

Thuật toán CKY – CKY Parsing Algorithm simulation

1. Giới thiệu thuật toán CKY

CKY (Coke- Kasami – Younger) là một thuật toán cải tiến của thuật toán phân tích cú pháp Bottom-Up (Bottom-Up Parsing là một chiến lược phân tích cú pháp bắt đầu từ các từ trong các chuỗi đầu vào và xây dựng các thành tố cú pháp).

CKY có thể tránh được những cách phân tích cú pháp không hợp lý so với thuật toán Bottom-Up thông thường. Do CKY sử dụng một hình thức văn phạm đặc biệt được gọi là Chomsky Normal Form (CNF). Các giải pháp trung gian của thuật toán được lưu trữ và chỉ triển khai những giải pháp trung gian nào có khả năng đóng góp vào việc phân tích đầy đủ cấu trúc cú pháp câu.

2. Mã giả thuật toán CKY

– CKY Algorithm (Recognition)

```
function CKY-Parse (words, grammar) returns table

  for j ← 1 to length(words) do: (loop over columns)

    table[j-1,j] ← {A | A → words[j] ∈ grammar} (add POS)

    for i ← j-2 downto 0 do: (loop over rows, backwards)

      for k ← i+1 to j-1 do: (loop over contents of cell)

        table[i,j] ← table[i,j] ∪

          {A | A → B C ∈ grammar,

            B ∈ table[i,k]

            C ∈ table[k,j] }
```

– CKY Parsing:

```

function CKY-Parse (words, grammar) returns parses

    for j ← 1 to length(words) do: (loop over columns)

        table[j-1,j] ← for all {A|A → words[j] ∈ grammar} (add all POS)

        for i ← j-2 downto 0 do: (loop over rows, backwards)

            for k ← i+1 to j-1 do: (loop over contents of cell)

                for all {A|A → B C}: (all productions)

                    back[i,j,A] ← { k,B,C } (add back pointer)

    return buildtree(back[1, length(words),S], table[1,LENGTH(words),S]
    (follow back pointer)

```

3. Các vấn đề của CKY

Tính hiệu quả:

- CKY có thể được thực hiện với thời gian: $O(n^3)$, trong đó n =số từ của câu.
- Sự phức tạp của các vòng lặp trong.
- Nhiều quy tắc hơn, thì ít hiệu quả hơn, nhưng điều này làm tăng một tỉ lệ hằng $L = r^2$ với r là số lượng của các biến (non-terminals).

Ngữ pháp đòi hỏi:

- Các thuật toán cơ bản đòi hỏi một ngữ pháp nhị phân: ngữ pháp CNF (Chomsky Normal Form).
- Thuật toán cơ bản có thể được mở rộng để phân tích cho CFGs tùy ý.
- Tuy nhiên, việc chuyển đổi thành ngữ pháp CNF dễ dàng và hiệu quả hơn so với phân tích với ngữ pháp tùy ý.
- Thuật toán Earley cho phép phân tích CFGs tùy ý.

4. Ngữ pháp Chomsky Normal Form

Một ngữ pháp phi ngữ cảnh mà RHS của mỗi quy tắc đưa ra là: 2 non-terminals hoặc 1 terminal. Chúng có thể là:

- Không quy tắc lẫn lộn (NP → the NN).
- Không có dạng NP → NNP, ngoại trừ dạng NN → dog.
- Về phải không có nhiều hơn 2 non-terminals như: VP → VBZ NP PP.

Bất kỳ CNG nào cũng có thể biến đổi thành một ngữ pháp tương đương yếu trong CNF tức là chúng tạo ra cũng một bộ chuỗi (các câu).

Quy ước đặt tên, ký hiệu trong văn phạm CNF:

- Sử dụng ký hiệu mới (nhị phân):
 - X_1, X_2, \dots, X_n
 - $S \rightarrow NP \ VP \ PUNC$ trở thành :
 - $S \rightarrow NP \ X_1,$
 - $X_1 \rightarrow VP \ PUNC$
- Xóa một ký hiệu:
 - $SBAR \rightarrow S, S \rightarrow NP \ VP$ trở thành
 - $SBAR \rightarrow NP \ VP$

Thuật toán CNF cải tiến:

1. Loại bỏ các unit-productions:

while there is a unit-production $A \rightarrow B$,

Remove $A \rightarrow B$.

For each $B \rightarrow u$, add $A \rightarrow u$.

2. Loại bỏ các terminals trong quy tắc lẫn lộn

For each production $A \rightarrow B_1 B_2 \dots B_k$, containing a terminal x

Add new non-terminal/production $X_1 \rightarrow x$ (unless it has already been added)

Replace every $B_i = x$ with X_1

3. Loại bỏ các quy tắc với nhiều hơn 2 nonterminals trên RHS (nhị phân)

For each rule p of form $A \rightarrow B_1 B_2 \dots B_k$

replace p with

$A \rightarrow B_1 X_1,$

$X_1 \rightarrow B_2 X_2,$

$X_2 \rightarrow B_3 X_3, \dots,$

$X_{(k-2)} \rightarrow B_{k-1} B_k$ (X_i là các biến mới)

5. Phân tích một câu đơn giản theo văn phạm CNF

Câu: •0 the •1 chef •2 eats •3 fish •4 with •5 the •6 chopsticks •7

Phân tích theo văn phạm CNF là:

