظهٔ به صورت شکل ۳-۴ به نمایش در آمده است.



شکل ۳-۴

یعنی $\text{Base}(\text{AUTO}) = 200$ و $w=4$ کلمه در خانهٔ حافظه برای AUTO باشد، آنگاه

$$\text{LOC}(\text{AUTO}[1932]) = 200, \quad \text{LOC}(\text{AUTO}[1933]) = 204, \quad \text{LOC}(\text{AUTO}[1934]) = 208, \dots$$

آدرس هر عنصر آرایه برای سال $K=1965$ را می‌توان با استفاده از رابطهٔ (۲-۴) به دست آورد.

$$\text{LOC}(\text{AUTO}[1965]) = \text{Base}(\text{AUTO}) + w(1965 - \text{lower bound}) = 200 + 4(1965 - 1932) = 332$$

مجدداً تأکید می‌کنیم که محتوای این عنصر را می‌توان بدون جستجوی هر عنصر دیگر آرایهٔ AUTO پیدا کرد.

توجه کنید: مجموعه‌ای از عناصر داده‌ای A را اندیس‌دار گویند هرگاه بتوان هر عنصر دلخواه A را که آن را A_k می‌نامیم در زمانی مستقل از K مکان‌یابی کرد یا مورد پردازش قرار داد. بحث بالا بیان می‌کند که

آرایه‌های خطی را می‌توان اندیس‌دار کرد. این یک خاصیت بسیار مهم آرایه‌های خطی است. در واقع لیستهای پیوندی که موضوع فصل بعد کتاب را تشکیل می‌دهد دارای چنین خاصیتی نیست.

✓ ۴-۴ پیمایش آرایه‌های خطی

فرض کنید A یک مجموعه از عناصر داده‌ای ذخیره شده در حافظه کامپیوتر است. هرگاه بخواهیم محتوای هر عنصر A را چاپ کنیم یا تعداد عناصر A را که دارای یک خاصیت داده شده هستند بشماریم این عمل را می‌توان با پیمایش A یعنی با دسترسی و پردازش هر عنصر A دقیقاً یکبار (که اغلب ملاقات نامیده می‌شود) انجام داد.

الگوریتم زیر یک آرایه خطی LA را پیمایش می‌کند. سادگی این الگوریتم از این واقعیت ناشی می‌شود که LA ساختار خطی دارد. ساختارهای خطی دیگری نظیر لیستهای پیوندی را می‌توان به سادگی پیمایش کرد. از طرف دیگر، پیمایش ساختارهای غیرخطی نظیر درختها و گرافها به میزان قابل توجهی پیچیده و مشکل است.

Algorithm 4.1: (Traversing a Linear Array) Here LA is a linear array with lower bound LB and upper bound UB . This algorithm traverses LA applying an operation $PROCESS$ to each element of LA .

1. [Initialize counter.] Set $K := LB$.
2. Repeat Steps 3 and 4 while $K \leq UB$.
3. [Visit element.] Apply $PROCESS$ to $LA[K]$.
4. [Increase counter.] Set $K := K + 1$.
- [End of Step 2 loop.]
5. Exit.

علاوه بر این در شکل دیگری از این الگوریتم می‌توان از یک حلقه تکرار $Repeat-For$ بجای حلقه $Repeat-While$ استفاده نمود.

Algorithm 4.1': (Traversing a Linear Array) This algorithm traverses a linear array LA with lower bound LB and upper bound UB .

1. Repeat for $K = LB$ to UB :
 Apply $PROCESS$ to $LA[K]$.
 [End of loop.]
2. Exit.