# From EBNF to PEG

Roman R. Redziejowski

Concurrency, Specification and Programming
Berlin 2012

A way to define grammar.

## EBNF: Extended Backus-Naur Form

A way to define grammar.

$$Literal = Decimal \mid Binary$$
$$Decimal = [0\text{-}9]^+ \text{ "." } [0\text{-}9]^*$$
$$Binary = [01]^+ \text{ "B"}$$

## Recursive-descent parsing

Parsing procedure for each equation and each terminal.

$$Literal \quad = \quad Decimal \quad | \quad Binary$$
$$Decimal \quad = \quad [0\text{-}9]^{+} \text{ "." } [0\text{-}9]^{*}$$
$$Binary \quad = \quad [01]^{+} \text{ "B"}$$

*Literal* calls *Decimal* or *Binary*.
*Decimal* calls repeatedly [0-9], then ".", then repeatedly [0-9].
*Binary* calls repeatedly [01], then "B".

## Recursive-descent parsing

Parsing procedure for each equation and each terminal.

$Literal$ = $Decimal$ | $Binary$
$Decimal$ = $[0\text{-}9]^+$ "." $[0\text{-}9]^*$
$Binary$ = $[01]^+$ "B"

*Literal* calls *Decimal* or *Binary*.
*Decimal* calls repeatedly [0-9], then ".", then repeatedly [0-9].
*Binary* calls repeatedly [01], then "B".

Problem: *Decimal* and *Binary* may start
with any number of 0's and 1's.
*Literal* cannot choose which procedure to call
by looking at any fixed distance ahead.

$Literal$ = $Decimal$ | $Binary$

$Decimal$ = $[0\text{-}9]^{+}$ "." $[0\text{-}9]^{*}$

$Binary$ = $[01]^{+}$ "B"

```
101B
^
```

$$Literal = Decimal \mid Binary$$
$$Decimal = [0\text{-}9]^{+} \text{ "." } [0\text{-}9]^{*}$$
$$Binary = [01]^{+} \text{ "B"}$$

```
101B
^
```

*Literal*

*Literal*   =   *Decimal*   |   *Binary*

*Decimal*   =   [0-9]$^+$ "." [0-9]$^*$

*Binary*   =   [01]$^+$ "B"

```
101B
^
```

*Literal*→*Decimal*

## Solution: Backtracking

$$Literal = Decimal \mid Binary$$
$$Decimal = [0\text{-}9]^+ \text{ "." } [0\text{-}9]^*$$
$$Binary = [01]^+ \text{ "B"}$$

```
101B
^
```

$Literal \rightarrow Decimal \rightarrow [0\text{-}9]$

$$Literal \quad = \quad Decimal \quad | \quad Binary$$
$$Decimal \quad = \quad [0\text{-}9]^{+} \quad "." \quad [0\text{-}9]^{*}$$
$$Binary \quad = \quad [01]^{+} \quad "B"$$

<span style="color:red">101</span>B
$\hat{}$

$Literal \rightarrow Decimal \rightarrow [0\text{-}9]$ : advance 3 times

$$Literal = Decimal \mid Binary$$
$$Decimal = [0\text{-}9]^+ \; "." \; [0\text{-}9]^*$$
$$Binary = [01]^+ \; "B"$$

101B
 ^

$Literal \rightarrow Decimal \rightarrow "."$

$$Literal = Decimal \mid Binary$$
$$Decimal = [0\text{-}9]^{+} \text{ "." } [0\text{-}9]^{*}$$
$$Binary = [01]^{+} \text{ "B"}$$

```
101B
^
```

*Literal*→*Decimal*→"." : fail, backtrack

$Literal$ = $Decimal$ | $Binary$

$Decimal$ = $[0\text{-}9]^+$ "." $[0\text{-}9]^*$

$Binary$ = $[01]^+$ "B"

```
101B
^
```

$Literal$

$Literal$ = $Decimal$ | $Binary$

$Decimal$ = $[0\text{-}9]^+$ "." $[0\text{-}9]^*$

$Binary$ = $[01]^+$ "B"

```
101B
^
```

$Literal \rightarrow Binary$

$$Literal \quad = \quad Decimal \quad | \quad Binary$$
$$Decimal \quad = \quad [0\text{-}9]^+ \text{ "." } [0\text{-}9]^*$$
$$Binary \quad = \quad [01]^+ \text{ "B"}$$

```
101B
^
```

