

Biçimsel Diller ve Otomata Teorisi

Hafta 11:
İçerikten Bağımsız Diller (I. Bölüm)

Plan

1. Düzenli İfade Uygulamaları
2. Chomsky Normal Form
3. CKY Algoritması

Düzenli İfade Uygulamaları

Düzenli ifadeler, uçsuz bucaksız gen dizilimlerinde spesifik bir bölüm aranırken sıklıkla kullanılan bir yöntemdir.

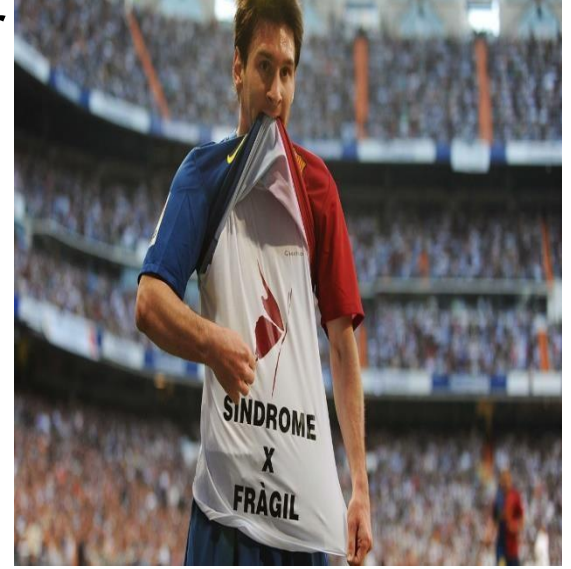
ör. Kırılgan X Sendromu (Fragile X Syndrome)

Kırılgan X sendromu, genetik bir bozukluktur ve zeka geriliğine yol açar. Erkeklerde her 250 kadınlarda ise her 800 kisiden birinde bu hastalığa neden olan gen bulunur.

Bu hastalığa neden olan DNA dizilimi şu şekildedir:

‘gcg’ nukleik asit üçlüsünün ardından ‘cgg’ veya

‘agg’ herhangi bir sayıda tekrarlar ve ardından ‘ctg’ nükleik asit üçlüsü gelir.



Bu D N A dizilimini şu düzenli ifade ile tanımlıyabiliriz:

$$g c g (c g g \cup a g g)^* c t g$$

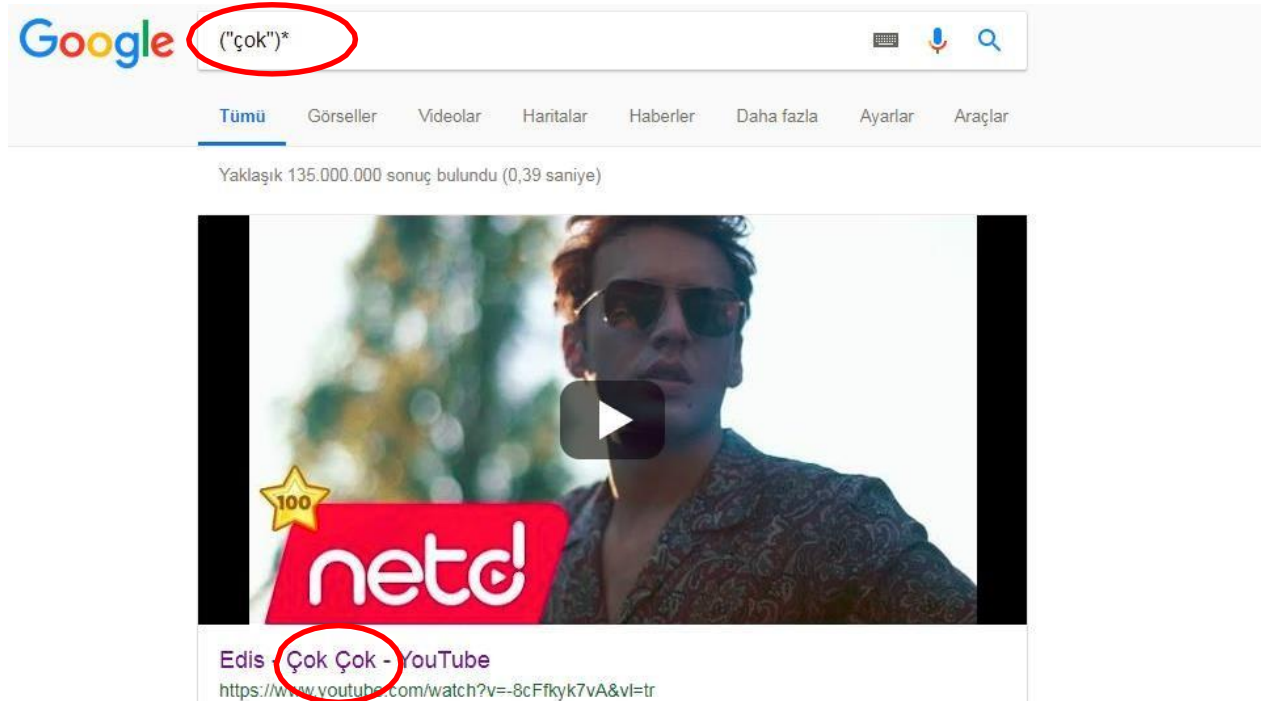
örneğin aşağıdaki D N A dizilimindeki bizim aradığımız bölüm:

