

CHƯƠNG TRÌNH DỊCH

Chương 4. Phân tích cú pháp Thuật toán Earley

TS. Phạm Văn Cảnh Khoa Công nghệ thông tin

Email: canh.phamvan@phenikaa-uni.edu.vn

Nội dung



- 1. Giới thiệu
- 2. Ý tưởng cơ bản
- 3. Mã minh họa
- 4. Ví dụ
- 5. Đánh giá thuật toán
- 6. Bài tập

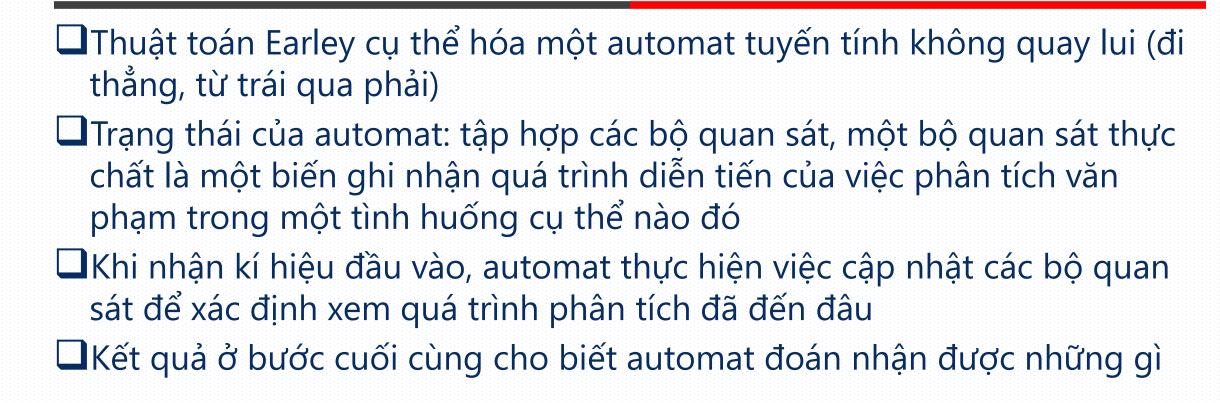
1. Giới thiệu



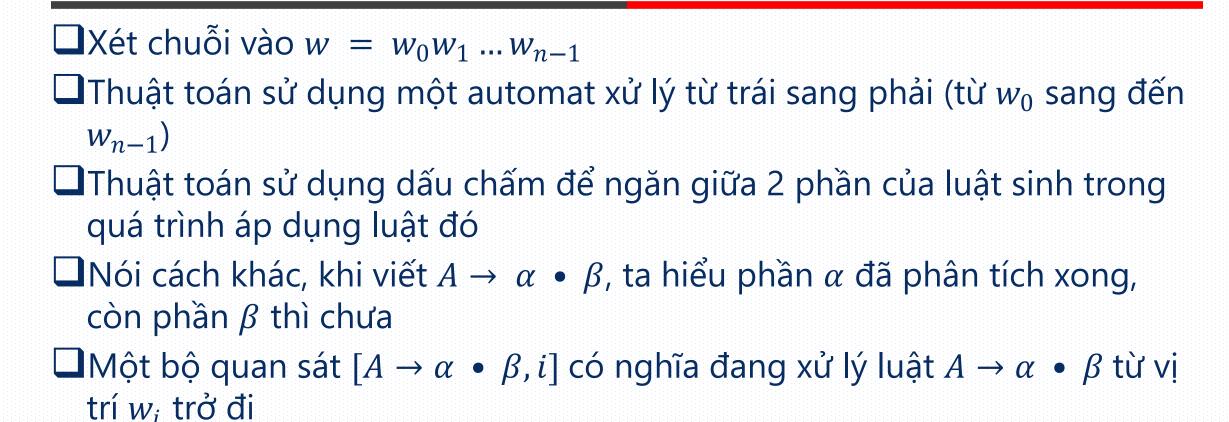
- ☐ Thuật toán Earley được giới thiệu năm 1968 bởi Jay Earley (nhà khoa học máy tính và tâm lý học, người Mỹ)
 - Công trình về phân tích văn phạm được đánh giá là một trong 25 bài báo xuất sắc nhất của tạp chí "Communications of the A.C.M" trong 1/4 thế kỷ
 - Earley nổi tiếng hơn trong ngành tâm lý học lâm sàng, chuyên về trị liệu nhóm, tác giả của Pattern System













