

Independent Submission
Internet-Draft
Intended status: Informational
Expires: 11 July 2026

M. Caldas
Independent
7 January 2026

CSV++ (CSV Plus Plus): Extension to RFC 4180 for Hierarchical Data
draft-mscaldas-csvpp-00

Abstract

This document specifies CSV++ (CSV Plus Plus), an extension to the Comma-Separated Values (CSV) format defined in RFC 4180. CSV++ adds support for repeating fields (one-to-many relationships) and hierarchical component structures while maintaining backward compatibility with standard CSV parsers. The extension uses declarative syntax in column headers to define array fields and nested structures, enabling representation of complex real-world data while preserving the simplicity and human-readability of CSV.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <https://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on 11 July 2026.

Copyright Notice

Copyright (c) 2026 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

| | |
|---|----|
| 1. Introduction | 3 |
| 1.1. Motivation | 3 |
| 1.2. Design Principles | 3 |
| 1.3. Requirements Language | 3 |
| 2. Conformance with RFC 4180 | 4 |
| 3. Field Separator Detection | 4 |
| 4. Array Fields (Repetitions) | 4 |
| 4.1. Syntax | 4 |
| 4.2. Examples | 5 |
| 4.3. Empty Values | 5 |
| 4.4. Escaping | 5 |
| 5. Structured Fields (Components) | 6 |
| 5.1. Syntax | 6 |
| 5.2. Examples | 6 |
| 6. Nested Structures | 6 |
| 6.1. Recursive Composition | 6 |
| 6.2. Examples | 7 |
| 6.3. Delimiter Selection Guidelines | 7 |
| 7. Parsing | 7 |
| 8. Implementation Considerations | 7 |
| 8.1. Validation | 8 |
| 8.2. Limits | 8 |
| 9. MIME Type and File Extension | 8 |
| 9.1. MIME Type | 8 |
| 9.2. File Extensions | 8 |
| 10. Security Considerations | 8 |
| 10.1. Delimiter Injection | 8 |
| 10.2. Complexity Attacks | 9 |
| 10.3. Encoding Issues | 9 |
| 11. IANA Considerations | 9 |
| 12. References | 9 |
| 12.1. Normative References | 9 |
| 12.2. Informative References | 9 |
| Appendix A. Grammar (ABNF) | 10 |
| Appendix B. Complete Examples | 10 |
| Acknowledgments | 10 |
| Author's Address | 10 |

1. Introduction

CSV++ extends the CSV format defined in [RFC4180] to support repeating fields (one-to-many relationships) and hierarchical component structures while maintaining backward compatibility with standard CSV parsers.

1.1. Motivation

Traditional CSV files represent flat, tabular data. However, real-world data often contains:

- * Repeated values (e.g., multiple phone numbers for one person)
- * Structured components (e.g., addresses with street, city, state, zip)
- * Nested hierarchies (e.g., addresses with multiple address lines)

CSV++ addresses these needs while keeping the simplicity and human-readability of CSV with a straightforward syntax.

1.2. Design Principles

1. ***Backward Compatibility:** Standard CSV parsers can read CSV++ files (though they won't interpret the enhanced structure)
2. ***Self-Documenting:** Structure is defined in column headers
3. ***Human Readable:** Data remains readable without special tools
4. ***Explicit Over Implicit:** Delimiters are declared, not assumed
5. ***Recursively Composable:** Structures can nest to any depth, though practical implementations SHOULD limit nesting to 3-4 levels for readability

1.3. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

2. Conformance with RFC 4180

CSV++ files MUST conform to [RFC4180] with these specifications:

- * Fields are separated by a delimiter (comma by default)
- * Records are separated by line breaks (CRLF or LF)
- * Fields containing special characters MUST be enclosed in double-quotes
- * Double-quotes within quoted fields MUST be escaped by doubling: ""
- * First record MAY be a header record per RFC 4180. However, CSV++ array and structure features REQUIRE headers to declare field types
- * MIME type: text/csv

3. Field Separator Detection

The field separator character is detected using the same rules as [RFC4180]. Parsers SHOULD auto-detect the field separator by:

1. Scanning the first line (header row)
2. Tracking bracket depth: [] and ()
3. Identifying characters that appear outside brackets (depth = 0)
4. Selecting the most common such character as the field separator
5. Common candidates: , (comma), \t (tab), | (pipe), ; (semicolon)

The comma (,) is the conventional field separator for CSV++ files.

