

# A proposal of a formal grammar for a subset of PDF

Guillaume Endignoux      Olivier Levillain      Jean-Yves Migeon

This document describes a grammar for a restricted subset of the Portable Document Format, which excludes – among others – incremental updates, object streams and linearization. The grammar is presented in BNF, in a top-down order – from the file structure to the character set.

## 1 File structure

The following rules describe the overall file structure.

```
 $\langle PDF \rangle ::= \langle header \rangle \langle body \rangle \langle xref \rangle \langle trailer \rangle \langle eof-marker \rangle$   
 $\langle header \rangle ::= \langle version \rangle \langle non-ascii-marker \rangle? \langle spaces \rangle$   
 $\langle version \rangle ::= \text{'\%PDF-1.'} \langle version-digit \rangle \langle eol \rangle$   
 $\langle non-ascii-marker \rangle ::= \text{'\%'} \langle non-ascii-char \rangle \langle non-ascii-char \rangle \langle non-ascii-char \rangle \langle non-ascii-char \rangle \langle eol \rangle$   
 $\langle body \rangle ::= \langle empty \rangle$   
           $| \langle body \rangle \langle indirect-object \rangle \langle spaces \rangle$   
 $\langle xref \rangle ::= \langle xref-header \rangle \langle xref-section \rangle$   
 $\langle xref-header \rangle ::= \text{'xref'} \langle eol \rangle \langle unsigned-int \rangle \langle space \rangle \langle unsigned-int \rangle \langle eol \rangle$   
 $\langle xref-section \rangle ::= \langle empty \rangle$   
           $| \langle xref-section \rangle \langle unsigned-int-10 \rangle \langle space \rangle \langle unsigned-int-5 \rangle \langle space \rangle (\text{'f'} | \text{'n'}) \langle xref-eol \rangle$   
 $\langle trailer \rangle ::= \text{'trailer'} \langle spaces \rangle \langle dictionary \rangle \langle eol \rangle \text{'startxref'} \langle eol \rangle \langle unsigned-int \rangle \langle eol \rangle$   
 $\langle eof-marker \rangle ::= \text{'%%EOF'} \langle eol \rangle?$ 
```

## 2 Objects

The following rules describe the syntax of direct and indirect objects.

```
 $\langle indirect-object \rangle ::= \langle indirect-object-header \rangle \langle direct-object \rangle \langle spaces \rangle \text{'endobj'}$   
           $| \langle indirect-object-header \rangle \langle dictionary \rangle \langle spaces \rangle \langle stream \rangle \langle spaces \rangle \text{'endobj'}$   
 $\langle indirect-object-header \rangle ::= \langle unsigned-int \rangle \langle spaces \rangle \langle unsigned-int \rangle \langle spaces \rangle \text{'obj'} \langle spaces \rangle$ 
```

$\langle \textit{direct-object} \rangle ::= \langle \textit{null} \rangle$   
 $\quad | \langle \textit{bool} \rangle$   
 $\quad | \langle \textit{int} \rangle$   
 $\quad | \langle \textit{real} \rangle$   
 $\quad | \langle \textit{string} \rangle$   
 $\quad | \langle \textit{name} \rangle$   
 $\quad | \langle \textit{reference} \rangle$   
 $\quad | \langle \textit{array} \rangle$   
 $\quad | \langle \textit{dictionary} \rangle$

$\langle \textit{reference} \rangle ::= \langle \textit{unsigned-int} \rangle \langle \textit{spaces} \rangle \langle \textit{unsigned-int} \rangle \langle \textit{spaces} \rangle \text{'R'}$

$\langle \textit{array} \rangle ::= \text{'['} \langle \textit{spaces} \rangle \langle \textit{array-content} \rangle \text{'}'}$

$\langle \textit{array-content} \rangle ::= \langle \textit{empty} \rangle$   
 $\quad | \langle \textit{direct-object} \rangle \langle \textit{spaces} \rangle \langle \textit{array-content} \rangle$

$\langle \textit{dictionary} \rangle ::= \text{'<'} \langle \textit{spaces} \rangle \langle \textit{dictionary-content} \rangle \text{'>'}$

$\langle \textit{dictionary-content} \rangle ::= \langle \textit{empty} \rangle$   
 $\quad | \langle \textit{dictionary-content} \rangle \langle \textit{key-value} \rangle$

$\langle \textit{key-value} \rangle ::= \langle \textit{name} \rangle \langle \textit{spaces} \rangle \langle \textit{direct-object} \rangle \langle \textit{spaces} \rangle$

### 3 Complex tokens

The following rules describe the syntax of complex tokens such as strings, names and streams.

$\langle \textit{string} \rangle ::= \langle \textit{string-literal} \rangle | \langle \textit{string-hexa} \rangle$

$\langle \textit{string-literal} \rangle ::= \text{'('} \langle \textit{string-content} \rangle \text{'}'}$

$\langle \textit{string-content} \rangle ::= \langle \textit{empty} \rangle$   
 $\quad | \langle \textit{string-content} \rangle \langle \textit{string-char} \rangle$

$\langle \textit{string-char} \rangle ::= \langle \textit{string-regular} \rangle$   
 $\quad | \text{'\'} (\text{'n'} | \text{'r'} | \text{'t'} | \text{'b'} | \text{'f'} | \text{'('} | \text{'('} | \text{'\'} | \langle \textit{eol} \rangle)$   
 $\quad | \text{'\'} \langle \textit{four-digit} \rangle \langle \textit{octal-digit} \rangle \langle \textit{octal-digit} \rangle$   
 $\quad | \langle \textit{string-literal} \rangle$