$Literal \rightarrow Binary \rightarrow [01]$

$$Literal = Decimal \mid Binary$$
$$Decimal = [0\text{-}9]^+ \ "." \ [0\text{-}9]^*$$
$$Binary = [01]^+ \ "B"$$

<span style="color:red">101</span>B
   ^

$Literal \rightarrow Binary \rightarrow [01]$ : advance 3 times

$$Literal = Decimal \mid Binary$$
$$Decimal = [0\text{-}9]^{+} \; "." \; [0\text{-}9]^{*}$$
$$Binary = [01]^{+} \; "B"$$

101B
^

$Literal \rightarrow Binary \rightarrow "B"$

*Literal* = *Decimal* | *Binary*

*Decimal* = [0-9]$^+$ "." [0-9]$^*$

*Binary* = [01]$^+$ "B"

101B
        ^

*Literal*→*Binary*→"B" : advance, return

## Solution: Backtracking

*Literal* = *Decimal* | *Binary*

*Decimal* = $[0\text{-}9]^{+}$ "." $[0\text{-}9]^{*}$

*Binary* = $[01]^{+}$ "B"

Backtracking solves the problem,
but may take exponential time.

Backtracking solves the problem,
but may take exponential time.

Solution: limited backtracking.
Never go back after one alternative succeeded.

Backtracking solves the problem,
but may take exponential time.

Solution: limited backtracking.
Never go back after one alternative succeeded.

- 1961 Brooker & Morris - Altas Compiler Compiler
- 1965 McClure - TransMoGrifier (TMG)
- 1972 Aho & Ullman - Top-Down Parsing Language (TDPL)
- ...
- 2004 Ford - Parsing Expression Grammar (PEG)

Backtracking solves the problem,
but may take exponential time.

Solution: limited backtracking.
Never go back after one alternative succeeded.

- 1961 Brooker & Morris - Altas Compiler Compiler
- 1965 McClure - TransMoGrifier (TMG)
- 1972 Aho & Ullman - Top-Down Parsing Language (TDPL)
- ...
- 2004 Ford - Parsing Expression Grammar (PEG)

It can work in linear time.

## PEG - Parsing Expression Grammar

Looks exactly like EBNF:

$$Literal = Decimal \; / \; Binary$$

$$Decimal = [0\text{-}9]^{+} \text{ "." } [0\text{-}9]^{*}$$

$$Binary = [01]^{+} \text{ "B"}$$

Looks exactly like EBNF:

$$Literal \quad = \quad Decimal \; / \; Binary$$
$$Decimal \quad = \quad [0\text{-}9]^+ \; "." \; [0\text{-}9]^*$$
$$Binary \quad = \quad [01]^+ \; "B"$$

Specification of a recursive-descent parser
with limited backtracking, where "/" means
an **ordered no-return choice**.

# PEG is not EBNF

|  | EBNF: | PEG: |
|---|---|---|
| A = ("a" / "aa") "b" | {ab, aab} | {ab} |
| A = ("aa" / "a") "ab" | {aaab, aab} | {aaab} |
| A = ("a" / "b"?) "a" | {aa, ba, a} | {aa, ba} |

|  | EBNF: | PEG: |
|---|---|---|
| A = ("a" / "aa") "b" | {ab, aab} | {ab} |
| A = ("aa" / "a") "ab" | {aaab, aab} | {aaab} |
| A = ("a" / "b"?) "a" | {aa, ba, a} | {aa, ba} |

Backtracking may examine input far ahead
so result may depend on context in front.

|  | EBNF: | PEG: |
|---|---|---|
| A = ("a" / "aa") "b" | {ab, aab} | {ab} |
| A = ("aa" / "a") "ab" | {aaab, aab} | {aaab} |
| A = ("a" / "b"?) "a" | {aa, ba, a} | {aa, ba} |

Backtracking may examine input far ahead
so result may depend on context in front.

