

# ABNF in ACL2

Alessandro Coglio

Kestrel Institute, Palo Alto, California, USA  
<http://www.kestrel.edu/~coglio>

Technical Report, April 2017

**Abstract.** Augmented Backus-Naur Form (ABNF) is a standardized formal grammar notation used in several Internet syntax specifications. This paper describes (i) a formalization of the syntax and semantics of the ABNF notation and (ii) a verified parser that turns ABNF grammar text into a formal representation usable in declarative specifications of parsers of ABNF-specified languages. This work has been developed in the ACL2 theorem prover.

## 1 Problem, Contribution, and Related Work

Augmented Backus-Naur Form (ABNF) is a standardized [12,18] formal grammar notation used in several Internet syntax specifications [8,13,21,16,11,10]. Since inadequate parsing may enable security exploits such as HTTP request smuggling [13, Section 9.5], formally verified parsers of ABNF-specified languages are of interest. It is important to ensure that the formal specifications against which the parsers are verified are faithful to the ABNF grammars.

The work described in this paper contributes to this goal by providing:

1. A formalization of the syntax and semantics of the ABNF notation.
2. A verified parser that turns ABNF grammar text (e.g. the grammar of HTTP [13,8]) into a formal representation usable in declarative specifications of parsers (e.g. an HTTP parser).

Since the syntax of ABNF itself is specified in ABNF, the approach to implement and verify this ABNF grammar parser should be applicable to other parsers of ABNF-specified languages, e.g. to an HTTP parser. The approach may also provide general ideas for building a tool that automatically generates, from suitable ABNF grammars, verified parsers with independently checkable proofs.

This work has been developed in the ACL2 theorem prover [3]. The development is available [1, [books/kestrel/abnf](#)] and thoroughly documented [2, [abnf](#)].<sup>1</sup> The development also includes a growing collection of operations on ABNF grammars, e.g. to check their well-formedness and to compose them, but this paper does not discuss these operations.

The author is not aware of other formalizations of the ABNF notation. Existing work on verified parsing [19,7,17] is based on different grammar notations,

---

<sup>1</sup> Some of the development's excerpts in this paper are slightly simplified for brevity.

and thus not directly applicable to ABNF parsing without trusted translations of notation. The APG parser generator [22] generates parsers from ABNF grammars, but without correctness proofs.

## 2 Background

### 2.1 ABNF

ABNF [12,18] adds conveniences and makes slight modifications to Backus-Naur Form (BNF) [6], without going beyond context-free grammars [14].

Instead of BNF's angle bracket notation for nonterminals, ABNF uses case-insensitive names consisting of letters, digits, and dashes, e.g. `HTTP-message` and `IPv6address`. ABNF includes an angle bracket notation for prose descriptions, e.g. `<host, see [RFC3986], Section 3.2.2>`, usable as last resort [12, Section 4] in the definiens of a nonterminal.

While BNF allows arbitrary terminals, ABNF uses only natural numbers as terminals, and denotes them via: (i) binary, decimal, or hexadecimal sequences, e.g. `%b1.11.1010`, `%d1.3.10`, and `%x.1.3.a` all denote the string '1 3 10'; (ii) binary, decimal, or hexadecimal ranges, e.g. `%x30-39` denotes any string 'n' with  $48 \leq n \leq 57$  (an ASCII digit); (iii) case-sensitive ASCII strings, e.g. `%s"Ab"` denotes the string '65 98'; and (iv) case-insensitive ASCII strings, e.g. `%i"ab"`, or just `"ab"`, denotes any string among '65 66', '65 98', '97 66', and '97 98'. ABNF terminals in suitable sets may represent ASCII or Unicode characters.

ABNF allows repetition prefixes  $n*m$ , where  $n$  and  $m$  are natural numbers in decimal notation; if absent,  $n$  defaults to 1, and  $m$  defaults to infinity. For example, `1*4HEXDIG` denotes one to four HEXDIGs, `*3DIGIT` denotes up to three DIGITs, and `1*OCTET` denotes one or more OCTETs. A single  $n$  prefix abbreviates  $n*n$ , e.g. `3DIGIT` denotes three DIGITs.

Instead of BNF's `|`, ABNF uses `/` to separate alternatives. Repetition prefixes have precedence over juxtaposition, which has precedence over `/`. Round brackets group things and override the aforementioned precedence rules, e.g. `*(WSP / CRLF WSP)` denotes strings obtained by repeating, zero or more times, either (i) a WSP or (ii) a CRLF followed by a WSP. Square brackets also group things but make them optional, e.g. `[":" port]` is equivalent to `0*1(":" port)`.

Instead of BNF's `::=`, ABNF uses `=` to define nonterminals, and `=/` to incrementally add alternatives to previously defined nonterminals. For example, the rule `BIT = "0" / "1"` is equivalent to `BIT = "0"` followed by `BIT =/ "1"`.

The syntax of ABNF itself is formally specified in ABNF [12, Section 4], after the syntax and semantics of ABNF are informally specified in natural language [12, Sections 1-3]. The syntax rules of ABNF prescribe the ASCII codes allowed for white space (spaces and horizontal tabs), line endings (carriage returns followed by line feeds), and comments (semicolons to line endings).