gcggcgtgtgtgcgagagagtgggtttaagctg**gcgcggaggcggctg**gcgcggaggctg

Google'da Duzenli Ifade Kullanimi

Google gibi arama motorlarında da duzenli ifadeler kullanilir.

| birleşme operatörü, " " bitistirme operatörü, * yıldız operatörüdür.



Icerikten Bagimsiz Dillerde Belirsizlik Durumu

Gecen hafta en sol ve en sag turetimi gormustuk. Buna gore bir kelime birden fazla ayni yolla turetilebiliyordu (en sag yada en sol turetimle). Bu durumda ise ayni kelimeye karsilik birden fazla turetim agaci olur. Ve bir kelimenin birden fazla anlami olur.

Belirsizlikle bas etmek icin etkili bir yontem verilen grammeri Chomksky Normal Formu (CNF) ye getirmektir.

CNF ile n uzunlugundaki bir kelime her zaman $2n - 1$ adimda uretilir; boylece bir kelimenin bir dilin elemani olup olmadigi (yada bu dilin grammeri tarafaindan uretilip uretilmedigi) kolayca bulunur.

Ayrica CYK algoritmasiyla verilen bir kelimenin verilen bir grammer tarafindan uretilip uretilmedigini test ederken, grammerin CNF formunda olmasi gerekir.

Chomsky Normal Formu (CNF)

Bir icerikten bagimsiz grammerin her bir kurali asagidaki uc formdan biri seklinde yazilabiliyorsa bu grammere Chomsky normal formuna sahiptir denir:

$$1. A \rightarrow B C$$

$$2. A \rightarrow a$$

$$3. S \rightarrow \varepsilon$$

Burada A , B , C birbirinden farkli degiskenler, a bir terminal ve S baslangic degiskenidir.

Teorem: Bir icerikten bagimsiz grammer tarafından uretilen bir icerikten bagimsiz dil, Chomsky normal formuna sahiptir.

Icerikten Bagimsiz Grammerleri CNF'ye Donusturma

Bir icerikten bagimsiz grammer $G = (V, \Sigma, R, S)$, CNF'ye donusturulurken asagidaki adimlar sirayla izlenir:

1. Yeni bir baslangic degiskeni S_1 tanimlanir ve bu degisken eski baslangic degiskenine (S 'ye) $S_1 \rightarrow S$ kurali ile baglanir.
2. Sag tarafinda ε terminalini iceren butun kurallar elenir (ortadan kaldirilir).

Diyelimki bir $A \rightarrow \varepsilon$, kuralimiz olsun ve bunu eyleyelim. Bu durumda 3 farkli guncelleme yapmamiz gerekebilir.

- i.* $B \rightarrow A$, kurali $B \rightarrow \varepsilon$ ye donusur.
- ii.* $B \rightarrow uAv$, kurali $B \rightarrow uv$ ya donusur. (burada u, v birer kelimedir.)
- iii.* $B \rightarrow uAvAw$, kurali, $B \rightarrow uvw, B \rightarrow uvAw$ ve $B \rightarrow uAvw$ kurallarina donusur.