$$A = \text{"a" } A \text{ "a" / "aa"} \qquad \text{EBNF: } a^{2n} \qquad \text{PEG: } a^{2^n}$$

In this case PEG = EBNF:

$$Literal \quad = \quad Decimal \; / \; Binary$$
$$Decimal \quad = \quad [0\text{-}9]^+ \; \text{"."} \; [0\text{-}9]^*$$
$$Binary \quad = \quad [01]^+ \; \text{"B"}$$

In this case PEG = EBNF:

$$
\begin{aligned}
\textit{Literal} &= \textit{Decimal} \;/\; \textit{Binary} \\
\textit{Decimal} &= [0\text{-}9]^{+} \; \text{"."} \; [0\text{-}9]^{*} \\
\textit{Binary} &= [01]^{+} \; \text{"B"}
\end{aligned}
$$

When does it happen?

Sérgio Queiroz de Medeiros

*Correspondência entre PEGs e Classes de Gramáticas Livres de Contexto.*

Ph.D. Thesis

Pontifícia Universidade Católica do Rio deJaneiro (2010).

Sérgio Queiroz de Medeiros

*Correspondência entre PEGs e Classes
de Gramáticas Livres de Contexto.*

Ph.D. Thesis

Pontifícia Universidade Católica
do Rio deJaneiro (2010).

If EBNF has LL(1) property then PEG = EBNF

But this is not LL(1):

$$Literal = Decimal \;/\; Binary$$
$$Decimal = [0\text{-}9]^+ \;\text{"."}\; [0\text{-}9]^*$$
$$Binary = [01]^+ \;\text{"B"}$$

But this is not LL(1):

$$Literal = Decimal \; / \; Binary$$
$$Decimal = [0\text{-}9]^+ \; "." \; [0\text{-}9]^*$$
$$Binary = [01]^+ \; "B"$$

Which means PEG = EBNF for a wider class.

But this is not LL(1):

$$Literal \quad = \quad Decimal \; / \; Binary$$
$$Decimal \; = \; [0\text{-}9]^{+} \; "." \; [0\text{-}9]^{*}$$
$$Binary \quad = \; [01]^{+} \; "B"$$

Which means PEG = EBNF for a wider class.
Let us find more about it.

## Simple grammar

Alphabet $\Sigma$ (the "terminals").

Set $N$ of names (the "nonterminals").

For each $A \in N$ one rule of the form:

- $A = e_1 \, e_2$   (*Sequence*) or
- $A = e_1 \mid e_2$   (*Choice*)

where $e_1, e_2 \in N \cup \Sigma \cup \{\varepsilon\}$.

Start symbol $S \in A$.

"Syntax expressions": $\mathbb{E} = N \cup \Sigma \cup \{\varepsilon\}$.

## Simple grammar

Alphabet $\Sigma$ (the "terminals").

Set $N$ of names (the "nonterminals").

For each $A \in N$ one rule of the form:

- $A = e_1 \, e_2$    (*Sequence*) or
- $A = e_1 \,|\, e_2$    (*Choice*)

where $e_1, e_2 \in N \cup \Sigma \cup \{\varepsilon\}$.

Start symbol $S \in A$.

"Syntax expressions": $\mathbb{E} = N \cup \Sigma \cup \{\varepsilon\}$.

Will consider two interpretations: EBNF and PEG.

$\mathcal{L}(e)$ – language of expression $e \in \mathbb{E}$.

- $\mathcal{L}(\varepsilon) = \{\varepsilon\}$
- $\mathcal{L}(a) = \{a\}$ for $a \in \Sigma$
- $\mathcal{L}(A) = \mathcal{L}(e_1)\mathcal{L}(e_2)$ for $A = e_1 e_2$
- $\mathcal{L}(A) = \mathcal{L}(e_1) \cup \mathcal{L}(e_2)$ for $A = e_1 \mid e_2$

Language defined by the grammar: $\mathcal{L}(S)$.

Relation $\overset{\text{BNF}}{\leadsto} \subseteq \mathbb{E} \times \Sigma^* \times \Sigma^*$, written $[e]\, x \overset{\text{BNF}}{\leadsto} y$.