### 2.2 ACL2

ACL2 [3] is a general-purpose theorem prover based on a quantifier-free untyped first-order logic of total functions that is an extension of a purely functional

```

(defun fact (n)
  (if (zp n)
      1
      (* n (fact (- n 1)))))

(defthm fact-above-arg
  (implies (natp n)
           (>= (fact n) n)))

(defchoose below (b) (n)
  (and (natp b)
       (< b (fact n))))

(defun-sk between (n)
  (exists (m)
    (and (natp m)
         (< (below n) m)
         (< m (fact n)))))

```

**Fig. 1.** Some simple examples of ACL2 functions and theorems.

subset of Common Lisp [15]. Predicates are functions and formulas are terms; they are false when their value is `nil`, and true when their value is `t` or anything non-`nil`.

In compliance with Lisp syntax: a function application is a parenthesized list of the function’s name followed by the arguments, e.g.  $x + 2 \times f(y)$  is written `(+ x (* 2 (f y)))`; names of constants start and end with `*`, e.g. `*limit*`; and comments extend from semicolons to line endings (like ABNF).

The user interacts with ACL2 by submitting a sequence of theorems, function definitions, etc. ACL2 attempts to prove theorems automatically, via algorithms similar to NQTHM [9], most notably simplification and induction. The user guides these proof attempts mainly by (i) proving lemmas for use by specific proof algorithms (e.g. rewrite rules for the simplifier) and (ii) supplying theorem-specific ‘hints’ (e.g. to case-split on certain conditions).

For example, the factorial function can be defined like `fact` in Fig. 1, where `zp` tests if `n` is 0 or not a natural number. Thus `fact` treats arguments that are not natural numbers as 0. ACL2 functions often handle arguments of the wrong type by explicitly or implicitly coercing them to the right type.<sup>2</sup>

To preserve consistency, recursive function definitions must be proved to terminate via a measure of the arguments that decreases in each recursive call according to a well-founded relation. For `fact`, ACL2 automatically finds a measure and proves that it decreases according to a standard well-founded relation, but sometimes the user has to supply a measure.

For example, a theorem saying that `fact` is above its argument can be introduced like `fact-above-arg` in Fig. 1, where `natp` tests if `n` is a natural number. ACL2 proves this theorem automatically,<sup>3</sup> finding and using an appropriate induction rule—the one derived from the recursive definition of `fact`, in this case.

Besides the discouraged ability to introduce arbitrary axioms, ACL2 provides consistency-preserving mechanisms to axiomatize new functions, such as indefinite description [5] functions. For example, a function constrained to be strictly below `fact` can be introduced like `below` in Fig. 1, where `b` is the variable bound by the indefinite description. This introduces the conservative axiom that, for every `n`, `(below n)` is a natural number less than `(fact n)`, if any exists—otherwise, `(below n)` is unconstrained.

<sup>2</sup> Since the logic is untyped, in ACL2 a ‘type’ is a subset of the universe of values.

<sup>3</sup> Provided that a standard arithmetic library [1, `books/arithmetic`] is loaded.

```

rulelist = 1*( rule / (*c-wsp c-nl) )
rule = rulename defined-as elements c-nl
defined-as = *c-wsp ("=" / "=/") *c-wsp
elements = alternation *c-wsp
alternation = concatenation *( *c-wsp "/" *c-wsp concatenation )
concatenation = repetition *(1*c-wsp repetition)
repetition = [repeat] element
element = rulename / group / option / char-val / num-val / prose-val
group = "(" *c-wsp alternation *c-wsp ")"
option = "[" *c-wsp alternation *c-wsp "]"
num-val = "%" (bin-val / dec-val / hex-val)
bin-val = "b" 1*BIT [ 1*("." 1*BIT) / ("-" 1*BIT) ]
dec-val = "d" 1*DIGIT [ 1*("." 1*DIGIT) / ("-" 1*DIGIT) ]
hex-val = "x" 1*HEXDIG [ 1*("." 1*HEXDIG) / ("-" 1*HEXDIG) ]
c-wsp = WSP / (c-nl WSP)
c-nl = comment / CRLF
CRLF = CR LF ; carriage return and line feed
WSP = SP / HTAB ; space or horizontal tab

```

**Fig. 2.** Some rules of the ABNF grammar of ABNF.

ACL2's Lisp-like macro mechanism provides the ability to extend the language with new constructs defined in terms of existing constructs. For instance, despite the lack of quantification in the logic, functions with top-level quantifications can be introduced. For example, the existence of a value strictly between `fact` and `below` can be packaged into a predicate like `between` in Fig. 1, where `defun-sk` is defined in terms of `defchoose` and `defun` [5].

### 3 ABNF Formalization

#### 3.1 Abstract Syntax

The formalization starts by defining an abstract syntax of ABNF, based on the ABNF rules that define the concrete syntax of ABNF.<sup>4</sup> To ease validation by inspection, this abstract syntax abstracts away only mostly lexical details (e.g. white space, comments, and defaults of repetition prefixes). ACL2's FTY macro library for introducing structured recursive types [2, `fty`] is used to define the abstract syntactic entities of ABNF—11 types in total.