- Khi automat xét đến kí hiệu w_m, có thể có nhiều phương án phân tích khác nhau, tất cả các phương án này đều được lưu lại để sử dụng trong các bước tính toán tiếp theo.
 Tập hợp S(m): tập các bộ quan sát dừng tại vị trí m
 Như vậy, nếu [A → α β,i] thuộc S(m) có nghĩa là dãy w_iw_{i+1} ... w_{m-1} được đoán nhận bởi phần α trong luật sinh A → α β
- Thuật toán cần phải sinh mọi thành phần trong S(m) trước khi chuyển sang kí hiệu w_{m+1}



- ☐ Thuật toán sẽ tính lần lượt S(0), S(1),..., S(n) ☐ Để dễ dàng thực hiện thuật toán, thuật toán bổ sung luật $P \to S$ vào tập luật (gọi là start rule) và bổ sung bộ $[P \to S, 0]$ vào S(0) ☐ Khi nhận kí hiệu w_m , automat sẽ bổ sung vào S(m) các bộ quan sát phù
- $lue{L}$ Khi nhạn kĩ hiệu w_m , automat sẽ bỏ sung vào S(m) các bộ quan sát phủ hợp, quá trình tính S(m) dừng khi không còn bộ quan sát nào có thể thêm vào
- □ Sau khi tính xong S(n), nếu bộ $[P \rightarrow S \bullet, 0]$ thuộc S(n) có nghĩa là dãy $w_0w_1 \dots w_{n-1}$ có thể sinh bởi S

2. Ý tưởng: 3 lệnh cơ bản



- **□Prediction (dự đoán):** với mọi bộ $[X \to \alpha \bullet Y \beta, j]$ thuộc S(k), ta tìm mọi luật sinh dạng $Y \to \gamma$ và bổ sung bộ $[Y \to \bullet \gamma, k]$ vào S(k)
- **□Scanning (xét duyệt):** với kí hiệu kết thúc a = wk, tìm mọi bộ $[X \to \alpha \bullet a \beta, j]$ thuộc S(k), bổ sung vào S(k+1) bộ $[X \to \alpha a \bullet \beta, j]$
- **□Completion (hoàn thành):** với mọi bộ $[X \to \gamma \bullet, j]$ thuộc S(k), tìm trong S(j) mọi bộ $[Y \to \alpha \bullet X \beta, i]$, bổ sung $[Y \to \alpha X \bullet \beta, i]$ vào S(k)





```
function EARLEY-PARSE(words, grammar)
    ENQUEUE((\gamma \rightarrow \bullet S, \theta), chart[\theta])
    for i \leftarrow from 0 to LENGTH(words) do
         for each state in chart[i] do
              if INCOMPLETE?(state) then
                  if NEXT-CAT(state) is a nonterminal then
                       PREDICTOR(state, i, grammar)
                  else do
                       SCANNER(state, i)
             else do
                  COMPLETER(state, i)
         end
    end
    return chart
```





```
procedure PREDICTOR((A \rightarrow \alpha \bullet B, i), j, grammar)
      for each (B \rightarrow \gamma) in GRAMMAR-RULES-FOR(B, grammar) do
           ADD-TO-SET((B \rightarrow \bullet \gamma, j), chart[j])
      end
procedure SCANNER((A \rightarrow \alpha \bullet B, i), j)
      if B \subset PARTS-OF-SPEECH(word[j]) then
           ADD-TO-SET((B \rightarrow word[j], i), chart[j + 1])
      end
procedure COMPLETER((B \rightarrow \gamma \bullet, j), k)
      for each (A \rightarrow \alpha \bullet B\beta, i) in chart[j] do
           ADD-TO-SET((A \rightarrow \alpha B \bullet \beta, i), chart[k])
      end
```



```
Bộ luật:
```

$$P \to S$$
 $S \to S + M \mid M$
 $M \to M * T \mid T$
 $T \to 1 \mid 2 \mid 3 \mid 4$

