

Sieve Email Filtering: Detecting Duplicate Deliveries
draft-ietf-appsawg-sieve-duplicate-00

Abstract

This document defines a new test command "duplicate" for the "Sieve" email filtering language. This test adds the ability to detect duplicate message deliveries. The main application for this new test is handling duplicate deliveries commonly caused by mailing list subscriptions or redirected mail addresses. The detection is normally performed by matching the message ID to an internal list of message IDs from previously delivered messages. For more complex applications, the "duplicate" test can also use the content of a specific header or other parts of the message.

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 <http://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 November 23, 2013.

Copyright Notice

Copyright (c) 2013 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 (<http://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 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

- [1. Introduction](#) [3](#)
- [2. Conventions Used in This Document](#) [3](#)
- [3. Test "duplicate"](#) [3](#)
- [4. Sieve Capability Strings](#) [6](#)
- [5. Examples](#) [6](#)
 - [5.1. Example 1](#) [6](#)
 - [5.2. Example 2](#) [7](#)
 - [5.3. Example 3](#) [7](#)
- [6. Security Considerations](#) [8](#)
- [7. IANA Considerations](#) [9](#)
- [8. References](#) [9](#)
 - [8.1. Normative References](#) [9](#)
 - [8.2. Informative References](#) [10](#)
- [Author's Address](#) [10](#)

1. Introduction

This is an extension to the Sieve filtering language defined by [RFC 5228](#) [[SIEVE](#)]. It adds a test to determine whether a certain message was seen before by the delivery agent in an earlier execution of the Sieve script. This can be used to detect and handle duplicate message deliveries.

Duplicate deliveries are a common side-effect of being subscribed to a mailing list. For example, if a member of the list decides to reply to both the user and the mailing list itself, the user will get one copy of the message directly and another through mailing list. Also, if someone cross-posts over several mailing lists to which the user is subscribed, the user will receive a copy from each of those lists. In another scenario, the user has several redirected mail addresses all pointing to his main mail account. If one of the user's contacts sends the message to more than one of those addresses, the user will likely receive more than a single copy. Using the "duplicate" extension, users have the means to detect and handle such duplicates, e.g. by discarding them, marking them as "seen", or putting them in a special folder.

Duplicate messages are normally detected using the Message-ID header field, which is required to be unique for each message. However, the "duplicate" test is flexible enough to use different (weaker) criteria for defining what makes a message a duplicate, for example based on the subject line or parts of the message body. Other applications of this new test command are also possible, as long as the tracked unique value is a string.

2. Conventions Used in This Document

The 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 [[KEYWORDS](#)].

Conventions for notations are as in [[SIEVE](#)] [Section 1.1](#), including use of the "Usage:" label for the definition of action and tagged arguments syntax.

3. Test "duplicate"

```
Usage: "duplicate" [":handle" <handle: string>]
           [":header" <header-name: string> /
           "uniqueid" <value: string>]
           [":seconds" <timeout: number>]
```


In its basic form, the "duplicate" test keeps track of which messages were seen before by this test during an earlier Sieve execution. Messages are identified by their message ID as contained in the Message-ID header. The "duplicate" test evaluates to "true" when the message was seen before and it evaluates to "false" when it was not.

As a side-effect, the "duplicate" test adds the message ID to an internal duplicate tracking list once the Sieve execution finishes successfully. This way, the same test will evaluate to "true" during the next Sieve execution. Implementations MUST prevent making any definitive modifications to the internal duplicate tracking list until the Sieve script execution finishes successfully. If failing script executions would add the message ID to the duplicate tracking list, all "duplicate" tests in the Sieve script would erroneously yield "true" for the next delivery attempt of the same message, which can -- depending on the action taken for a duplicate -- easily lead to discarding the message without further notice.

However, deferring the definitive modification of the tracking list to the end of a successful Sieve script execution is not without problems. It can cause a race condition when a duplicate message is delivered in parallel before the tracking list is updated. This way, a duplicate message could be missed by the "duplicate" test. More complex implementations could use a locking mechanism to prevent this problem. But, irrespective of what implementation is chosen, situations in which the "duplicate" test erroneously yields "true" MUST be prevented at all costs.

The "duplicate" test MUST only check for duplicates amongst message ID values encountered in previous executions of the Sieve script; it MUST NOT consider ID values encountered earlier in the current Sieve script execution as potential duplicates. This means that all "duplicate" tests in a Sieve script execution, including those located in scripts included using the "include" [[INCLUDE](#)] extension, MUST always yield the same result if the arguments are identical.

Implementations SHOULD limit the number of entries in the duplicate tracking list. When limiting the number of entries, implementations SHOULD discard the oldest ones first.