For example, the concrete syntax of numeric terminal notations such as `%d1.3.10` and `%x30-39` is defined in ABNF by `num-val` in Fig. 2, where the definitions of `BIT`, `DIGIT`, and `HEXDIG` are not shown but should be obvious from the names. In ACL2, the corresponding abstract syntax is formalized by `num-val` in Fig. 3, where `fty::deftagsum` introduces a tagged sum type (disjoint union), `num-val` is the name of the type, `:direct` tags direct notations such as `%d1.3.10` whose only component `get` is a list of natural numbers (type `nat-list`, whose recognizer is `nat-listp`) such as `(1 3 10)`, and `:range` tags range notations such as `%x30-39` whose components `min` and `max` are natural numbers (type

<sup>4</sup> The meta circularity of the definition of the concrete syntax of ABNF in ABNF is broken by the human in the loop, who defines an abstract syntax of ABNF in ACL2.

```

(fty::deftagsum num-val
  (:direct ((get nat-list)))
  (:range ((min nat) (max nat))))

(fty::defprod rule
  ((name rulename)
   (incremental bool)
   (definiens alternation)))

(fty::deflist rulelist
  :elt-type rule)

(fty::deftypes alt/conc/rep/elem
  (fty::deflist alternation
   :elt-type concatenation)
  (fty::deflist concatenation
   :elt-type repetition)
  (fty::defprod repetition
   ((range repeat-range)
    (element element)))
  (fty::deftagsum element
   (:rulename ((get rulename)))
   (:group ((get alternation)))
   (:option ((get alternation)))
   (:char-val ((get char-val)))
   (:num-val ((get num-val)))
   (:prose-val ((get prose-val))))))

```

**Fig. 3.** Some excerpts of the abstract syntax of ABNF formalized in ACL2.

`nat`, whose recognizer is `natp`) such as 48 and 57. This type definition introduces: a recognizer `num-val-p` for the type; constructors `num-val-direct` and `num-val-range`; destructors `num-val-direct->get`, `num-val-range->min`, and `num-val-range->max`; and several theorems about these functions.

As another example, the concrete syntax of alternations, concatenations, repetitions, etc. is defined in ABNF by `alternation` and mutually recursive companions in Fig. 2. In ACL2, the corresponding abstract syntax is formalized by `alternation` and companions in Fig. 3, where `fty::deftypes` introduces mutually recursive types, `fty::deflist` introduces a type of lists over the element type that appears after `:elt-type`, `fty::defprod` introduces a product type similar to a `fty::deftagsum` summand, `repeat-range` is a type for repetition prefixes such as `1*` or `3*6`, `rulename` is a type for rule names, `char-val` is a type for string terminal notations such as `%s"Ab"`, and `prose-val` is a type for prose notations such as `<host, see [RFC3986], Section 3.2.2>`. These type definitions introduce recognizers, constructors, destructors, and theorems.

As a third example, the concrete syntax of grammars (i.e. lists of rules) is defined in ABNF by `rulelist` in Fig. 2. In ACL2, the corresponding abstract syntax is formalized by `rulelist` in Fig. 3, where the `incremental` component of `rule` is a boolean that says whether the rule is incremental (`=/`) or not (`=`).

## 3.2 Semantics

A grammar describes how a sequence of natural numbers (terminals) can be organized in tree structures according to the grammar's rules. Thus, the semantics of the abstract syntactic entities is formalized via matching relations with trees.

Since a single terminal notation like `%d1.3.10` or `%s"Ab"` denotes multiple natural numbers in sequence, it is convenient to use lists of natural numbers, instead of individual natural numbers, to label leaves of trees. A rule name (non-terminal) can label the root of a (sub)tree, with branches for one of the concatenations of the alternation that defines the rule name. Since a concatenation is a sequence of repetitions, and each repetition may denote multiple instances of

```

(fty::deftypes trees
  (fty::deftagsum tree
    (:leafterm ((get nat-list)))
    (:leafrule ((get rulename)))
    (:nonleaf ((rulename? maybe-rulename) (branches tree-list-list))))
  (fty::deflist tree-list :elt-type tree)
  (fty::deflist tree-list-list :elt-type tree-list))

(defun tree-match-num-val-p (tree num-val)
  (and (tree-case tree :leafterm)
    (let ((nats (tree-leafterm->get tree)))
      (num-val-case num-val
        :direct (equal nats num-val.get)
        :range (and (equal (len nats) 1)
                     (<= num-val.min (car nats))
                     (<= (car nats) num-val.max))))))

(mutual-recursion
  (defun tree-list-list-match-alternation-p (treess alt rules) ...)
  (defun tree-list-list-match-concatenation-p (treess conc rules) ...)
  (defun tree-list-match-repetition-p (trees rep rules) ...)
  (defun tree-list-match-element-p (trees elem rules) ...)
  (defun tree-match-element-p (tree elem rules)
    (element-case elem
      :rulename (tree-case tree
        :leafterm nil
        :leafrule (equal tree.get element.get)
        :nonleaf (and (equal tree.rulename? element.get)
                       (let ((alt (lookup-rulename element.get rules)))
                         (tree-list-list-match-alternation-p
                          tree.branches alt rules))))
      :group ...
      :option ...
      :char-val ...
      :num-val (tree-match-num-val-p tree element.get)
      :prose-val t))

(defun parse-treep (tree string rulename rules)
  (and (treep tree)
    (tree-match-element-p tree (element-rulename rulename) rules)
    (equal (tree->string tree) string)))

```

Fig. 4. Some excerpts of the semantics of ABNF formalized in ACL2.

its element, the branches are organized into a list of lists: the outer list matches the list of repetitions that form the concatenation, and each inner list matches the element instances of the corresponding repetition. Rule names can also label leaves, to represent the tree structure of strings that include nonterminals. Round-bracketed groups and square-bracketed options are like anonymous rules: roots of (sub)trees that match groups and options are not labeled by rule names, but have lists of lists of branches for concatenations from the alternations inside the brackets, in the same way as named rules; additionally, a square-bracketed option is allowed to have an empty list of lists of branches, to represent the absence of the option.