$\langle \textit{string-hexa} \rangle ::= \text{'<'} \langle \textit{hexa-content} \rangle \text{'>'}$

$\langle \textit{hexa-content} \rangle ::= \langle \textit{spaces} \rangle$   
 $\quad | \langle \textit{hexa-content} \rangle \langle \textit{hexa-char} \rangle$

$\langle \textit{hexa-char} \rangle ::= \langle \textit{hexa-digit} \rangle \langle \textit{spaces} \rangle \langle \textit{hexa-digit} \rangle \langle \textit{spaces} \rangle$

$\langle \textit{name} \rangle ::= \text{'/'} \langle \textit{name-content} \rangle$

$\langle \textit{name-content} \rangle ::= \langle \textit{empty} \rangle$   
 $\quad | \langle \textit{name-content} \rangle \langle \textit{name-char} \rangle$

$$\begin{aligned}
\langle \textit{name-char} \rangle &::= \langle \textit{name-regular} \rangle \\
&| \text{'\#'} \langle \textit{hexa-digit} \rangle \langle \textit{hexa-digit} \rangle \\
\langle \textit{stream} \rangle &::= \text{'stream'} \langle \textit{eol} \rangle \langle \textit{stream-content} \rangle \text{'endstream'} \\
\langle \textit{stream-content} \rangle &::= \langle \textit{empty} \rangle \\
&| \langle \textit{stream-content} \rangle \langle \textit{any-char} \rangle
\end{aligned}$$

## 4 Simple tokens

The following rules describe the syntax of simple tokens such as numbers.

$$\begin{aligned}
\langle \textit{null} \rangle &::= \text{'null'} \\
\langle \textit{bool} \rangle &::= \text{'true'} | \text{'false'} \\
\langle \textit{int} \rangle &::= \langle \textit{sign} \rangle? \langle \textit{unsigned-int} \rangle \\
\langle \textit{real} \rangle &::= \langle \textit{sign} \rangle? (\langle \textit{digits} \rangle \text{'.'} \langle \textit{digits} \rangle? | \text{'.'} \langle \textit{digits} \rangle) \\
\langle \textit{unsigned-int} \rangle &::= \langle \textit{digits} \rangle \\
\langle \textit{digits} \rangle &::= \langle \textit{digit} \rangle \\
&| \langle \textit{digits} \rangle \langle \textit{digit} \rangle \\
\langle \textit{unsigned-int-10} \rangle &::= \langle \textit{unsigned-int-5} \rangle \langle \textit{unsigned-int-5} \rangle \\
\langle \textit{unsigned-int-5} \rangle &::= \langle \textit{digit} \rangle \langle \textit{digit} \rangle \langle \textit{digit} \rangle \langle \textit{digit} \rangle \langle \textit{digit} \rangle \\
\langle \textit{xref-eol} \rangle &::= \langle \textit{space} \rangle \langle \textit{cr} \rangle | \langle \textit{space} \rangle \langle \textit{lf} \rangle | \langle \textit{cr} \rangle \langle \textit{lf} \rangle \\
\langle \textit{eol} \rangle &::= \langle \textit{cr} \rangle | \langle \textit{lf} \rangle | \langle \textit{cr} \rangle \langle \textit{lf} \rangle \\
\langle \textit{any-space} \rangle &::= \langle \textit{space} \rangle | \langle \textit{cr} \rangle | \langle \textit{lf} \rangle | \langle \textit{tab} \rangle \\
\langle \textit{spaces} \rangle &::= \langle \textit{empty} \rangle \\
&| \langle \textit{spaces} \rangle \langle \textit{any-space} \rangle
\end{aligned}$$

## 5 Character set

The following rules describe character classes.

$$\begin{aligned}
\langle \textit{sign} \rangle &::= \text{'+'} | \text{'-'} \\
\langle \textit{version-digit} \rangle &::= [\text{'0'}-\text{'7'}] \\
\langle \textit{hexa-digit} \rangle &::= [\text{'0'}-\text{'9'}\text{'a'}-\text{'f'}\text{'A'}-\text{'F'}] \\
\langle \textit{octal-digit} \rangle &::= [\text{'0'}-\text{'7'}] \\
\langle \textit{four-digit} \rangle &::= [\text{'0'}-\text{'3'}]
\end{aligned}$$

$\langle digit \rangle ::= [ '0' - '9' ]$

$\langle non-ascii-char \rangle ::= [ '\text{x}80' - '\text{xFF}' ]$

$\langle space \rangle ::= '\text{x}20'$

$\langle cr \rangle ::= '\text{x}0D'$

$\langle lf \rangle ::= '\text{x}0A'$

$\langle tab \rangle ::= '\text{x}09'$

$\langle any-char \rangle ::= [ '\text{x}00' - '\text{xFF}' ]$

$\langle string-regular \rangle ::= \langle any-char \rangle - ( '(' \mid ')' \mid '\backslash' )$

$\langle regular \rangle ::= \langle any-char \rangle - ( '\text{x}00' \mid '\text{x}09' \mid '\text{x}0A' \mid '\text{x}0C' \mid '\text{x}0D' \mid '\text{x}20' \mid '(' \mid ')' \mid '<' \mid '>' \mid '[' \mid ']' \mid '\{ ' \mid '\}' \mid '/' \mid '\%' )$

$\langle name-regular \rangle ::= \langle regular \rangle - '#'$