Also, implementations SHOULD let entries in the tracking list expire after a short period of time. The user can explicitly control the length of this expiration time by means of the ":seconds" argument, which is always specified in seconds. If the ":seconds" argument is omitted, an appropriate default value MUST be used. A default expiration time of around 7 days is deemed to be appropriate. Sites SHOULD impose a maximum limit on the expiration time. If that limit is exceeded, the maximum value MUST silently be substituted;

Bosch

Expires November 23, 2013

[Page 4]

exceeding the limit MUST NOT produce an error. If the ":seconds" argument is zero, the "duplicate" test MUST yield "false" unconditionally.

By default, the content of the message's Message-ID header field is used as the unique ID for duplicate tracking. For more complex applications, the "duplicate" test can also be used to detect duplicate deliveries based on other message text. Then, the tracked unique ID can be an arbitrary string value extracted from the message. By adding the ":header" argument with a message header field name, the content of the specified header field can be used as the tracked unique ID instead of the default Message-ID header. Alternatively, the tracked unique ID can be specified explicitly using the ":uniqueid" argument. The ":header" and ":uniqueid" arguments are mutually exclusive and specifying both for a single "duplicate" test command MUST trigger an error.

If the tracked unique ID value is extracted directly from a message header field, i.e. when the ":uniqueid" argument is not used, leading and trailing whitespace (see [Section 2.2 of RFC 5228](#) [[SIEVE](#)]) MUST first be trimmed from the value before performing the actual duplicate verification. When the ":uniqueid" argument is used, such normalization concerns are the responsibility of the user.

If the header field specified using the ":header" argument exists multiple times in the message, only the first occurrence MUST be used for duplicate tracking. If the specified header field is not present in the message, the "duplicate" test MUST yield "false" unconditionally. In that case the duplicate tracking list is left unmodified by this test, since no unique ID value is available. The same rules apply with respect to the Message-ID header field for the basic "duplicate" test without a ":header" or ":uniqueid" argument, since that header field could also be missing or occurring multiple times.

The string parameter of the ":uniqueid" argument can be composed from arbitrary text extracted from the message using the "variables" [[VARIABLES](#)] extension. To extract text from the message body, the "foreverypart" and "extracttext" [[SIEVE-MIME](#)] extensions need to be used as well. This provides the user with detailed control over what identifies a message as a duplicate.

Note that the "duplicate" test does not support either the "index" [[DATE-INDEX](#)], or "mime" [[SIEVE-MIME](#)] extensions directly, meaning that none of the ":index", ":mime:" or associated arguments are added to the "duplicate" test when these extensions are active. The ":uniqueid" argument can be used in combination with the "variables" [[VARIABLES](#)] extension to achieve the same result indirectly.

Bosch

Expires November 23, 2013

[Page 5]

The tracked unique ID value MUST be matched case-sensitively, irrespective of whether it originates from a header or is specified explicitly using the ":uniqueid" argument. To achieve case-insensitive behavior, the "set" command added by the "variables" [[VARIABLES](#)] extension can be used in combination with the ":uniqueid" argument to normalize the tracked unique ID value to upper or lower case.

Using the ":handle" argument, the duplicate test can be employed for multiple independent purposes. The message is recognized as a duplicate only when the tracked unique ID was seen before in an earlier script execution by a "duplicate" test with the same ":handle" argument.

NOTE: The necessary mechanism to track duplicate messages is very similar to the mechanism that is needed for tracking duplicate responses for the "vacation" [[VACATION](#)] action. One way to implement the necessary mechanism for the "duplicate" test is therefore to store a hash of the tracked unique ID and, if provided, the ":handle" argument.

[4.](#) Sieve Capability Strings

A Sieve implementation that defines the "duplicate" test command will advertise the capability string "duplicate".

[5.](#) Examples

[5.1.](#) Example 1

In this basic example message duplicates are detected by tracking the Message-ID header. Duplicate deliveries are stored in a special folder contained in the user's Trash folder. If the folder does not exist, it is created automatically using the "mailbox" [[MAILBOX](#)] extension. This way, the user has a chance to recover messages when necessary. Messages that are not recognized as duplicates are stored in the user's inbox as normal.

```
require ["duplicate", "fileinto", "mailbox"];

if duplicate {
    fileinto :create "Trash/Duplicate";
}
```


5.2. Example 2

This example shows a more complex use of the "duplicate" test. The user gets network alerts from a set of remote automated monitoring systems. Multiple notifications can be received about the same event from different monitoring systems. The Message-ID of these messages is different, because these are all distinct messages from different senders. To avoid being notified multiple times about the same event the user writes the following script:

```
require ["duplicate", "variables", "imap4flags",
        "fileinto"];

if header :matches "subject" "ALERT: *" {
  if duplicate :seconds 60 :uniqueid "${1}" {
    setflag "\\seen";
  }
  fileinto "Alerts";
}
```

The subjects of the notification message are structured with a predictable pattern which includes a description of the event. In the script above the "duplicate" test is used to detect duplicate alert events. The message subject is matched against a pattern and the event description is extracted using the "variables" [\[VARIABLES\]](#) extension. If a message with that event in the subject was received before, but more than a minute ago, it is not detected as a duplicate due to the specified ":seconds" argument. In the event of a duplicate, the message is marked as "seen" using the "imap4flags" [\[IMAP4FLAGS\]](#) extension. All alert messages are put into the "Alerts" mailbox irrespective of whether those messages are duplicates or not.