Formally, (lists of (lists of)) trees are recursively defined by `tree` and companions in Fig. 4, where `maybe-rulename` is a type consisting of rule names and `nil`—the latter is used for roots not labeled by rule names. A function `tree->string` collects the natural numbers and rule names at the leaves of a



```
(defchoose parse-http (tree) (string)
  (and (parse-treep tree string *http-message* *http-grammar*)
       (disambiguatep tree)))
```

**Fig. 6.** A sketch of a declarative specification of an HTTP parser.

any tree; predicates external to grammars can be used to define the meaning of specific prose notations.

The assertions in the previous paragraph are formalized by `tree-list-list-match-alternation-p` and companions in Fig. 4, where `mutual-recursion` introduces mutually recursive `defuns`, `element-case` and `tree-case` are case analysis binders analogous to `num-val-case` in `tree-match-num-val-p`, and `lookup-rulename` collects from the rules of a grammar all the alternatives that define a rule name. Termination is proved via a lexicographic measure consisting of the size of the trees followed by the size of the abstract syntactic entities.

Given the rules of a grammar and a rule name, a parse tree for a string is a tree that matches the rule name and that has the string at the leaves, as formalized by `parse-treep` in Fig. 4. This predicate can be used to write declarative specifications of parsers of ABNF-specified languages. For instance, an HTTP parser can be specified by something like `parse-http` in Fig. 6, where `*http-grammar*` is a constant of type `rulelist` (Fig. 3) representing the rules of the ABNF grammar of HTTP [13,8], `*http-request*` is a constant of type `rulename` (Fig. 3) representing the top-level rule name `HTTP-message`, and the predicate `disambiguatep` states disambiguating restrictions, since the grammar of HTTP is ambiguous. The function `parse-http` returns concrete syntax trees, because grammars do not specify abstract syntax; a practical HTTP parser can be specified as the composition of `parse-http` followed by a suitable HTTP syntax abstraction function.

### 3.3 Concrete Syntax

The concrete syntax of ABNF is formalized in ACL2 using the rules of the ABNF grammar of ABNF, but “written in abstract syntax” because the concrete syntax is not available before it is formalized. Since the FTY constructors of the abstract syntax are somewhat verbose, some specially crafted and named functions and macros are defined, and used to write abstract syntactic entities in a way that looks more like concrete syntax, easing not only their writing, but also their validation by inspection. For example, the rules `group` and `num-val` (Fig. 2) are written in abstract syntax as shown in Fig. 7.

After transcribing the 40 rules of the ABNF grammar of ABNF to this form, a constant `*abnf-grammar*` consisting of their list (analogous to `*http-grammar*` in Fig. 6) is defined. Since grammars, i.e. values of type `rulelist` (Fig. 3), are endowed with semantics (Sect. 3.2), this constant provides a formalization of the concrete syntax of ABNF in ACL2.

The link between the concrete and abstract syntax of ABNF is formalized by 51 executable ACL2 functions that abstract trees that match syntactic enti-

```

(def-rule-const *group*
  (/_ "(" (*_ *c-wsp*) *alternation* (*_ *c-wsp*) ")"))

(def-rule-const *num-val*
  (/_ "%" (!_ (/_ *bin-val*) (/_ *dec-val*) (/_ *hex-val*))))

```

**Fig. 7.** Some excerpts of the concrete syntax of ABNF formalized in ACL2.

ties of ABNF to abstract syntactic entities of ABNF. For example, a function `abstract-num-val` abstracts a tree that matches the rule name `num-val` (Fig. 2) to a value of type `num-val` (Fig. 3). This function calls other abstraction functions on its subtrees, e.g. a function `abstract-*bit` that abstracts a list of trees that matches `*BIT` to the big endian value of their bits. As another example, the top-level abstraction function `abstract-rulelist` abstracts a tree that matches the rule name `rulelist` (Fig. 2) to the corresponding value of type `rulelist` (Fig. 3)—a grammar.

## 4 ABNF Grammar Parser

When specifying a parser of an ABNF-specified language as sketched in Fig. 6, a constant like `*http-grammar*` can be built by manually transcribing the grammar, as done for `*abnf-grammar*` (Sect. 3.3). But this transcription can be performed automatically, by running (i) the grammar parser described in the rest of this section, which produces a parse tree that matches `rulelist`, followed by (ii) `abstract-rulelist` (Sect. 3.3) on the resulting parse tree.