$[e]\, xy \overset{\text{BNF}}{\leadsto} y$ means "$xy$ has prefix $x \in \mathcal{L}(e)$".

Or: parsing procedure for $e$, applied to $xy$ consumes $x$".

$w \in \mathcal{L}(S) \iff [S]\, w\$ \overset{\text{BNF}}{\leadsto} \$$
where \$ is "end of text" marker.

$[e]\ x \overset{\text{BNF}}{\leadsto} y$ holds if and only if
it can be proved using these inference rules:

$$\overline{[\varepsilon]\ x \overset{\text{BNF}}{\leadsto} x} \qquad \overline{[a]\ ax \overset{\text{BNF}}{\leadsto} x}$$

$$\frac{A = e_1 e_2 \quad [e_1]\ xyz \overset{\text{BNF}}{\leadsto} yz \quad [e_2]\ yz \overset{\text{BNF}}{\leadsto} z}{[A]\ xyz \overset{\text{BNF}}{\leadsto} z}$$

$$\frac{A = e_1 | e_2 \quad [e_1]\ xy \overset{\text{BNF}}{\leadsto} y}{[A]\ xy \overset{\text{BNF}}{\leadsto} y}$$

$$\frac{A = e_1 | e_2 \quad [e_2]\ xy \overset{\text{BNF}}{\leadsto} y}{[A]\ xy \overset{\text{BNF}}{\leadsto} y}$$

## Example of proof

Grammar: $S = aX$, $\quad X = S|b$ $\qquad$ Proof of $aab \in \mathcal{L}(S)$

$$\dfrac{}{[b]\ b\$ \overset{\text{BNF}}{\leadsto} \$}$$

$$\dfrac{}{[a]\ ab\$ \overset{\text{BNF}}{\leadsto} b\$} \qquad \dfrac{X = S|b \quad [b]\ b\$ \overset{\text{BNF}}{\leadsto} \$}{[X]\ b\$ \overset{\text{BNF}}{\leadsto} \$}$$

$$\dfrac{S = aX \quad [a]\ ab\$ \overset{\text{BNF}}{\leadsto} b\$ \quad [X]\ b\$ \overset{\text{BNF}}{\leadsto} \$}{[S]\ ab\$ \overset{\text{BNF}}{\leadsto} \$}$$

$$\dfrac{}{[a]\ aab\$ \overset{\text{BNF}}{\leadsto} ab\$} \qquad \dfrac{X = S|b \quad [S]\ ab\$ \overset{\text{BNF}}{\leadsto} \$}{[X]\ ab\$ \overset{\text{BNF}}{\leadsto} \$}$$

$$\dfrac{S = aX \quad [a]\ aab\$ \overset{\text{BNF}}{\leadsto} ab\$ \quad [X]\ ab\$ \overset{\text{BNF}}{\leadsto} \$}{[S]\ aab\$ \overset{\text{BNF}}{\leadsto} \$}$$

## PEG interpretation

Elements of $\mathbb{E}$ are parsing procedures
that consume input or return "failure".

- $\varepsilon$ returns success without consuming input.
- $a$ consumes $a$ if input starts with $a$.
  Otherwise returns failure.
- $A = e_1 \, e_2$ calls $e_1$ then $e_2$.
  If any of them failed, backtracks and returns failure.
- $A = e_1 \mid e_2$ calls $e_1$.
  If $e_1$ succeeded, returns success.
  If $e_1$ failed, calls $e_2$ and returns its result.

## "Natural semantics" (after Medeiros)

Relation $\overset{\text{PEG}}{\leadsto} \subseteq \mathbb{E} \times \Sigma^* \times (\Sigma^* \cup \text{fail})$, written $[e]\, x \overset{\text{PEG}}{\leadsto} y$.

- $[e]\, xy \overset{\text{PEG}}{\leadsto} y$ means "$e$ consumes prefix $x$ of $xy$".
- $[e]\, x \overset{\text{PEG}}{\leadsto} \text{fail}$ means "$e$ applied to $x$ returns failure".

$w$ accepted by the grammar iff $[S]\, w\$ \overset{\text{PEG}}{\leadsto} \$$.