Chuỗi
$$w = 2 + 3 * 4$$

// start rule



```
S(0): \bullet 2 + 3 * 4
     (1) P \rightarrow \bullet S
                                     (0)
                                               # start rule
                                               # predict từ (1)
     (2) S \rightarrow \bullet S + M
                                     (0)
     (3) S \rightarrow \bullet M
                                     (0)
                                               # predict từ (1)
                                     (0)
                                               # predict từ (3)
     (4) \quad \mathsf{M} \to \bullet \; \mathsf{M} \; * \; \mathsf{T}
     (5) \quad \mathsf{M} \to \bullet \ \mathsf{T}
                                               # predict từ (3)
                                     (0)
                                               # predict từ (5)
     (6) T \rightarrow \bullet number
                                     (0)
S(1): 2 \bullet + 3 * 4
     (1) T \rightarrow number \bullet
                                     (0)
                                               # scan từ S(0)(6)
     (2)
                                     (0)
                                               # complete từ (1) và S(0)(5)
            M \rightarrow T \bullet
                                     (0)
                                               # complete từ (2) và S(0)(4)
           M \rightarrow M \bullet * T
                                     (0)
                                               # complete từ (2) và S(0)(3)
            S \rightarrow M \bullet
     (5) S \rightarrow S \bullet + M
                                     (0)
                                               # complete từ (4) và S(0)(2)
     (6) P \rightarrow S \bullet
                                               # complete từ (4) và S(0)(1)
                                     (0)
```



```
S(1): 2 \bullet + 3 * 4
     (1) T \rightarrow number \bullet
                                   (0)
                                             # scan từ S(0)(6)
     (2)
           M \rightarrow T \bullet
                                   (0)
                                             # complete từ (1) và S(0)(5)
     (3) \quad M \rightarrow M \bullet * T
                                   (0)
                                             # complete từ (2) và S(0)(4)
     (4) S \rightarrow M \bullet
                                   (0)
                                             # complete từ (2) và S(0)(3)
     (5) S \rightarrow S \bullet + M
                                   (0)
                                             # complete từ (4) và S(0)(2)
     (6) P \rightarrow S \bullet
                                   (0)
                                             # complete từ (4) và S(0)(1)
S(2): 2 + \bullet 3 * 4
     (1) S \rightarrow S + \bullet M
                                   (0)
                                             # scan từ S(1)(5)
           M \rightarrow \bullet M * T
                                   (2)
                                             # predict từ (1)
     (3) \quad \mathsf{M} \to \bullet \ \mathsf{T}
                                   (2)
                                             # predict từ (1)
                                   (2)
                                             # predict từ (3)
     (4) T \rightarrow \bullet number
```



```
S(2): 2 + \bullet 3 * 4
     (1) S \rightarrow S + \bullet M
                                  (0)
                                           # scan từ S(1)(5)
     (2) M \rightarrow M * T
                                  (2)
                                           # predict từ (1)
                                  (2)
                                           # predict từ (1)
     (3) \quad \mathsf{M} \to \bullet \ \mathsf{T}
     (4) T \rightarrow \bullet number
                                  (2)
                                           # predict từ (3)
S(3): 2 + 3 \bullet * 4
     (1) T \rightarrow number \bullet
                                           # scan từ S(2)(4)
                                  (2)
     (2)
                                  (2)
                                           # complete từ (1) và S(2)(3)
           M \rightarrow T \bullet
     (3)
                                  (2)
                                            # complete từ (2) và S(2)(2)
          M \rightarrow M \bullet * T
     (4)
                                  (0)
                                            # complete từ (2) và S(2)(1)
          S \rightarrow S + M \bullet
     (5) S \rightarrow S \bullet + M
                                  (0)
                                            # complete từ (4) và S(0)(2)
     (6) P \rightarrow S \bullet
                                  (0)
                                            # complete từ (4) và S(0)(1)
```



```
S(3): 2 + 3 • * 4
           T → number •
                                 (2)
                                           # scan từ S(2)(4)
                                 (2)
                                           # complete từ (1) và S(2)(3)
           M \rightarrow T \bullet
          M \rightarrow M \bullet * T
                                 (2)
                                           # complete từ (2) và S(2)(2)
                                           # complete từ (2) và S(2)(1)
                                 (0)
           S \rightarrow S + M \bullet
                                 (0)
                                           # complete từ (4) và S(0)(2)
     (5) S \rightarrow S \bullet + M
                                 (0)
                                           # complete từ (4) và S(0)(1)
     (6) P \rightarrow S \bullet
S(4): 2 + 3 * • 4
     (1) \quad \mathsf{M} \to \mathsf{M}^* \bullet \mathsf{T}
                                 (2)
                                           # scan từ S(3)(3)
     (2) T \rightarrow \bullet number
                                 (4)
                                           # predict từ (1)
```



```
S(4): 2 + 3 * • 4
    (1) \quad \mathsf{M} \to \mathsf{M}^* \bullet \mathsf{T}
                              (2)
                                        # scan từ S(3)(3)
                              (4)
    (2) T \rightarrow \bullet number
                                        # predict từ (1)
S(5): 2 + 3 * 4 \bullet
    (1) T \rightarrow number \bullet
                               (4)
                                        # scan từ S(4)(2)
                              (2)
    (2) M \rightarrow M * T \bullet
                                        # complete từ (1) và S(4)(1)
    (3) \quad M \rightarrow M \bullet * T
                           (2)
                                        # complete từ (2) và S(2)(2)
                           (0)
    (4) S \rightarrow S + M \bullet
                                        # complete từ (2) và S(2)(1)
                            (0)
    (5) S \rightarrow S \bullet + M
                                        # complete từ (4) và S(0)(2)
    (6) P \rightarrow S \bullet
                                (0)
                                        # complete từ (4) và S(0)(1)
```