Icerikten Bagimsiz Grammerleri CNF'ye Donusturma

3 . $A \rightarrow B$ formundaki tum unit (birli) kurallar elenir.

Aslinda burda elenen B 'dir. Su halde B 'nin yerine (varsa) sol

tarafinda B olan bir kuralin sag tarafindakiler yazilir. Ornegin

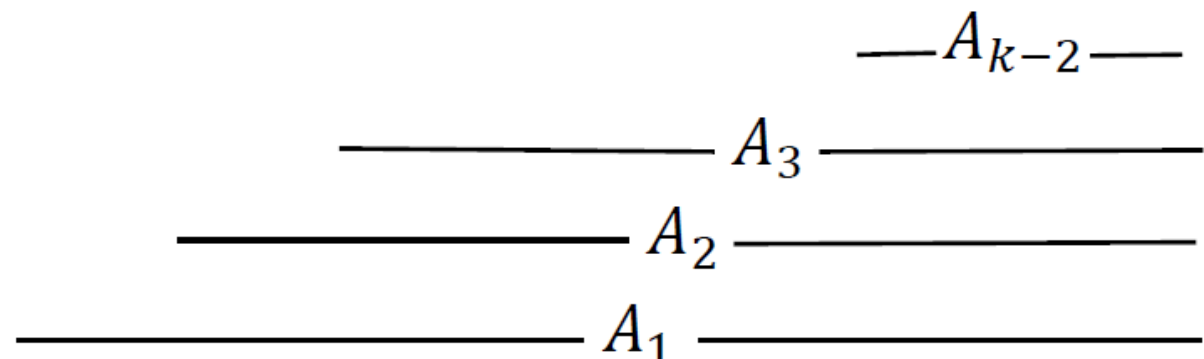
$A \rightarrow B$ ve $B \rightarrow uC|D$ olsun. Bu durumda yeni kural $A \rightarrow uC|D$

olur.

4.Sag tarafinda ikiden fazla sembol olan tum kurallar elenir.

Bunun icin

$$A \rightarrow u_1 \quad u_2 \quad u_3 \quad u_4 \quad \dots \quad u_{k-2} \quad u_{k-1} \quad u_k$$



Icerikten Bagimsiz Grammerleri CNF'ye Donusturma

$$A \rightarrow u_1 A_1$$

$$A_1 \rightarrow u_2 A_2$$

$$A_2 \rightarrow u_3 A_3$$

...

$$A_{k-3} \rightarrow u_{k-2} A_{k-2}$$

$$A_{k-2} \rightarrow u_{k-1} u_k$$

ornegin $k = 3$ icin, $A \rightarrow u_1 u_2 u_3$ icin $A \rightarrow u_1 A_1$, $A_1 \rightarrow u_2 u_3$.

5. Son olarak okun $A \rightarrow u_1 u_2$ formundaki kurallar elenir, ki burada hem u_1 hem de u_2 nin ikisi birden degisken degildir. Yani ornegin $u_1 u_2$; $0B$, $B0$, 00 , 10 gibi ifadeler olabilir. Burada u_1 ve u_2 turune gore 3 durum dusunulur:

i. u_1 terminal ($u_1 \in \Sigma$), u_2 degisken ($u_2 \in V$): Bu durumda iki yeni kural eklenir: $A \rightarrow A_1 u_2$, $A_1 \rightarrow u_1$. A_1 yeni bir degiskendir V' 'ye eklenir. ($V' = V \cup A_1$)

Icerikten Bagimsiz Grammerleri CNF'ye Donusturma

ii. u_2 terminal, u_1 degisken: Bu durumda iki yeni kural eklenir:

$A \rightarrow u_1 A_1$ ve $A_1 \rightarrow u_2$. A_1 yeni bir degiskendir V 'ye eklenir.

iii. u_1 ve u_2 terminal: Bu durumda uc yeni kural eklenir:

$A \rightarrow A_1 A_2$, $A_1 \rightarrow u_1$ ve $A_2 \rightarrow u_2$. A_1 ve A_2 yeni degiskenlerdir V 'ye eklenir ($V = V \cup A_1 \cup A_2$).

or. $G = (\{A, B\}, \{0,1\}, R, A)$ grammeri icin R kurallari soyle olsun

$$A \rightarrow BAB|B| \varepsilon$$

$$B \rightarrow 00|\varepsilon$$

bu grammeri CNF'ye donusturelim.