S → NP VBZ

S → NP VP

VP → VP PP

VP → VBZ NP

VP → VBZ PP

VP → VBZ NNS

VP → VBZ VP

VP → VBP PP

NP → DT NN

NP → DT NNS

PP → IN NP

DT → the

NN → chef

NNS → fish

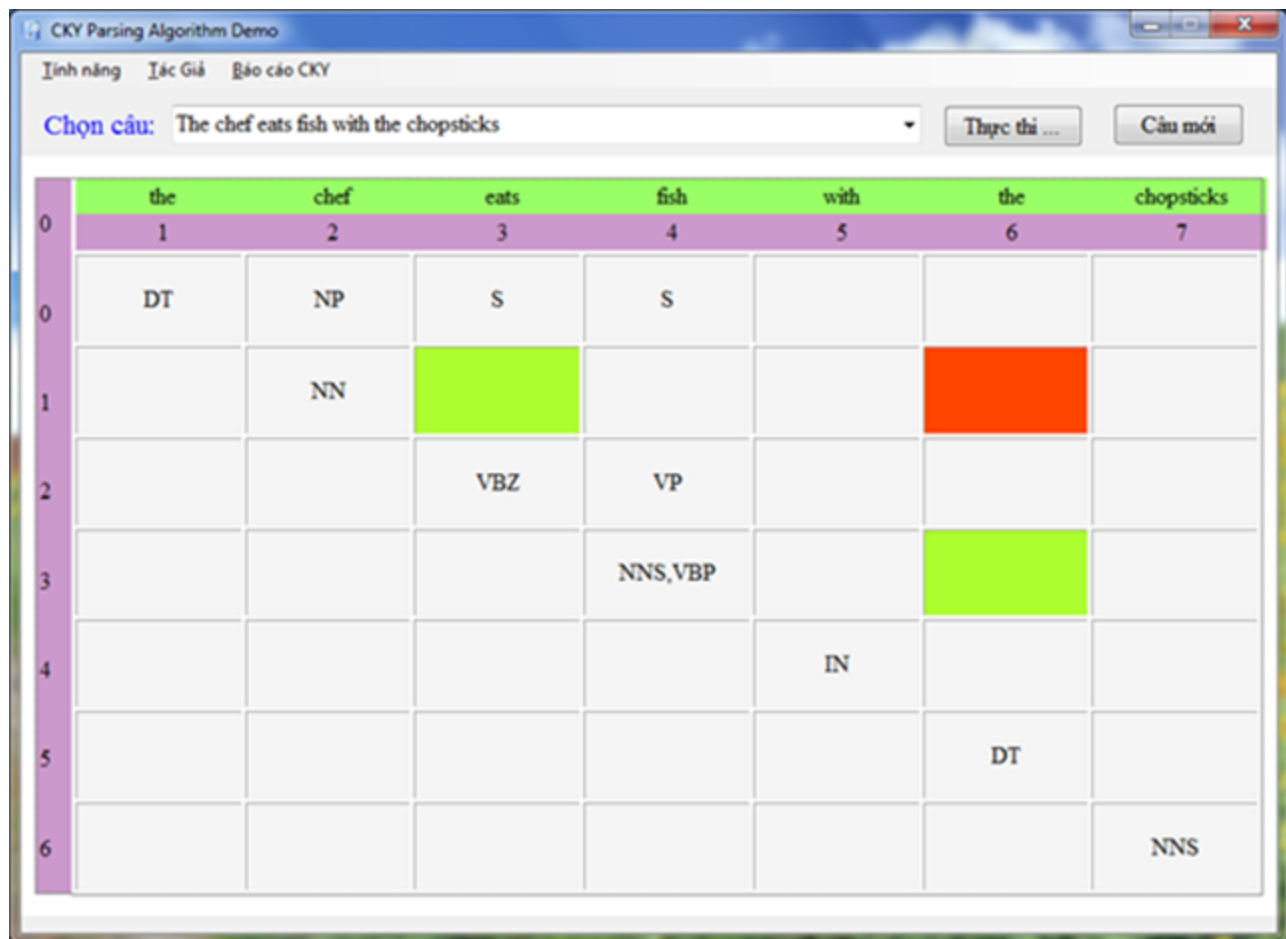
NNS → chopsticks

VBP → fish

VBZ → eats

IN → with

6. Demo CKY algorithm



BÀI TẬP

1. Cho văn phạm G:

$$S \rightarrow AB \mid XB$$

$$T \rightarrow AB \mid XB \quad X \rightarrow AT$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Chỉ ra quá trình thực hiện thuật toán CYK với $w = \mathbf{aaabbb}$

2. Cho văn phạm G:

$$S \rightarrow AA \mid AS \mid b$$

$$A \rightarrow SA \mid AS \mid a$$

Chỉ ra quá trình thực hiện thuật toán CYK với $w = \mathbf{abaab}$

3. Sử dụng thuật toán CYK để chỉ ra cây phân tích cho chuỗi $(5+7)*3$ thuộc văn phạm G

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid \text{số}$$

4. Chỉ ra cây phân tích của chuỗi **true and not false** sinh bởi thuật toán CYK với tập luật văn phạm G

$$E \rightarrow E \text{ and } T \mid T$$

$$T \rightarrow T \text{ or } F \mid F$$

$$F \rightarrow \text{not } F \mid (E) \mid \text{true} \mid \text{false}$$

Lưu ý: Cài đặt thuật toán CKY và kiểm tra kết quả thực hiện các câu trên