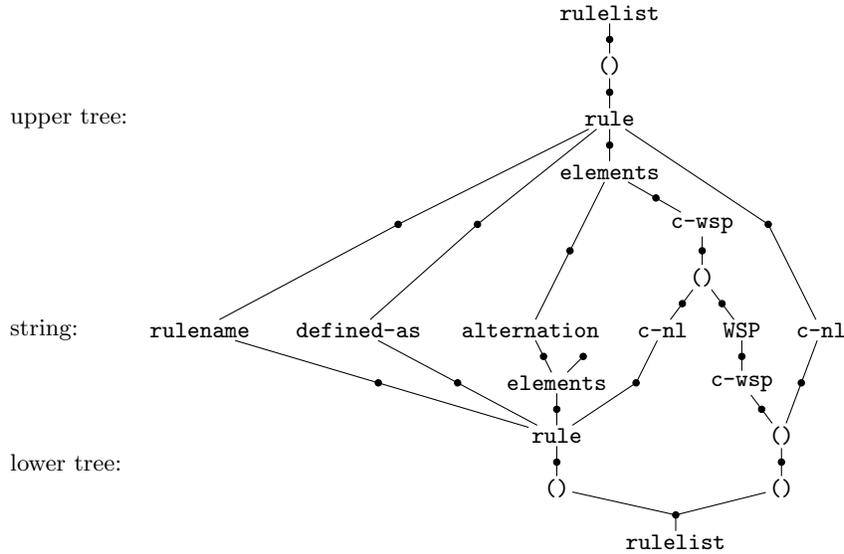
Since the grammar parser is verified as described below, and this automatic transcription process operates on the actual grammar text (e.g. copied and pasted from an Internet standard document), the resulting formal parsing specification is quite faithful to the grammar.

Running this process on the ABNF grammar of ABNF produces the same value as the manually built `*abnf-grammar*`. This provides a validation.

### 4.1 Implementation

The ABNF grammar of ABNF is ambiguous, as shown by the two different parse trees for the same string (of nonterminals, for brevity) in Fig. 8: the first `c-nl` can either end a `rule` (lower tree) or form, with the `WSP`, a `c-wsp` under `elements` (upper tree). The ambiguity only affects where certain comments, white space, and line endings go in the parse trees; it does not affect the abstract syntax, and thus the semantics, of ABNF. The parser resolves the ambiguity by always parsing as many consecutive `c-wsps` as possible, as in the upper tree in Fig. 8.

Aside from this ambiguity, the ABNF grammar of ABNF is mostly LL(1), with some LL(2) and LL(\*) parts [4,20]. The parser is implemented as a recursive descent with backtracking. Backtracking is expected to be limited in reasonable grammars—the parser runs very quickly on popular grammars [8,13,21,16,11,10].



**Fig. 8.** An example showing the ambiguity of the ABNF grammar of ABNF.

The parser consists of 85 executable ACL2 functions. There is a parsing function for each rule, and parsing functions for certain groups, options, repetitions, and terminal notations. ACL2's *Seq* macro library for stream processing [2, `seq`] is used to define these functions. Each function takes a list of natural numbers to parse as input, and, consistently with *Seq*, returns (i) an indication of success (`nil`, i.e. no error) or failure (an error message), (ii) a (list of) parse tree(s) if successful, and (iii) the remaining natural numbers in the input.

For example, the parsing function for CRLF (Fig. 2) is `parse-crlf` in Fig. 9, where `parse-cr` parses a carriage return, yielding a parse tree that is bound to `tree-cr`, then `parse-lf` parses a line feed, yielding a parse tree that is bound to `tree-lf`, and finally `return` returns (i) `nil` (success), (ii) a CRLF parse tree with the two subtrees, and (iii) the remaining input after the carriage return and line feed. If `parse-cr` or `parse-lf` fails, `parse-crlf` fails. The threading of the input and the propagation of the failures is handled by the `seq` macro behind the scenes.

As another example, the parsing function for WSP (Fig. 2) is `parse-wsp` in Fig. 9, where `parse-sp` attempts to parse a space, returning a WSP parse tree with a SP subtree if successful, and otherwise `parse-htab` attempts to parse a horizontal tab, returning a WSP parse tree with a HTAB subtree if successful. If both `parse-sp` and `parse-htab` fail, `parse-wsp` fails. The backtracking is handled by the `seq-backtrack` macro behind the scenes.

As a third example, the function for \*BIT is `parse-*bit` in Fig. 9, which uses `parse-bit` to parse as many bits as possible, eventually returning the corresponding list of BIT parse trees. Termination is proved by the decrease of the

```

(defun parse-crlf (input)
  (seq input
    (tree-cr := (parse-cr input))
    (tree-lf := (parse-lf input))
    (return (tree-nonleaf *crlf* (list (list tree-cr) (list tree-lf))))))

(defun parse-wsp (input)
  (seq-backtrack input
    ((tree := (parse-sp input))
     (return (tree-nonleaf *wsp* (list (list tree)))))
    ((tree := (parse-htab input))
     (return (tree-nonleaf *wsp* (list (list tree))))))

(defun parse-*bit (input)
  (seq-backtrack input
    ((tree := (parse-bit input))
     (trees := (parse-*bit input))
     (return (cons tree trees)))
    ((return nil))))

(defun parse-grammar (input)
  (b* ((mv error? tree? rest) (parse-rulelist input)))
  (cond (error? nil)
        (rest nil)
        (t tree?)))

