Network Working Group Internet-Draft Updates: <u>4180</u> (if approved) Intended status: Standards Track Expires: November 4, 2013 M. Hausenblas DERI, NUI Galway E. Wilde EMC Corporation J. Tennison Open Data Institute May 3, 2013

## URI Fragment Identifiers for the text/csv Media Type draft-hausenblas-csv-fragment-03

#### Abstract

This memo defines URI fragment identifiers for text/csv MIME entities. These fragment identifiers make it possible to refer to parts of a text/csv MIME entity, identified by row, column, or cell. Fragment identification can use single items, or ranges.

### Note to Readers

This draft should be discussed on the apps-discuss mailing list  $[\underline{11}]$ .

Online access to all versions and files is available on github  $[\underline{12}]$ .

Status of this Memo

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

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 <u>http://datatracker.ietf.org/drafts/current/</u>.

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 November 4, 2013.

#### Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal

Hausenblas, et al. Expires November 4, 2013 [Page 1]

# Provisions Relating to IETF Documents

(<u>http://trustee.ietf.org/license-info</u>) 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 Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

## Table of Contents

$\underline{1}.  \text{Introduction}  .  .  .  .  .  .  .  .  .  $				
<u>1.1</u> . What is text/csv?	<u>3</u>			
<pre>1.2. Why text/csv Fragment Identifiers?</pre>	<u>3</u>			
<u>1.2.1</u> . Motivation	<u>3</u>			
<u>1.2.2</u> . Use Cases	<u>4</u>			
<u>1.3</u> . Incremental Deployment	<u>4</u>			
<u>1.4</u> . Notation Used in this Memo	<u>4</u>			
$\underline{2}$ . Fragment Identification Methods				
<u>2.1</u> . Row-based selection	<u>5</u>			
<u>2.2</u> . Column-based selection	<u>5</u>			
<pre>2.3. Cell-based selection</pre>	<u>6</u>			
<u>2.4</u> . Multi-Selections	<u>6</u>			
<u>3</u> . Fragment Identification Syntax	7			
<u>4</u> . Fragment Identifier Processing	7			
<u>4.1</u> . Syntax Errors in Fragment Identifiers	7			
<u>4.2</u> . Semantics of Fragment Identifiers	7			
5. IANA Considerations	<u>8</u>			
<u>6</u> . Security Considerations	<u>8</u>			
$\underline{7}$ . Implementation Status	<u>9</u>			
<u>8</u> . Change Log	<u>9</u>			
<u>8.1</u> . From -02 to -03	<u>9</u>			
<u>8.2</u> . From -01 to -02	<u>9</u>			
<u>8.3</u> . From -00 to -01	<u>10</u>			
<u>9</u> . References	<u>10</u>			
<u>9.1</u> . Normative References	<u>10</u>			
<u>9.2</u> . Non-Normative References	<u>10</u>			
Appendix A. Acknowledgements	<u>11</u>			
Authors' Addresses	11			

### **<u>1</u>**. Introduction

This memo updates the text/csv media type defined in <u>RFC 4180</u> [1] by defining URI fragment identifiers for text/csv MIME entities.

This section gives an introduction to the general concepts of text/ csv MIME entities and URI fragment identifiers, and discusses the need for fragment identifiers for text/csv and deployment issues. <u>Section 2</u> discusses the principles and methods on which this memo is based. <u>Section 3</u> defines the syntax, and <u>Section 4</u> discusses processing of text/csv fragment identifiers.

### **<u>1.1</u>**. What is text/csv?

Internet Media Types (often referred to as "MIME types") as defined in <u>RFC 2045</u> [2] and <u>RFC 2046</u> [3] are used to identify different types and sub-types of media. The text/csv media type is defined in <u>RFC</u> <u>4180</u> [1], using US-ASCII [8] as the default character encoding (other character encodings can be used as well). Apart from a media type parameter for specifying the character encoding ("charset"), there is a second media type parameter ("header") that indicates whether there is a header row in the CSV document or not.

## **<u>1.2</u>**. Why text/csv Fragment Identifiers?

URIS are the identification mechanism for resources on the Web. The URI syntax specified in <u>RFC 3986</u> [4] optionally includes a so-called "fragment identifier", separated by a number sign ("#"). The fragment identifier consists of additional reference information to be interpreted by the client after the retrieval action has been successfully completed. The semantics of a fragment identifier is a property of the media type resulting from a retrieval action, regardless of the URI scheme used in the URI reference. Therefore, the format and interpretation of fragment identifiers is dependent on the media type of the retrieval result.