$[e]\ x \stackrel{\text{PEG}}{\leadsto} Y$ holds if and only if
it can be proved using these inference rules:

$$\frac{}{[\varepsilon]\ x \stackrel{\text{PEG}}{\leadsto} x} \qquad \frac{}{[a]\ ax \stackrel{\text{PEG}}{\leadsto} x} \qquad \frac{b \neq a}{[b]\ ax \stackrel{\text{PEG}}{\leadsto} \text{fail}} \qquad \frac{}{[a]\ \varepsilon \stackrel{\text{PEG}}{\leadsto} \text{fail}}$$

$$\frac{A = e_1 e_2 \quad [e_1]\ xyz \stackrel{\text{PEG}}{\leadsto} yz \quad [e_2]\ yz \stackrel{\text{PEG}}{\leadsto} Z}{[A]\ xyz \stackrel{\text{PEG}}{\leadsto} Z}$$

$$\frac{A = e_1 e_2 \quad [e_1]\ x \stackrel{\text{PEG}}{\leadsto} \text{fail}}{[A]\ x \stackrel{\text{PEG}}{\leadsto} \text{fail}} \qquad \frac{A = e_1 | e_2 \quad [e_1]\ xy \stackrel{\text{PEG}}{\leadsto} y}{[A]\ xy \stackrel{\text{PEG}}{\leadsto} y}$$

$$\frac{A = e_1 | e_2 \quad [e_1]\ x \stackrel{\text{PEG}}{\leadsto} \text{fail} \quad [e_2]\ xy \stackrel{\text{PEG}}{\leadsto} Y}{[A]\ xy \stackrel{\text{PEG}}{\leadsto} Y}$$

where $Y$ is $y$ or fail and $Z$ is $z$ or fail .

By induction on the height of proof trees for
$[S]$ $w\$$ $\overset{\text{PEG}}{\leadsto}$ $\$$ and $[S]$ $w\$$ $\overset{\text{BNF}}{\leadsto}$ $\$$:

- $[S]$ $w\$$ $\overset{\text{PEG}}{\leadsto}$ $\$$ $\Rightarrow$ $[S]$ $w\$$ $\overset{\text{BNF}}{\leadsto}$ $\$$. (Medeiros)

By induction on the height of proof trees for
$[S]$ $w\$ \overset{PEG}{\leadsto} \$$ and $[S]$ $w\$ \overset{BNF}{\leadsto} \$$:

- $[S]$ $w\$ \overset{PEG}{\leadsto} \$$ $\Rightarrow$ $[S]$ $w\$ \overset{BNF}{\leadsto} \$$. (Medeiros)

- $[S]$ $w\$ \overset{BNF}{\leadsto} \$$ $\Rightarrow$ $[S]$ $w\$ \overset{PEG}{\leadsto} \$$
  if for every Choice $A = e_1 | e_2$ holds
  $\mathcal{L}(e_1) \cap \mathrm{Pref}(\mathcal{L}(e_2) \, \mathrm{Tail}(A)) = \varnothing$.

By induction on the height of proof trees for
$[S]$ $w\$ \overset{\text{PEG}}{\leadsto} \$$ and $[S]$ $w\$ \overset{\text{BNF}}{\leadsto} \$$:

- $[S]$ $w\$ \overset{\text{PEG}}{\leadsto} \$$ $\Rightarrow$ $[S]$ $w\$ \overset{\text{BNF}}{\leadsto} \$$. (Medeiros)

- $[S]$ $w\$ \overset{\text{BNF}}{\leadsto} \$$ $\Rightarrow$ $[S]$ $w\$ \overset{\text{PEG}}{\leadsto} \$$
  if for every Choice $A = e_1 | e_2$ holds
  $\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2)\,\text{Tail}(A)) = \varnothing$.

(Tail($A$) is any possible continuation after $A$:
$y \in \text{Tail}(A)$ iff proof tree of $[S]$ $w\$ \overset{\text{BNF}}{\leadsto} \$$ for some $w$
contains partial result $[A]$ $xy\$ \overset{\text{BNF}}{\leadsto} y\$$.)

Let us say that Choice $A = e_1|e_2$ is "safe" to mean
$\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2)\,\text{Tail}(A)) = \varnothing$.