```

**Fig. 9.** Some excerpts of the ABNF grammar parser in ACL2.

length of the input. This function never fails: when no bits can be parsed, the empty list `nil` of parse trees is returned.

The functions for **alternation** and mutually recursive companions (Fig. 2) are mutually recursive. Their termination is proved via a lexicographic measure consisting of the size of the input followed by an ordering of these functions—the length of the input alone is insufficient to prove termination, because, for example, `parse-alternation` calls `parse-concatenation` on the same input.

The top-level parsing function is `parse-grammar` in Fig. 9, where `b*` binds the results of `parse-rulelist` to the three variables in the triple `(mv ...)`. The function checks that there is no remaining input, returning just the parse tree if successful (or `nil`, i.e. no parse tree, if a failure occurs). There is also a wrapper function `parse-grammar-from-file` that takes a file name as input and calls `parse-grammar` on the file’s contents.

## 4.2 Verification

The correctness of the parser consists of:

- Soundness: the parser recognizes only ABNF grammars.
- Completeness: the parser recognizes all ABNF grammars (almost; see below).

The soundness theorem is `parse-treep-of-parse-grammar` in Fig. 10, where `*rulelist*` represents the rule name `rulelist` (Fig. 2). The theorem says that if `parse-grammar` succeeds (i.e. it does not return `nil`), it returns a parse tree that organizes the input into the tree structure of a grammar (i.e. a list of rules).

This soundness theorem is proved via two theorems for each of the parsing functions that return triples:

```

(defthm parse-treep-of-parse-grammar
  (implies (and (nat-listp input)
                (parse-grammar input))
            (parse-treep
             (parse-grammar input) input *rulelist* *abnf-grammar*)))

(defthm input-decomposition-of-parse-crlf
  (implies (and (nat-listp input)
                (not (mv-nth 0 (parse-crlf input))))
            (equal (append (tree->string (mv-nth 1 (parse-crlf input)))
                           (mv-nth 2 (parse-crlf input)))
                    input)))

(defthm tree-match-of-parse-crlf
  (implies (and (nat-listp input)
                (not (mv-nth 0 (parse-crlf input))))
            (tree-match-element-p (mv-nth 1 (parse-crlf input))
                                  (element-rulename *crlf*)
                                  *abnf-grammar*)))

```

**Fig. 10.** Some excerpts of the ABNF grammar parser soundness proof in ACL2.

- Input decomposition: if the function succeeds, the string at the leaves of the returned parse tree(s) consists of the natural numbers parsed from the input.
- Tree matching: if the function succeeds, the returned parse tree/trees is/are consistent with the syntactic entity that the function is intended to parse.

For example, the input decomposition theorem of `parse-crlf` is `input-decomposition-of-parse-crlf` in Fig. 10, where `mv-nth` extracts the components (zero-indexed) of the triple returned by `parse-crlf`. The theorem says that if `parse-crlf` succeeds (i.e. its first result is `nil`, not an error), joining the string at the leaves of the returned tree with the returned remaining input yields the original input.

Each input decomposition theorem is proved by expanding the parsing function and using the input decomposition theorems of the called functions as rewrite rules. For instance, in the input decomposition theorem of `parse-crlf`, expanding `parse-crlf` turns the `(append ...)` into one involving `parse-cr` and `parse-lf`, making their input decomposition theorems applicable.

As another example, the tree matching theorem of `parse-crlf` is `tree-match-of-parse-crlf` in Fig. 10. The theorem says that if `parse-crlf` succeeds (as in the input decomposition theorem), the returned parse tree matches CRLF—which `parse-crlf` is intended to parse.

Each tree matching theorem is proved by expanding the parsing function and the tree matching predicate, and using the tree matching theorems of the called functions as rewrite rules. For instance, in the tree matching theorem of `parse-crlf`, expanding `parse-crlf` and `tree-match-element-p` turns the conclusion into the fact that the subtrees match CR and LF when `parse-cr` and `parse-lf` succeed, making their tree matching theorems applicable.

The input decomposition and tree matching theorems of the recursive functions are proved by induction on their recursive definitions.

```

(defthm parse-grammar-when-tree-match
  (implies (and (treep tree)
                (tree-match-element-p
                 tree (element-rulename *rulelist*) *abnf-grammar*)
                (tree-terminatedp tree)
                (tree-rulelist-restriction-p tree))
           (equal (parse-grammar (tree->string tree)) tree)))

(defthm parse-wsp-when-tree-match
  (implies (and (treep tree)
                (nat-listp rest-input)
                (tree-match-element-p
                 tree (element-rulename *wsp*) *abnf-grammar*)
                (tree-terminatedp tree))
           (equal (parse-wsp (append (tree->string tree) rest-input))
                  (mv nil tree rest-input))))

(defthm parse-*bit-when-tree-list-match
  (implies (and (tree-listp trees)
                (nat-listp rest-input)
                (tree-list-match-repetition-p
                 trees (*_ *bit*) *abnf-grammar*)
                (tree-list-terminatedp trees)
                (mv-nth 0 (parse-bit rest-input)))
           (equal (parse-*bit (append (tree-list->string trees) rest-input))
                  (mv nil trees rest-input))))

(defthm fail-sp-when-match-htab
  (implies (and (tree-match-element-p
                tree (element-rulename *htab*) *abnf-grammar*)
                (tree-terminatedp tree))
           (mv-nth 0 (parse-sp (append (tree->string tree) rest-input)))))

(defthm constraints-from-parse-sp
  (implies (not (mv-nth 0 (parse-sp input)))
           (equal (car input) 32)))

(defthm constraints-from-tree-match-htab
  (implies (and (tree-match-element-p
                tree (element-rulename *htab*) *abnf-grammar*)
                (tree-terminatedp tree))
           (equal (car (tree->string tree)) 9)))

(defun-sk pred-alternation (input)
  (forall (tree rest-input)
    (implies
      (and (treep tree)
           (nat-listp rest-input)
           (tree-match-element-p
            tree (element-rulename *alternation*) *abnf-grammar*)
           (tree-terminatedp tree)
           ... ; 8 parsing failure hypotheses on rest-input
           (equal input (append (tree->string tree) rest-input)))
      (equal (parse-alternation (append (tree->string tree) rest-input))
             (mv nil tree rest-input)))))

(defthm parse-alternation-when-tree-match-lemma
  (pred-alternation input))