## **<u>1.2.1</u>**. Motivation

Similar to the motivation in <u>RFC 5147</u> [9], which defines fragment identifiers for plain text files, referring to specific parts of a resource can be very useful, because it enables users and applications to create more specific references. Users can create references to the part they really are interested in or want to talk about, rather than always pointing to a complete resource. Even though it is suggested that fragment identification methods are specified in a media type's registration (see [10]), many media types do not have fragment identification methods associated with them.

Fragment identifiers are only useful if supported by the client, because they are only interpreted by the client. Therefore, a new fragment identification method will require some time to be adopted by clients, and older clients will not support it. However, because the URI still works even if the fragment identifier is not supported (the resource is retrieved, but the fragment identifier is not interpreted), rapid adoption is not highly critical to ensure the success of a new fragment identification method.

## 1.2.2. Use Cases

Fragment identifiers for text/csv as defined in this memo make it possible to refer to specific parts of a text/csv MIME entity. Use cases include, but are not limited to, selecting a part for visual rendering, stream processing, making assertions about a certain value (provenance, confidence, comments, etc.), or data integration.

#### **<u>1.3</u>**. Incremental Deployment

As long as text/csv fragment identifiers are not supported universally, it is important to consider the implications of incremental deployment. Clients (for example, Web browsers) not supporting the text/csv fragment identifier described in this memo will work with URI references to text/csv MIME entities, but they will fail to understand the identification of the sub-resource specified by the fragment identifier, and thus will behave as if the complete resource was referenced. This is a reasonable fallback behavior, and in general users should take into account the possibility that a program interpreting a given URI will fail to interpret the fragment identifier part. Since fragment identifier evaluation is local to the client (and happens after retrieving the MIME entity), there is no reliable way for a server to determine whether a requesting client is using a URI containing a fragment identifier.

## **<u>1.4</u>**. Notation Used in this Memo

The capitalized key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in <u>RFC</u> <u>2119</u> [<u>5</u>].

#### **2**. Fragment Identification Methods

This memo specifies fragment identification using following methods: "row" for row selections, "col" for columns selections, and "cell" for cell selections.

Throughout the sections below, the following example table in CSV (having 7 rows, including one header row, and 3 columns) is used: date,temperature,place 2011-01-01,1,Galway 2011-01-02,-1,Galway 2011-01-03,0,Galway 2011-01-01,6,Berkeley 2011-01-02,8,Berkeley 2011-01-03,5,Berkeley

### 2.1. Row-based selection

To select a specific record, the "row" scheme followed by a single number is used (the first row has the position 0). http://example.com/data.csv#row=3

The above CSV fragment identifies the fourth row: 2011-01-03,0,Galway

Fragments can also select ranges of rows: http://example.com/data.csv#row=4-6

The above CSV fragment identifies three consecutive rows: 2011-01-01,6,Berkeley 2011-01-02,8,Berkeley 2011-01-03,5,Berkeley

The value "\*" can be used to indicate the last row, so the previous URI is equivalent to: http://example.com/data.csv#row=4-\*

#### **<u>2.2</u>**. Column-based selection

To select values from a certain column, the "col" scheme, followed by a position: http://example.com/data.csv#col=1

The above CSV fragment addresses the second column, identifying the column: temperature 1 -1 0 6 8 5

The "col" scheme can also be used to identify ranges of columns:

http://example.com/data.csv#col=0-1

The above CSV fragment addresses the first and second column: date,temperature 2011-01-01,1 2011-01-02,-1 2011-01-03,0 2011-01-01,6 2011-01-02,8 2011-01-03,5

The value "\*" can be used to indicate the last column.