I. Yeni bir baslangic degiskeni olusturup bunu A 'ya gonderelim.

$$S \rightarrow A$$

$$A \rightarrow BAB|B| \varepsilon$$

$$B \rightarrow 00|\varepsilon$$

2. ε içeren tüm kuralları eleyelim.

ε içeren kurallar $A \rightarrow \varepsilon$ ve $B \rightarrow \varepsilon$ dir.

$A \rightarrow \varepsilon$ kuralını eledigimizde sağ tarafında A içeren tüm kurallar bundan etkilenir ki bunlar: $S \rightarrow A$ ve $A \rightarrow BAB$ dir. $S \rightarrow A$ ve $A \rightarrow \varepsilon$ olduğundan kuralı $S \rightarrow \varepsilon$ kuralını elde ederiz. $A \rightarrow BAB$ ve $A \rightarrow \varepsilon$ olduğundan $A \rightarrow BB$ elde edilir. Bu yeni kurallarla:

$$\begin{aligned} S &\rightarrow A \mid \varepsilon \\ A &\rightarrow BAB \mid B \mid BB \\ B &\rightarrow 00 \mid \varepsilon \end{aligned}$$

$B \rightarrow \varepsilon$ kuralını eleyelim. İlk olarak $A \rightarrow BAB$ kuralı için 3 yeni kural ekleriz (bakınız 2. adım durum iii.) Bu kurallar:

$A \rightarrow BA$, $A \rightarrow AB$ ve $A \rightarrow A$ kurallarıdır.

İkinci olarak $A \rightarrow B$ yi düşünelim. $B \rightarrow \varepsilon$ için $A \rightarrow \varepsilon$ kuralı elde edilir (bakınız 2. adım durum i.) Fakat bu yeni kuralı eklemeyiz çünkü bu kuralı zaten yukarıda elemistik.

Ucuncu olarak $A \rightarrow B$ yi dusunelim. $B \rightarrow \varepsilon$ icin yine 2. adim iii. durum geregi, $A \rightarrow B$, $A \rightarrow B$ ve $A \rightarrow \varepsilon$ yeni kurallari elde edilir. Bunlari icinden $A \rightarrow B$ kuralini bir defa aliriz, $A \rightarrow \varepsilon$ kuralini almayiz cunku zaten bu kurali daha once elemistik. Sonuc olarak 2. adim sonrasi su kurallari elde ederiz:

$$\begin{aligned} S &\rightarrow A \mid \varepsilon \\ A &\rightarrow B \mid A \mid B \mid B \mid B \mid A \mid B \mid B \mid A \mid A \\ B &\rightarrow \emptyset \end{aligned}$$

3. Unit kuralların elenmesi. Sahip oldugumuz unit kurallar $A \rightarrow A$ $S \rightarrow A$ ve $A \rightarrow B$ dir. $A \rightarrow A$ kuralini yeni bir kural eklemeyen eleyebiliriz. Kurallarimiz

$$\begin{aligned} S &\rightarrow A \mid \varepsilon \\ A &\rightarrow B \mid A \mid B \mid B \mid B \mid A \mid B \mid B \mid A \\ B &\rightarrow \emptyset \end{aligned}$$

$S \rightarrow A$ yi elemek için, bu kuralda A gordugumuz yere A dan turetilebilen $B A B \mid B \mid B B \mid A B \mid B A$ yapilarini yaziyoruz. Yeni kurallarimiz

$$S \rightarrow B A B \mid B \mid B B \mid A B \mid B A \mid \varepsilon$$

$$A \rightarrow B A B \mid B \mid B B \mid A B \mid B A$$

$$B \rightarrow 00$$

En son olarak $A \rightarrow B$ yi eliyoruz. Bunun için solunda B olan tek kural $B \rightarrow 00$ yi kullaniyoruz. 3. adım sonunda kurallarimiz:

$$S \rightarrow B A B \mid 00 \mid B B \mid A B \mid B A \mid \varepsilon$$

$$A \rightarrow B A B \mid 00 \mid B B \mid A B \mid B A$$

$$B \rightarrow 00$$

4. Sag tarafinda ikiden fazla sembol olan tum kuralların elenmesi.

Bu kurallar $S \rightarrow B A B$ ve $A \rightarrow B A B$ dir.