4. Array Fields (Repetitions)

4.1. Syntax

A field containing repeated values is declared in the header using square brackets:

```
column_name[delimiter]  
column_name[ ]
```

Where:

- * `column_name` - The name of the field
- * `[delimiter]` - Optional: The character used to separate repeated values
- * `[]` - Empty brackets use the default array delimiter

Delimiter Resolution:

1. If delimiter is specified: `phone[]` uses `|`
2. If empty brackets: `phone[]` uses the tilde (`~`) as default delimiter

The tilde (`~`) is recommended as the default array delimiter to avoid conflicts with common data characters and the field separator.

4.2. Examples

```
id,name,phone[|],email[;]
1,John,555-1234|555-5678|555-9012,john@work.com;john@home.com
2,Jane,555-4444,jane@company.com
```

Figure 1: Arrays with Explicit Delimiters

```
id,name,phone[],email[]
1,John,555-1234~555-5678~555-9012,john@work.com~john@home.com
2,Jane,555-4444,jane@company.com
```

Figure 2: Arrays with Default Delimiters

4.3. Empty Values

Empty values in repetitions are represented by consecutive delimiters:

```
id,tags[|]
1,urgent||priority
```

Figure 3

This represents three tags: "urgent", "" (empty), "priority"

4.4. Escaping

If the repetition delimiter appears in the data, the entire field MUST be quoted per [RFC4180]:

```
id,notes[]
1,"First note|with|pipes|Second note contains | character"
```

Figure 4

5. Structured Fields (Components)

5.1. Syntax

A field containing structured components is declared using parentheses:

```
column_name[repetition_delim]component_delim(
    comp1 component_delim comp2 ...)
column_name[]component_delim(comp1 component_delim comp2 ...)
column_name[(comp1 component_delim comp2 ...)]
column_name(comp1 component_delim comp2 ...)
```

Component Delimiter Resolution:

1. If specified before (: address^(...) uses ^
2. If omitted: address(...) uses the caret (^) as default delimiter

The caret (^) is recommended as the default component delimiter to avoid conflicts with common data characters.

5.2. Examples

```
id,name,geo^(lat^lon)
1,Location A,34.0522^-118.2437
2,Location B,40.7128^-74.0060
```

Figure 5: Simple Structure

```
id,name,address[~]^(street^city^state^zip)
1,John,123 Main St^Los Angeles^CA^90210~456 Oak Ave^New York^NY^10001
2,Jane,789 Pine St^Boston^MA^02101
```

Figure 6: Repeated Structures

6. Nested Structures

6.1. Recursive Composition

Structures can nest arbitrarily deep. Component names can themselves be arrays or structures. Within component names in (...), array and structure syntax applies recursively.

6.2. Examples

```
id,name,address[~]^(type^lines[;]^city^state^zip)
1,John,home^123 Main;Apt 4^LA^CA^90210~work^456 Oak^NY^NY^10001
```

Figure 7: Array Within Structure

```
id,location^(name^coords:(lat:lon))
1,Office^34.05:-118.24
2,Home^40.71:-74.00
```

Figure 8: Structure Within Structure

6.3. Delimiter Selection Guidelines

To maintain readability and parseability:

1. ***REQUIRED:*** Use different delimiters at each nesting level. Nested structures **MUST** use different component delimiters than their parent
2. Use visually distinct delimiters at each level
3. ***Recommended progression:*** ~ -> ^ -> ; -> :
4. Avoid using the field separator as a component delimiter
5. Document delimiter choices for complex schemas
6. ***Recommendation:*** Limit nesting to 3-4 levels maximum

7. Parsing

CSV++ parsers process files in two phases:

1. ***Header Parsing:*** Parse column headers to identify field types (simple, array, or structured) and extract delimiter information
2. ***Data Parsing:*** For each data row, split fields according to their declared type, respecting [RFC4180] quoting rules for nested delimiters

The ABNF grammar in Appendix A provides a formal specification. Implementations **MUST** handle arbitrary nesting depth up to their documented limits.