Let us say that Choice $A = e_1 | e_2$ is "safe" to mean
$\mathcal{L}(e_1) \cap \mathrm{Pref}(\mathcal{L}(e_2) \mathrm{Tail}(A)) = \varnothing$.

> The two interpretations are equivalent
> if every Choice in the grammar is safe.

$$\mathcal{L}(e_1) \cap \mathrm{Pref}(\mathcal{L}(e_2)\,\mathrm{Tail}(A)) = \varnothing$$

$$\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2)\,\text{Tail}(A)) = \varnothing$$

Requires $\varepsilon \notin \mathcal{L}(e_1)$.

$$\mathcal{L}(e_1) \cap \mathrm{Pref}(\mathcal{L}(e_2)\,\mathrm{Tail}(A)) = \varnothing$$

Requires $\varepsilon \notin \mathcal{L}(e_1)$.

Depends on context.

$$\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2) \text{Tail}(A)) = \varnothing$$

Requires $\varepsilon \notin \mathcal{L}(e_1)$.

Depends on context.

Difficult to check: $\mathcal{L}(e_1)$, $\mathcal{L}(e_2)$, and $\text{Tail}(A)$
can be any context-free languages.

Intersection of context-free languages is in general
undecidable.

$$\mathcal{L}(e_1) \cap \mathrm{Pref}(\mathcal{L}(e_2)\,\mathrm{Tail}(A)) = \varnothing$$

Requires $\varepsilon \notin \mathcal{L}(e_1)$.

Depends on context.

Difficult to check: $\mathcal{L}(e_1)$, $\mathcal{L}(e_2)$, and $\mathrm{Tail}(A)$
can be any context-free languages.

Intersection of context-free languages is in general
undecidable.

Can be "approximated" by stronger conditions.

Consider $A = e_1|e_2$.

$\text{FIRST}(e_1)$, $\text{FIRST}(e_2)$:
sets of possible first letters
of words in $\mathcal{L}(e_1)$ respectively $\mathcal{L}(e_2)$.

Consider $A = e_1 | e_2$.

$\mathrm{FIRST}(e_1)$, $\mathrm{FIRST}(e_2)$:
sets of possible first letters
of words in $\mathcal{L}(e_1)$ respectively $\mathcal{L}(e_2)$.

If $\mathcal{L}(e_1), \mathcal{L}(e_2)$, do not contain $\varepsilon$,
$\mathrm{FIRST}(e_1) \cap \mathrm{FIRST}(e_2) = \varnothing$
implies $\mathcal{L}(e_1) \cap \mathrm{Pref}(\mathcal{L}(e_2)\,\mathrm{Tail}(A)) = \varnothing$.

Consider $A = e_1 | e_2$.

$\text{FIRST}(e_1)$, $\text{FIRST}(e_2)$:
sets of possible first letters
of words in $\mathcal{L}(e_1)$ respectively $\mathcal{L}(e_2)$.

If $\mathcal{L}(e_1), \mathcal{L}(e_2)$, do not contain $\varepsilon$,
$\text{FIRST}(e_1) \cap \text{FIRST}(e_2) = \varnothing$
implies $\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2)\,\text{Tail}(A)) = \varnothing$.

This is LL(1) for grammar without $\varepsilon$.
Each choice in such grammar is safe.
The two interpretations are equivalent.

To go beyond LL(1), we shall look at
first *expressions* rather than first *letters*.

$$S = X \mid Y$$
$$X = Z \mid V$$
$$Y = W X$$
$$Z = a \mid b$$
$$V = b \mid T$$
$$W = d \mid U$$
$$T = c V$$
$$U = c W$$

$S = X \mid Y$

$X = Z \mid V$

$Y = W X$

$Z = a \mid b$

$V = b \mid T$

$W = d \mid U$

$T = c V$

$U = c W$



$\text{FIRST}(X) = \{a, b, c\}$

$S = X \mid Y$

$X = Z \mid V$

$Y = W X$

$Z = a \mid b$

$V = b \mid T$

$W = d \mid U$

$T = c V$

$U = c W$



$\text{FIRST}(X) = \{a, b, c\}$

$\text{FIRST}(Y) = \{c, d\}$

$S = X \mid Y$

$X = Z \mid V$

$Y = W X$

$Z = a \mid b$

$V = b \mid T$

$W = d \mid U$

$T = c V$

$U = c W$



$\text{FIRST}(X) = \{a, b, c\}$

$\text{FIRST}(Y) = \{c, d\}$

$\{a, b, c\} \cap \{c, d\} \neq \varnothing$: $S = X \mid Y$ is not LL(1).