$S \rightarrow B A B$ i elemek için $S \rightarrow B A_1$ ve $A_1 \rightarrow A B$ kurallarini ekliyoruz.

$A \rightarrow B A B$ i elemek için $A \rightarrow B A_2$ ve $A_2 \rightarrow A B$ kurallarini ekliyoruz. 4. adim sonunda elde ettigimiz durumlar:

$$S \rightarrow B A_1 \mid 00 \mid B B \mid A B \mid B A \mid \varepsilon$$

$$A \rightarrow B A_2 \mid 00 \mid B B \mid A B \mid B A$$

$$B \rightarrow 00$$

$$A_1 \rightarrow A B$$

$$A_2 \rightarrow A B$$

5. Bu adimda sag tarafinda iki tane sembol olan (fakat ikisi birden degisken olmayan) kurallari eliyecemiz. Bu ornek icin bu kurallar, $S \rightarrow 00$, $A \rightarrow 00$ ve $B \rightarrow 00$ kurallari. 5. adim iii. durum geregi yeni ekleyecegimiz kurallar sunlar olur:

$$S \rightarrow A_3 A_3, \quad A_3 \rightarrow 0$$

$$A \rightarrow A_4 A_4, \quad A_4 \rightarrow 0$$

$$B \rightarrow A_5 A_5, \quad A_5 \rightarrow 0$$

Sonuc olarak asagidaki CNF grammerine erisilir:

$$S \rightarrow BA_1 \mid A_3A_3 \mid BB \mid AB \mid BA \mid \varepsilon$$

$$A \rightarrow BA_2 \mid A_4A_4 \mid BB \mid AB \mid BA$$

$$B \rightarrow A_5A_5$$

$$A_1 \rightarrow AB$$

$$A_2 \rightarrow AB$$

$$A_3 \rightarrow 0$$

$$A_4 \rightarrow 0$$

$$A_5 \rightarrow 0$$

or. $G = (\{A, B, S\}, \{a, b\}, R, S)$ grammeri icin R kurallari soyle olsun:

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

bu grammeri CNF ye donusturelim.

1. Yeni bir baslangic degiskeni olusturup bunu S_0 ye gonderelim.

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow A S A \mid a B \\ A &\rightarrow B \mid S \\ B &\rightarrow b \mid \varepsilon \end{aligned}$$

2. ε iceren kuralların elenmesi.

$B \rightarrow \varepsilon$ kuralını eleyelim. Yeni kurallar $A \rightarrow \varepsilon$ ve $S \rightarrow a$ olur.

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow A S A \mid a B \mid a \\ A &\rightarrow B \mid S \mid \varepsilon \\ B &\rightarrow b \end{aligned}$$

$A \rightarrow \varepsilon$ kuralını eleyelim. Yeni kurallar $S \rightarrow A S$ ve $S \rightarrow S A$ ve $S \rightarrow S$ olur.

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow A S A \mid a B \mid a \mid A S \mid S A \mid S \\ A &\rightarrow B \mid S \\ B &\rightarrow b \end{aligned}$$

3. Unit kurallari eylelim. $S \rightarrow S$ i direkt eleyebiliriz. $S_0 \rightarrow S$ kurali $S_0 \rightarrow ASA \mid aB \mid a \mid AS \mid SA \mid S$ olur.

$$S_0 \rightarrow ASA \mid aB \mid a \mid AS \mid SA$$

$$S \rightarrow ASA \mid aB \mid a \mid AS \mid SA$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

$A \rightarrow B$ kurali kurali $A \rightarrow b$ olur.

$$S_0 \rightarrow ASA \mid aB \mid a \mid AS \mid SA$$

$$S \rightarrow ASA \mid aB \mid a \mid AS \mid SA$$

$$A \rightarrow b \mid S$$

$$B \rightarrow b$$

$A \rightarrow S$ kurali $A \rightarrow ASA \mid aB \mid a \mid AS \mid SA \mid S$ olur.