## 2.3. Cell-based selection

To select particular fields, use the "cell" scheme, followed by a row number, a comma, and a column number. http://example.com/data.csv#cell=3,0

The above CSV fragment addresses the field in the first column within the fourth row, yielding: 2011-01-03

It is also possible to select cell-based fragments that have more than just one cell, in which case the cell selection uses the same range syntax as for row and column range selections. For these selections, the syntax uses the upper-lefthand cell as the starting point of the selection, followed by a minus sign, and then the lowerrighthand cell as the end point of the selection. http://example.com/data.csv#cell=3,0-5,1

The above CSV fragment selects a region that starts at the fourth row and the first column, and ends at the sixth row and the second column: 2011-01-03,0 2011-01-01,6 2011-01-02,8

## 2.4. Multi-Selections

Row, column, and cell selections can make more than one selection, in which case the individual selections are separated by semicolons. In these cases, the resulting fragment may be a disjoint fragment, such as the selection "#row=2;5" for the example CSV, which would select the third and the sixth row. It is up to the user agent to decide how to handle disjoint fragments, but since they are allowed, user agents should be prepared to handle disjoint fragments.

### **<u>3</u>**. Fragment Identification Syntax

The syntax for the text/csv fragment identifiers is as follows.

The following syntax definition uses ABNF as defined in  $\frac{\text{RFC} 4234}{\text{IGI}}$ , including the rule DIGIT.

NOTE: In the descriptions that follow, specified text values MUST be used exactly as given, using exactly the indicated lower-case letters. In this respect, the ABNF usage differs from [6].

csv-fragment	=	rowsel / colsel / cellsel
rowsel	=	"row=" singlespec 0*( ";" singlespec)
colsel	=	"col=" singlespec 0*( ";" singlespec)
cellsel	=	"cell=" cellspec 0*( ";" cellspec)
singlespec	=	position [ "-" position ]
cellspec	=	<pre>cellrow "," cellcol [ "-" cellrow "," cellcol ]</pre>
cellrow	=	position
cellcol	=	position
position	=	number / "*"
number	=	1*( DIGIT )

#### **<u>4</u>**. Fragment Identifier Processing

Applications implementing support for the mechanism described in this memo MUST behave as described in the following sections.

## 4.1. Syntax Errors in Fragment Identifiers

If a fragment identifier contains a syntax error (i.e., does not conform to the syntax specified in <u>Section 3</u>), then it MUST be ignored by clients. Clients MUST NOT make any attempt to correct or guess fragment identifiers. Syntax errors MAY be reported by clients.

### **4.2**. Semantics of Fragment Identifiers

Rows and columns in CSV are counted from zero. Positions thus refer to the rows and columns starting from position 0, which identifies the first row or column of a CSV. The special character "\*" can be used to refer to the last row or column of a CSV, thus allowing fragment identifiers to easily identify ranges that extend to the last row or column.

If single selections refer to non-existing rows or columns (i.e., beyond the size of of the CSV), they MUST be ignored.

Internet-Draft

text/csv Fragment Identifiers

If ranges extend beyond the size of the CSV (by extending to row or columns beyond the size of the CSV), they MUST be interpreted to only extend to the actual size of the CSV.

If selections of ranges of rows or columns or selections of cell ranges are specified in a way so that they select "inversely" (i.e., "#row=10-5" or "#cell=10,10-5,5"), they MUST be ignored.

Each specification of an identified region is processed independently, and ignored specifications (because of reason listed in the previous paragraphs) to not cause the whole fragment identifier to fail, they just mean that this single specification is ignored. For the example file, the fragment identifier "#row=0-1;4-3;12-15" does identify the first two rows: the second specification is an "inverse" specification and thus ignored, and the third specification targets rows beyond the actual size of the CSV and is also ignored.

The complete fragment identifier identifies all the successfully evaluated identified parts as a single identified fragment. This fragment can be disjoint because of multiple selections. Multiple selections also can result in overlapping individual parts, and it is up to the user agent how to process such a fragment, and whether the individual parts are still made accessible (i.e., visualized in visual user agents), or are presented as one unit. For example, the fragment identifier "#row=2-5;3-4" contains a second identified part that is completely contained in the first identified part. Whether a user agent maintains this selection as two parts, or simply signals that the identified fragment spans from the third to the sixth row, is up for the user agent to decide.

#### 5. IANA Considerations

Note to RFC Editor: Please change this section to read as follows after the IANA action has been completed: "IANA has added a reference to this specification in the text/csv Media Type registration."

IANA is requested to update the registration of the MIME Media type text/csv at <a href="http://www.iana.org/assignments/media-types/text/">http://www.iana.org/assignments/media-types/text/</a> with the fragment identifier defined in this memo by adding a reference to this memo (with the appropriate RFC number once it is known).