$S = X \mid Y$

$X = Z \mid V$

$Y = W X$

$Z = a \mid b$

$V = b \mid T$

$W = d \mid U$

$T = c V$

$U = c W$

$$S = X \mid Y$$
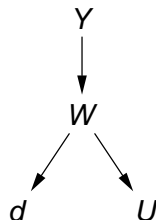$$X = Z \mid V$$
$$Y = W X$$
$$Z = a \mid b$$
$$V = b \mid T$$
$$W = d \mid U$$
$$T = c V$$
$$U = c W$$



Each word in $\mathcal{L}(X)$ has a prefix in $\{a, b\} \cup \mathcal{L}(T) = a \cup c^* b$.
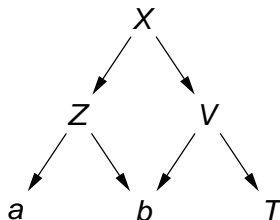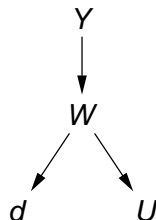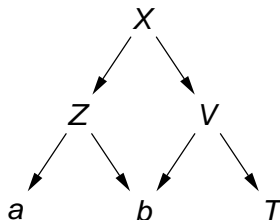
$S = X \mid Y$

$X = Z \mid V$

$Y = W X$

$Z = a \mid b$

$V = b \mid T$

$W = d \mid U$

$T = c V$

$U = c W$

Each word in $\mathcal{L}(X)$ has a prefix in $\{a, b\} \cup \mathcal{L}(T) = a \cup c^* b$.
Each word in $\mathcal{L}(Y)$ has a prefix in $\{d\} \cup \mathcal{L}(U) = d^* b$.

Each word in $\mathcal{L}(X)$ has a prefix in $a \cup c^* b$.
Each word in $\mathcal{L}(Y)$ has a prefix in $d^* b$.

$$\mathcal{L}(X) = (a \cup c^* b)(\ldots)$$
$$\mathcal{L}(Y) = (c^* d)(\ldots)$$

Each word in $\mathcal{L}(X)$ has a prefix in $a \cup c^*b$.
Each word in $\mathcal{L}(Y)$ has a prefix in $d^*b$.

$$\mathcal{L}(X) = (a \cup c^*b)(\ldots)$$
$$\mathcal{L}(Y) = (c^*d)(\ldots)$$

$$\mathcal{L}(X) \cap \mathrm{Pref}(\mathcal{L}(Y)\,\mathrm{Tail}(S))$$
$$= (a \cup c^*b)(\ldots) \;\cap\; \mathrm{Pref}((c^*d)(\ldots)(\ldots))$$

Each word in $\mathcal{L}(X)$ has a prefix in $a \cup c^* b$.
Each word in $\mathcal{L}(Y)$ has a prefix in $d^* b$.

$$\mathcal{L}(X) = (a \cup c^* b)(\ldots)$$
$$\mathcal{L}(Y) = (c^* d)(\ldots)$$

$$\mathcal{L}(X) \cap \mathrm{Pref}(\mathcal{L}(Y)\,\mathrm{Tail}(S))$$
$$= (a \cup c^* b)(\ldots) \cap \mathrm{Pref}((c^* d)(\ldots)(\ldots))$$

No word in $a \cup c^* b$ is a prefix of word in $c^* d$ and vice-versa.

Each word in $\mathcal{L}(X)$ has a prefix in $a \cup c^* b$.
Each word in $\mathcal{L}(Y)$ has a prefix in $d^* b$.