$$S_0 \rightarrow ASA \mid aB \mid a \mid AS \mid SA \mid S$$

$$S \rightarrow ASA \mid aB \mid a \mid AS \mid SA \mid S$$

$$A \rightarrow b \mid ASA \mid aB \mid a \mid AS \mid SA \mid S$$

$$B \rightarrow b$$

4. Sağ tarafında ikiden fazla sembol olan kuralları eleyelim.

$S_0 \rightarrow ASA$ yerine $S_0 \rightarrow AA_1$ ve $A_1 \rightarrow SA$

$S \rightarrow ASA$ yerine $S \rightarrow AA_1$

$A \rightarrow ASA$ yerine $A \rightarrow AA_1$

yazılır.

$$S_0 \rightarrow AA_1 | aB | a | AS | SA | S$$

$$S \rightarrow AA_1 | aB | a | AS | SA | S$$

$$A \rightarrow b | AA_1 | aB | a | AS | SA | S$$

$$B \rightarrow b$$

$$A_1 \rightarrow SA$$

5. $S_0 \rightarrow aB$ yerine $S_0 \rightarrow A_2B$ ve $A_2 \rightarrow a$

$S \rightarrow aB$ yerine $S \rightarrow A_2B$; $A \rightarrow aB$ yerine $A \rightarrow A_2B$

$$S_0 \rightarrow AA_1 | A_2B | a | AS | SA | S$$

$$S \rightarrow AA_1 | A_2B | a | AS | SA | S$$

$$A \rightarrow b | AA_1 | A_2B | a | AS | SA | S$$

$$B \rightarrow b$$

$$A_1 \rightarrow SA, A_2 \rightarrow a$$

Cocke-Younger-Kasami (CYK) Algoritması

CYK algoritması Chomsky normal formunda verilen bir grammerin, verilen bir kelimeyi üretip üretmeyeceğine karar verilirken kullanılır. Yani giriş olarak CNF formunda bir

$G = (V, \Sigma, R, S)$ grammeri ve bir w kelimesi alır. Çıkış olarak eğer bu w kelimesi G grammerinden üretilirse "evet", üretilemezse "hayır" yazısını yazdırır.

ör. $G = (\{S, A, B, C\}, \{a, b\}, R, S)$ ve CNF formundaki R kuralları

$$S \rightarrow A B \mid B C$$

$$A \rightarrow B A \mid a$$

$$B \rightarrow C C \mid b$$

$$C \rightarrow A B \mid a$$

şeklinde verilen içerikten bağımsız grammer $w = b a a b a$ kelimesini üretir mi?

$$S \rightarrow AB|BC$$

$$A \rightarrow BA|a$$

$$B \rightarrow CC|b$$

$$C \rightarrow AB|a$$

| | | | | | |
|---|--------------------------------|-------------------------------|--------------------------|--------------------------|-------------------------|
| 5 | $\text{baaba} \diagup S, A, C$ | | | | |
| 4 | $\text{baab} \diagup -$ | $\text{aaba} \diagup S, A, C$ | | | |
| 3 | $\text{baa} \diagup -$ | $\text{aab} \diagup B$ | $\text{aba} \diagup B$ | | |
| 2 | $\text{ba} \diagup S, A$ | $\text{aa} \diagup B$ | $\text{ab} \diagup S, C$ | $\text{ba} \diagup S, A$ | |
| 1 | $\text{b} \diagup B$ | $\text{a} \diagup A, C$ | $\text{a} \diagup A, C$ | $\text{b} \diagup B$ | $\text{a} \diagup A, C$ |
| | b | a | a | b | a |

$$5. b \leftarrow B, a \leftarrow A, C$$

$$4. ba \leftarrow (B)(A, C) = BA, BC \leftarrow A, S$$

$$aa \leftarrow (A, C)(A, C) = AA, AC, CA, CC \leftarrow B \text{ (burada yalnızca } CC' \text{ ye ulaşılabilir)}$$

$$ab \leftarrow (A, C)(B) = AB, CB \leftarrow S, C$$

$$3. baa, (b, aa) \text{ ve } (ba, a) \text{ şeklinde ayrıştırılır.}$$

$$ba, a \leftarrow (S, A)(A, C) = SA, SC, AA, AC \text{ (bu ifadelerin hiç birine ulaşamaz)}$$

$$b, aa \leftarrow (B)(B) = BB \text{ (bu ifadeye ulaşamaz, yani buna giden bir ok yoktur)}$$