## <u>6</u>. Security Considerations

The fact that software implementing fragment identifiers for CSV and software not implementing them differs in behavior, and the fact that

different software may show documents or fragments to users in different ways, can lead to misunderstandings on the part of users. Such misunderstandings might be exploited in a way similar to spoofing or phishing.

Implementers and users of fragment identifiers for CSV text should also be aware of the security considerations in <u>RFC 3986</u> [4] and <u>RFC 3987</u> [7].

#### 7. Implementation Status

Note to RFC Editor: Please remove this section before publication.

As explained in a draft currently under development <<u>http://tools.ietf.org/html/draft-sheffer-running-code</u>>, this section contains information about implementation status, so that reviews of the draft document can take implementation reports into account as well. If you are implementing this draft, please contact this draft's authors. Any implementation status reports are intended for draft publications only; the section will be removed when the draft is published in RFC form.

#### 8. Change Log

Note to RFC Editor: Please remove this section before publication.

#### 8.1. From -02 to -03

- o Added section on "Implementation Status" (Section 7).
- o Added examples of ranges of rows and columns.
- o Corrected errors in examples.

#### 8.2. From -01 to -02

- o Removed slices ("#where:") as fragment identification method.
- Removed any special support for headers, which means that they are now treated as a regular (the first) row (if a header row is present).
- o Changed semantics and syntax to allow multiple selection of rows, columns, and cells, and to allow ranges of rows and columns.

### 8.3. From -00 to -01

- o Added cell-based selections.
- o Added Jeni Tennison as author; updated Erik Wilde's affiliation to EMC.

#### 9. References

## <u>9.1</u>. Normative References

- [1] Shafranovich, Y., "Common Format and MIME Type for Comma-Separated Values (CSV) Files", <u>RFC 4180</u>, October 2005.
- [2] Freed, N. and N. Borenstein, "Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies", <u>RFC 2045</u>, November 1996.
- [3] Freed, N. and N. Borenstein, "Multipurpose Internet Mail Extensions (MIME) Part Two: Media Types", <u>RFC 2046</u>, November 1996.
- [4] Berners-Lee, T., Fielding, R., and L. Masinter, "Uniform Resource Identifier (URI): Generic Syntax", <u>RFC 3986</u>, January 2005.
- [5] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>RFC 2119</u>, March 1997.
- [6] Crocker, D. and P. Overell, "Augmented BNF for Syntax Specifications: ABNF", <u>RFC 4234</u>, October 2005.
- [7] Duerst, M. and M. Suignard, "Internationalized Resource Identifiers (IRI)", <u>RFC 3987</u>, January 2005.

## 9.2. Non-Normative References

- [8] ANSI X3.4-1986, "Coded Character Set 7-Bit American National Standard Code for Information Interchange", STD 63, <u>RFC 3629</u>, 1992.
- [9] Wilde, E. and M. Duerst, "URI Fragment Identifiers for the text/plain Media Type", <u>RFC 5147</u>, April 2008.
- [10] Freed, N., Klensin, J., and T. Hansen, "Media Type Specifications and Registration Procedures", <u>BCP 13</u>, <u>RFC 6838</u>, January 2013.

## URIS

- [11] <<u>https://www.ietf.org/mailman/listinfo/apps-discuss</u>>
- [12] <<u>https://github.com/dret/I-D/tree/master/csv-fragment</u>>

## Appendix A. Acknowledgements

Thanks for comments and suggestions provided by Richard Cyganiak, Ian Davis, and Gannon Dick.

Authors' Addresses

Michael Hausenblas DERI, NUI Galway IDA Business Park Galway Ireland

Phone: +353-91-495730 Email: michael.hausenblas@deri.org URI: <u>http://sw-app.org/about.html</u>

Erik Wilde EMC Corporation 6801 Koll Center Parkway Pleasanton, CA 94566 U.S.A.

Phone: +1-925-6006244
Email: erik.wilde@emc.com
URI: http://dret.net/netdret/

Jeni Tennison Open Data Institute 65 Clifton Street London EC2A 4JE U.K.

Phone: +44-797-4420482 Email: jeni@jenitennison.com URI: <u>http://www.jenitennison.com/blog/</u>