```

**Fig. 11.** Some excerpts of the ABNF grammar parser completeness proof in ACL2.

The soundness theorem `parse-treep-of-parse-grammar` is proved from the input decomposition and tree matching theorems of `parse-rulelist`, and the fact that `parse-grammar` fails if there is remaining input.

Since the ABNF grammar of ABNF is ambiguous (Fig. 8) but the parser returns a single parse tree at a time, absolute completeness is unprovable. Completeness is provable relatively to trees consistent with how the parser resolves the ambiguity. A predicate `tree-rulelist-restriction-p` formalizes these restrictions on trees: each (`*c-wsp c-nl`) subtree, except the one (if any) that starts a `rulelist`, must not start with `WSP`.

The completeness theorem is `parse-grammar-when-tree-match` in Fig. 11. The theorem says that if a terminated tree matches a `rulelist` (i.e. it is a concrete syntactic representation of a grammar) and is consistent with how the parser resolves the ambiguity, `parse-grammar` succeeds on the string at the leaves of the tree and returns the tree.

This main completeness theorem is proved via an auxiliary completeness theorem for each of the parsing functions that return triples. The formulation of these auxiliary theorems is analogous to the main one, but with additional complexities: in the conclusions, the parsing functions are applied to the string at the leaves of the tree(s) joined with some remaining input; this makes these theorems usable as rewrite rules, and enables the addition of certain critical hypotheses to these theorems.

For example, the completeness theorem of `parse-wsp` is `parse-wsp-when-tree-match` in Fig. 11. As another example, the completeness theorem of `parse-bit` is `parse-bit-when-tree-list-match` in Fig. 11. The hypothesis that `parse-bit` fails on `rest-input` is critical: without it, `parse-bit` could parse another bit from `rest-input`, and return a longer list of trees than `trees`.

Each auxiliary completeness theorem is proved by expanding the parsing function and the tree matching predicate, using the completeness theorems of the called functions as rewrite rules, and also using, as needed, certain disambiguation theorems.

The need and nature of these disambiguation theorems, in simple form, are illustrated by examining the proof of the completeness theorem of `parse-wsp`. The hypothesis that `tree` matches `WSP` expands to two cases:

1. The subtree matches the `SP` alternative of `WSP`. Then the completeness theorem of `parse-sp` applies, `parse-sp` succeeds returning the subtree, and `parse-wsp` succeeds returning `tree`.
2. The subtree matches the `HTAB` alternative of `WSP`. For the completeness theorem of `parse-htab` to apply, `parse-sp` must be shown to fail so that `parse-wsp` reduces to `parse-htab` and the proof proceeds as in the `SP` case.

The theorem saying that `parse-sp` fails on the string at the leaves of a terminated tree matching `HTAB`, is `fail-sp-when-match-htab` in Fig. 11. This theorem is proved via two theorems saying that `parse-sp` and `HTAB` induce incompatible constraints (ASCII codes) on the same value at the start of the input: the two theorems are `constraints-from-parse-sp` and `constraints-from-tree-match-htab` in Fig. 11, where the `input` in the former is instantiated with the `(append ...)` in the latter, when proving `fail-sp-when-match-htab`.

There are 26 parsing constraint theorems similar to the one for `parse-sp`, and 49 tree matching constraint theorems similar to the one for `HTAB`. There

are 87 disambiguation theorems similar to `fail-sp-when-match-htab`: they say that certain parsing functions fail when trees match certain syntactic entities, effectively showing that the parser can disambiguate all the alternatives in the ABNF grammar of ABNF, including deciding when to stop parsing unbounded repetitions. The disambiguation theorems are used to prove not only some completeness theorems, but also other disambiguation theorems. Some disambiguation theorems include critical parsing failure hypotheses similarly to the completeness theorem of `parse-*bit`. Many disambiguation theorems show incompatible constraints just on the first one or two natural numbers in the input, corresponding to LL(1) and LL(2) parts of the grammar. But for LL(\*) parts of the grammar, the disambiguation theorems show incompatible constraints on natural numbers that follow unbounded prefixes of the input; to “go past” these prefixes in the proofs of these disambiguation theorems, certain completeness theorems are used in turn.