Bộ [P → S •, 0] thuộc S(5), như vậy có thể kết luận chuỗi w được suy dẫn từ P

5. Đánh giá thuật toán



- Nhiều phiên bản cài đặt sau này có sửa đổi chút ít so với thuật toán gốc (thuật toán được giới thiệu trong slide này cũng không phải thuật toán gốc)
- Là một sự kết hợp thông minh của 3 trường phái
 - Tiếp cận top-down (bước prediction)
 - Tiếp cận bottom-up (bước scanning và completion)
 - Quy hoạch động (lưu lại trạng thái để dùng lại)
- ☐ Không bị hạn chế văn phạm đầu vào
 - O Do là top-down nên không bị hạn chế bởi suy dẫn rỗng
 - Dùng quy hoạch động không bị hạn chế bởi ký hiệu đệ quy (hoặc đệ quy trái)

5. Đánh giá thuật toán



Làm việc trực tiếp với luật dạng CFG: không cần phải tách thành các luật chuẩn chomsky, vì vậy kích cỡ tập luật không quá lớn \square Trong tình huống tổng quát: có độ phức tạp tính toán $O(n^3)$ với n là độ dài chuỗi vào \square Nếu văn phạm không có nhập nhằng: độ phức tạp tính toán cỡ $O(n^2)$ Nếu văn phạm đơn giản (dạng LL, LR,...): độ phức tạp cận tuyến tính $\sim O(n)$ Thực hiện đặc biệt tốt nếu văn phạm đệ quy trái

6. Bài tập