3. aab , (a, ab) ve (aa, b) şeklinde ayrıştırılır.

$$a, ab \leftarrow (A, C)(S, C) = AS, AC, CS, CC \leftarrow B$$

$$aa, b \leftarrow (B)(B) = BB \text{ (bu ifadeye ulaşamaz, yani buna giden bir ok yoktur)}$$

aba , (a, ba) ve (ab, a) şeklinde ayrıştırılır.

$$a, ba \leftarrow (A, C)(S, A) = AS, AA, CS, CA \text{ (hiç birine ulaşamaz).}$$

$$ab, a \leftarrow (S, C)(A, C) = SA, SC, CA, CC \leftarrow B$$

4. $baab$, (b, aab) , (ba, ab) , (baa, b) şeklinde ayrıştırılır.

$$b, aab \leftarrow (B)(B) = BB$$

$$ba, ab \leftarrow (S, A)(S, C) = SS, SC, AS, AC \text{ (bu ifadelere ulaşamaz)}$$

$$baa, b \leftarrow -, B$$

$aaba$, (a, aba) , (aa, ba) , (aab, a) şeklinde ayrıştırılır.

$$a, aba \leftarrow (A, C)(B) = AB, CB \leftarrow S, C$$

$$aa, ba \leftarrow (B)(S, A) = BS, BA \leftarrow A$$

$$aab, a \leftarrow (B)(A, C) = BA, BC \leftarrow A, S$$

4. $baaba, (b, aaba), (ba, aba), (baa, ba), (baab, a)$ şeklinde ayrıştırılır.

$b, aaba \leftarrow (B)(S, A, C) = BS, BA, BC \leftarrow A, S$

$ba, aba \leftarrow (S, A)(B) = SB, AB \leftarrow S, C$

$baa, ba \leftarrow -, SA$

$baab, a \leftarrow -(A, C)$

Sonuc olarak piramitin tepesinde S başlangıç değişkeni olduğu için, $w = baaba$ kelimesi verilen gramerden üretilir. Başka bir deyişle başlangıç değişkeni S den w kelimesine ulaşılabilir.

ör. $G = (\{A, B, X, T, S\}, \{a, b\}, R, S)$ ve R kuralları

$$S \rightarrow AB|XB|\varepsilon$$

$$T \rightarrow AB|XB$$

$$X \rightarrow AT$$

$$A \rightarrow a$$

$$B \rightarrow b$$

şeklinde verilen içerikten bağımsız gramer $w = aabb$ kelimesini üretir mi?

$$S \rightarrow AB|XB|\varepsilon$$

$$T \rightarrow AB|XB$$

$$X \rightarrow AT$$

$$A \rightarrow a$$

$$B \rightarrow b$$

| | | | | |
|---|-----------------|---------------|------------|-----------|
| 4 | $aabb$ / S, T | | | |
| 3 | aab / X | abb / $-$ | | |
| 2 | aa / $-$ | ab / S, T | bb / $-$ | |
| 1 | a / A | a / A | b / B | b / B |
| | a | a | b | b |

$$4. b \leftarrow B, a \leftarrow A$$

$$3. aa \leftarrow AA \text{ (bu ifadeye ulaşılamaz)}$$

$$aa \leftarrow AB \leftarrow ST, bb \leftarrow BB$$

$$2. aab, (a, ab) \text{ ve } (aa, b) \text{ şeklinde ayrıştırılır.}$$

$$a, ab \leftarrow (A)(S, T) = AS, AT \leftarrow X, aa, b \leftarrow -B$$

$$abb, (a, bb) \text{ ve } (ab, b) \text{ şeklinde ayrıştırılır.}$$

$$a, bb \leftarrow (A)-, ab, b \leftarrow (S, T)(B) = SB, TB \text{ (bu ifadelere ulaşılamaz)}$$

$$1. aabb, (a, abb), (aa, bb) \text{ ve } (aab, b) \text{ şeklinde ayrıştırılır.}$$

$$a, abb \leftarrow (A), -$$

$$aa, bb \leftarrow - -$$

$$aab, b \leftarrow (X)(B) = XB \leftarrow (S)T$$