Since the auxiliary completeness theorems call the parsing functions not on variables but on (`append . . .`) terms, induction on the recursive parsing functions is not readily applicable [9, Chapt. 15]. For the singly recursive functions like `parse-*bit`, induction on the list of trees is used. For the mutually recursive functions like `parse-alternation`, an analogous induction on the (lists of (lists of)) trees seems unwieldy due to the number (10) of mutually recursive parsing functions. Instead, the desired completeness assertions are packaged into predicates like `pred-alternation` in Fig. 11, where the tree and remaining input are universally quantified and a new variable `input` is equated to the argument of the parsing function. Then theorems like `parse-alternation-when-tree-match-lemma` in Fig. 11 are proved by induction on the recursive parsing functions (now applicable to the variable `input`), from which the desired completeness theorems readily follow.

The main completeness theorem of `parse-grammar` is proved from the auxiliary completeness theorem of `parse-rulelist` and the fact that the absence of remaining input fulfills the parsing failure hypotheses on the remaining input.

All the theorems and proofs overviewed in this subsection are discussed in detail in the documentation [2, `abnf`]. Even the short overview above should convey that the completeness proof is considerably more laborious than the soundness proof, perhaps because the completeness proof must show that the parser can reconstruct any parse tree from its string at the leaves, while the soundness proof must show that the parser can just construct one appropriate parse tree when it succeeds. The overview should also convey that the proof methods used are rather systematic and general, and therefore should be applicable to other parsers of ABNF-specified languages, and possibly to the automatic generation of verified parsers from suitable ABNF grammars.

## Acknowledgements

This work was supported by DARPA under Contract No. FA8750-15-C-0007.

## References

1. ACL2 theorem prover and community books: Source code, <http://github.com/acl2/acl2>
2. ACL2 theorem prover and community books: User manual, <http://www.cs.utexas.edu/~moore/acl2/manuals/current/manual>
3. ACL2 theorem prover: Web page, <http://www.cs.utexas.edu/users/moore/acl2>
4. Aho, A.V., Lam, M.S., Sethi, R., Ullman, J.D.: Compilers: Principles, Techniques, and Tools. Pearson, 2nd edn. (2007)
5. Avigad, J., Zach, R.: The epsilon calculus. In: Zalta, E.N. (ed.) The Stanford Encyclopedia of Philosophy. Metaphysics Research Lab, Stanford University, summer 2016 edn. (2016), <https://plato.stanford.edu/archives/sum2016/entries/epsilon-calculus/>
6. Backus, J.W., Bauer, F.L., Green, J., Katz, C., McCarthy, J., Perlis, A.J., Rutishauser, H., Samelson, K., Vauquois, B., Wegstein, J.H., van Wijngaarden, A., Woodger, M.: Report on the algorithmic language ALGOL 60. Communications of the ACM 3(5), 299–314 (May 1960)
7. Barthwal, A., Norrish, M.: Verified, executable parsing. In: Proc. 18th European Symposium on Programming (ESOP). pp. 160–174 (2009)
8. Berners-Lee, T., Fielding, R., Masinter, L.: Uniform Resource Identifier (URI): Generic syntax. Request for Comments (RFC) 3986 (Jan 2005)
9. Boyer, R.S., Moore, J.S.: A Computational Logic. Academic Press (1979)
10. Bray, T.: The JavaScript Object Notation (JSON) data interchange format. Request for Comments (RFC) 7159 (Mar 2014)
11. Crispin, M.: Internet Message Access Protocol - version 4rev1. Request for Comments (RFC) 3501 (Mar 2003)
12. Crocker, D., Overell, P.: Augmented BNF for syntax specifications: ABNF. Request for Comments (RFC) 5234 (Jan 2008)
13. Fielding, R., Reschke, J.: Hypertext Transfer Protocol (HTTP/1.1): Message syntax and routing. Request for Comments (RFC) 7230 (Jun 2014)
14. Hopcroft, J.E., Motwani, R., Ullman, J.D.: Introduction to Automata Theory, Languages, and Computation. Pearson, 3rd edn. (2006)
15. Kaufmann, M., Moore, J.S.: A precise description of the ACL2 logic. Tech. rep., Department of Computer Sciences, University of Texas at Austin (1998)
16. Klensin, J.: Simple Mail Transfer Protocol. Request for Comments (RFC) 5321 (Oct 2008)
17. Koprowski, A., Binszok, H.: TRX: A formally verified parser interpreter. Logical Methods in Computer Science 7(2), 1–26 (2011)
18. Kyzivat, P.: Case-sensitive string support in ABNF. Request for Comments (RFC) 7405 (Dec 2014)
19. Nipkow, T.: Verified lexical analysis. In: Proc. 11th International Conference on Theorem Proving in Higher-Order Logics (TPHOL). pp. 1–15 (1998)
20. Parr, T., Fisher, K.: *LL(\*)*: The foundation of the ANTLR parser generator. In: Proc. 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI). pp. 425–436 (2011)
21. Resnick, P.: Internet Message Format. Request for Comments (RFC) 5322 (Oct 2008)
22. Thomas, L.D.: APG: ABNF Parser Generator, <http://www.coasttocoastresearch.com>