$\mathcal{L}(X) = (a \cup c^* b)(\ldots)$
$\mathcal{L}(Y) = (c^* d)(\ldots)$

$\mathcal{L}(X) \cap \mathrm{Pref}(\mathcal{L}(Y)\,\mathrm{Tail}(S))$
    $= (a \cup c^* b)(\ldots) \cap \mathrm{Pref}((c^* d)(\ldots)(\ldots))$

No word in $a \cup c^* b$ is a prefix of word in $c^* d$ and vice-versa.

The intersection is empty: $S = X | Y$ is safe.

# Some terminology

*X* starts with *a*, *b*, or *T*:
"*X* has *a*, *b*, and *T* as possible first expressions".

$$\{a, b, T\} \sqsubseteq X$$

$X$ starts with $a$, $b$, or $T$:
"$X$ has $a$, $b$, and $T$ as possible first expressions".

$$\{a, b, T\} \sqsubseteq X$$

No word in $a$, $b$, or $T$ is a prefix of a word in $d$ or $U$
and vice-versa:
"$\{a, b, T\}$ and $\{d, U\}$ are exclusive".

$$\{a, b, T\} \asymp \{d, U\}$$

If $\varepsilon \notin e_1$ and $\varepsilon \notin e_2$
and there exist $\mathrm{FIRST}_1 \sqsubseteq e_1$, $\mathrm{FIRST}_2 \sqsubseteq e_2$
such that $\mathrm{FIRST}_1 \asymp \mathrm{FIRST}_2$
then $A = e_1 | e_2$ is safe.

The two interpretations of an $\varepsilon$-free grammar
are equivalent if for every Choice $A = e_1|e_2$,
$e_1$ and $e_2$ have exclusive sets of first expressions.

- (Good news) Grammar with $\varepsilon$ is easy to handle. This involves first expressions of Tail($A$), that are obtained using the classical computation of FOLLOW.

- (Good news) Grammar with $\varepsilon$ is easy to handle. This involves first expressions of Tail($A$), that are obtained using the classical computation of FOLLOW.
- (Good news) The results for simple grammar are easily extended to full EBNF / PEG.

## Final remarks

- (Good news) Grammar with $\varepsilon$ is easy to handle. This involves first expressions of Tail($A$), that are obtained using the classical computation of FOLLOW.
- (Good news) The results for simple grammar are easily extended to full EBNF / PEG.
- (Good news) The possible sets of first expressions are easily obtained in a mechanical way.

## Final remarks

- (Good news) Grammar with $\varepsilon$ is easy to handle. This involves first expressions of Tail($A$), that are obtained using the classical computation of FOLLOW.
- (Good news) The results for simple grammar are easily extended to full EBNF / PEG.
- (Good news) The possible sets of first expressions are easily obtained in a mechanical way.
- (Bad news) Checking that they are exclusive is not easy: it is undecidable in general case (but we may hope first expressions are simple enough to be decidable.)

S = (aa|a)b    (that is: S = Xb, X = aa|a.)

$S = (aa|a)b$    (that is: $S = Xb$, $X = aa|a$.)

$\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2)\,\text{Tail}(X)) = aa \cap \text{Pref}(ab) = \varnothing$.

S = (aa|a)b    (that is: S = Xb, X = aa|a.)

$\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2) \text{Tail}(X)) = \text{aa} \cap \text{Pref}(ab) = \varnothing.$

X is safe. Both interpretations accept {aab,ab}.

$S = (aa|a)b$  (that is: $S = Xb$, $X = aa|a$.)

$\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2)\,\text{Tail}(X)) = aa \cap \text{Pref}(ab) = \varnothing$.

X is safe. Both interpretations accept {aab,ab}.

Sets of first expressions in X: {aa} and {a}. Not exclusive!

$S = (aa|a)b$    (that is: $S = Xb$, $X = aa|a$.)

$\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2) \text{Tail}(X)) = aa \cap \text{Pref}(ab) = \varnothing$.

X is safe. Both interpretations accept {aab,ab}.

Sets of first expressions in X: {aa} and {a}. Not exclusive!

There is more to squeeze out of $\mathcal{L}(e_1) \cap \text{Pref}(\mathcal{L}(e_2) \text{Tail}(A))$.

Thanks for your attention!