5.3. Example 3

This example shows how the "duplicate" test can be used to limit the frequency of notifications sent using the "enotify" [\[NOTIFY\]](#) extension. Consider the following scenario: a mail user receives XMPP notifications [\[NOTIFY-XMPP\]](#) about new mail through Sieve, but sometimes a single contact sends many messages in a short period of time. Now the user wants to prevent being notified of all of those messages. The user wants to be notified about messages from each person at most once per 30 minutes and writes the following script:


```
require ["variables", "envelope", "enotify", "duplicate"];

if envelope :matches "from" "*" { set "sender" "${1}"; }
if header :matches "subject" "*" { set "subject" "${1}"; }

if not duplicate :seconds 1800 :uniqueid "${sender}")
{
  notify :message "[SIEVE] ${sender}: ${subject}"
    "xmpp:user@im.example.com";
}
```

The example shown above uses the message envelope sender rather than the Message-ID header as the unique ID for duplicate tracking.

The example can be extended to allow multiple messages from the same sender in close succession as long as the discussed subject is different. This can be achieved as follows:

```
require ["variables", "envelope", "enotify", "duplicate"];

if envelope :matches "from" "*" { set "sender" "${1}"; }
if header :matches "subject" "*" { set "subject" "${1}"; }

# account for 'Re:' prefix
if string :comparator "i;ascii-casemap"
  :matches "${subject}" "Re:*"
{
  set "subject" "${1}";
}
if not duplicate :seconds 1800
  :uniqueid "${sender} ${subject}")
{
  notify :message "[SIEVE] ${sender}: ${subject}"
    "xmpp:user@im.example.com";
}
```

This uses a combination of the message envelope sender and the subject of the message as the unique ID for duplicate tracking.

6. Security Considerations

A flood of unique messages could cause the list of tracked message ID values to grow indefinitely. Implementations therefore SHOULD implement limits on the number and lifespan of entries in that list.

7. IANA Considerations

The following template specifies the IANA registration of the Sieve extension specified in this document:

To: iana@iana.org
Subject: Registration of new Sieve extension

Capability name: duplicate
Description: Adds test 'duplicate' that can be used to test whether a particular message is a duplicate, i.e. whether a copy of it was seen before by the delivery agent that is executing the Sieve script.
RFC number: this RFC
Contact address: Sieve mailing list <sieve@ietf.org>

This information should be added to the list of sieve extensions given on <http://www.iana.org/assignments/sieve-extensions>.

8. References

8.1. Normative References

[DATE-INDEX]

Freed, N., "Sieve Email Filtering: Date and Index Extensions", [RFC 5260](#), July 2008.

[INCLUDE] Daboo, C. and A. Stone, "Sieve Email Filtering: Include Extension", [RFC 6609](#), May 2012.

[KEYWORDS]

Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

[SIEVE] Guenther, P. and T. Showalter, "Sieve: An Email Filtering Language", [RFC 5228](#), January 2008.

[SIEVE-MIME]

Hansen, T. and C. Daboo, "Sieve Email Filtering: MIME Part Tests, Iteration, Extraction, Replacement, and Enclosure", [RFC 5703](#), October 2009.

[VARIABLES]

Homme, K., "Sieve Email Filtering: Variables Extension", [RFC 5229](#), January 2008.

8.2. Informative References

[IMAP4FLAGS]

Melnikov, A., "Sieve Email Filtering: Imap4flags Extension", [RFC 5232](#), January 2008.

[MAILBOX]

Melnikov, A., "The Sieve Mail-Filtering Language -- Extensions for Checking Mailbox Status and Accessing Mailbox Metadata", [RFC 5490](#), March 2009.

[NOTIFY]

Melnikov, A., Leiba, B., Segmuller, W., and T. Martin, "Sieve Email Filtering: Extension for Notifications", [RFC 5435](#), January 2009.

[NOTIFY-XMPP]

Saint-Andre, P. and A. Melnikov, "Sieve Notification Mechanism: Extensible Messaging and Presence Protocol (XMPP)", [RFC 5437](#), January 2009.

[VACATION]

Showalter, T. and N. Freed, "Sieve Email Filtering: Vacation Extension", [RFC 5230](#), January 2008.

Author's Address

Stephan Bosch
Enschede
NL

Email: stephan@rename-it.nl