8. Implementation Considerations

8.1. Validation

Implementations SHOULD validate:

- * Matching number of components across repeated structures
- * Proper bracket nesting in headers
- * Delimiter conflicts (same delimiter at multiple levels)
- * MUST reject: Nested structures using the same component delimiter as their parent
- * Reasonable nesting depth (recommend warning beyond 3-4 levels)

8.2. Limits

Implementations MAY impose reasonable limits on:

- * Nesting depth (recommended minimum: 10 levels)
- * Number of components per structure (recommended minimum: 100)
- * Number of repetitions per array (recommended minimum: 1000)

9. MIME Type and File Extension

9.1. MIME Type

CSV++ files use the text/csv media type defined in [RFC4180].

9.2. File Extensions

- * .csv - Standard extension (recommended for compatibility)
- * .csvpp - MAY be used to explicitly indicate CSV++ format
- * .csvplus - Alternative explicit extension

10. Security Considerations

10.1. Delimiter Injection

Malicious data could attempt to inject delimiters to break parsing. Implementations MUST respect [RFC4180] quoting. Quoted fields MUST be parsed as literal values. Delimiters inside quotes MUST NOT be interpreted as separators.

10.2. Complexity Attacks

Deeply nested or highly repetitive structures could cause excessive memory consumption or CPU exhaustion during parsing.

Mitigations:

- * Implement depth limits
- * Implement size limits
- * Use streaming parsers for large files
- * Validate headers before processing data

10.3. Encoding Issues

Files SHOULD use UTF-8 encoding. Implementations SHOULD detect and handle encoding issues. BOM (Byte Order Mark) MAY be present.

11. IANA Considerations

This document has no IANA actions.

CSV++ files use the text/csv media type defined in [RFC4180]. The format is fully backward compatible with standard CSV parsers.

12. References

12.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC4180] Shafranovich, Y., "Common Format and MIME Type for Comma-Separated Values (CSV) Files", RFC 4180, DOI 10.17487/RFC4180, October 2005, <<https://www.rfc-editor.org/info/rfc4180>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

12.2. Informative References

Appendix A. Grammar (ABNF)

```

csvpp-file      = header-row data-rows

header-row      = field *(field-sep field) CRLF
data-rows       = *(data-row CRLF)
data-row        = value *(field-sep value)

field           = simple-field / array-field /
                 struct-field / array-struct-field
simple-field      = name
array-field      = name "[" [delimiter] "]"
struct-field     = name [component-delim] "(" component-list ")"
array-struct-field = name "[" [delimiter] "]"
                 [component-delim] "(" component-list ")"

component-list   = component *(component-delim component)
component        = simple-field / array-field /
                 struct-field / array-struct-field

name            = 1*field-char
field-char      = ALPHA / DIGIT / "_" / "-"
delimiter       = CHAR
component-delim = CHAR

value           = quoted-value / unquoted-value
quoted-value    = DQUOTE *(textdata / escaped-quote) DQUOTE
unquoted-value  = *textdata
escaped-quote   = DQUOTE DQUOTE
textdata       = <any character except DQUOTE, CRLF, or field-sep>

```

Appendix B. Complete Examples

```

id,cust,items[~]^(sku^name^qty^price^opts[;]:(k:v))
1,Alice,S1^Shirt^2^20^sz:M;col:blu~S2^Pant^1^50^sz:32

```

Figure 9: E-commerce Order

Acknowledgments

This specification was inspired by the HL7 Version 2.x delimiter hierarchy and the need for a simple, human-readable format for hierarchical data that maintains compatibility with existing CSV tools.

Author's Address

Marcelo Caldas
Independent
Roswell, Georgia
United States of America
Email: mscaldas@gmail.com