

The HTTP Distribution and Replication Protocol

[Submitted to W3C August 25, 1997](#)

The latest version of this document:

<http://www.w3.org/TR/NOTE-drp>

This version of this document:

<http://www.w3.org/TR/NOTE-drp-19970825>

Editors:

[Arthur van Hoff](#), [Marimba Inc.](#)

[John Giannandrea](#), [Netscape Inc.](#)

[Mark Hapner](#), [Sun Microsystems Inc.](#)

[Steve Carter](#), [Novell Inc.](#)

[Milo Medin](#), [At Home Corp.](#)

Status of this Document

This document is a NOTE made available by the W3 Consortium for discussion only. This indicates no endorsement of its content, nor that the Consortium has, is, or will be allocating any resources to the issues addressed by the NOTE.

This document is a submission to W3C from Marimba, Netscape, Sun Microsystems, Novell, and @Home. Please see [Acknowledged Submissions to W3C](#) regarding its disposition.

Abstract

This document provides a specification of a protocol for the efficient replication of data over HTTP.

Table of Contents

1. Introduction
2. The HTTP Distribution and Replication Protocol
 - 2.1 Content Identifiers
 - 2.2 Index Format
 - 2.3 Index Retrieval
 - 2.4 Content Based Addressing
 - 2.5 Differential Downloads
 - 2.6 HTTP Proxy Caching

MICROSOFT
Exhibit 1005

1 2.7 Backward Compatibility

2
3 3. Conclusion

4 3.1 References

5
6 **Appendices**

7 A. Index DTD

8
9 **1. Introduction**

10
11 This document provides a complete description of the HTTP Distribution and Replication Protocol
12 (DRP). The goal of the DRP protocol is to significantly improve the efficiency and reliability of data
13 distribution over HTTP. Here is a more detailed list of goals:

- 14
 - 15 • Provide a widely applicable mechanism for the efficient replication of data and content.
 - 16 • Improve the efficiency and reliability of caching servers and proxies.
 - 17 • Significantly improve the efficiency of subscription-based data and content distribution.
 - 18 • Improved reliability and versioning for the distribution of mission-critical applications and data.
 - 19 • Create an interoperability standard for replicating content between clients and servers from
20 different vendors.
 - 21 • Keep the protocol simple to avoid incompatibilities.
 - 22 • Provide functionality that is complementary to existing standards.
 - 23 • Provide functionality that is backward compatible with existing HTTP servers and proxies.
 - 24 • Provide functionality that can be deployed anywhere where HTTP is available today.

25 The HTTP Distribution and Replication Protocol was designed to efficiently replicate a hierarchical set
26 of files to a large number of clients. No assumption is made about the content or type of the files; they
27 are simply files in some hierarchical organization.

28 After the initial download a client can keep the data up-to-date using the DRP protocol. Using DRP the
29 client can download only the data that has changed since the last time it checked. Downloading only
30 the differences is much more efficient because data typically evolves slowly over time, and because
31 changes are usually restricted to only a subset of the data. One of the goals of the DRP protocol is to
32 avoid downloading the same data more than once.

33 The data that is distributed can consist of more than just HTML pages; it can consist of any kind of
34 code or content. The DRP protocol provides strong guarantees about data versioning. This is a
35 requirement for distributing mission critical applications, because getting the correct version of each
36 component is a necessary to ensure a complete application will work correctly. Correct file versioning
37 using existing HTTP requests is problematic because there is currently no reliable mechanism for
38 identifying a particular version of a file.

39 The DRP protocol uses content identifiers to automatically share resources that are requested more
40 than once. This eliminates redundant transfers of commonly used resources. The content identifiers
41 used in the DRP protocol are based on widely accepted checksum technology.

42 The DRP protocol uses a data structure called an *index*, which is currently specified using the
43 eXtensible Markup Language (XML). Because the index describes meta data, we anticipate using the
44 Resource Description Framework (RDF), in a future version of the DRP protocol specification. XML
45 is used in the interim because the RDF standard was not finalized at the time of writing.

46 The DRP protocol relies on existing HTTP/1.1 functionality to achieve better performance when

1 replicating files, and it is backward compatible with existing HTTP/1.0 and HTTP/1.1 proxies so that it
2 can be deployed immediately. To further improve the download efficiency and caching behavior of
3 HTTP, the proposal introduces new HTTP header information that may become part of the core HTTP
4 protocol specification in the future. Because the DRP protocol is layered on HTTP, it will benefit from
5 ongoing efforts in the HTTP-NG working group.

6 The DRP protocol is based on technology which was originally developed, implemented, and deployed
7 by [Marimba Inc.](#) Although it was originally designed for software distribution, it is widely applicable
8 in many different areas, and it enables the large scale automatic distribution of software, data, and
9 content to many different clients.

10 2. The HTTP Distribution and Replication Protocol

11 2.1 Content Identifiers

12 The DRP protocol uses content identifiers to identify individual pieces of content. A content identifier
13 is a token that can be used to uniquely identify any piece of data or content. It can also be used to
14 determine whether two pieces of content are identical with great accuracy.

15 A content identifier consists of one or more Uniform Resource Identifiers ([URI](#)) separated by commas.
16 The URIs are combined together into a single content identifier and form a globally unique identifier.
17 Here is the syntax of a content-identifier:

```
18 content-identifier = URI ( "," URI )*
```

19 Typically a checksum algorithm is used to generate a content identifier for a piece of content. One of
20 the checksum algorithms which can be used is the Message Digest algorithm from RSA. The MD5
21 algorithm is a well-known algorithm for computing a 128-bit checksum for any file or object. See
22 <http://www.rsa.com/pub/rfc1321.txt> for details on the implementation of the MD5 algorithm.

23 The likelihood of two different files producing the same MD5 checksum is very small (about 1 in
24 2^{64}), and as such, the MD5 checksum of a file can be used to construct a reliable content identifier
25 that is very likely to uniquely identify the file. The reverse is also true. If two files have the same MD5
26 checksum, it is very likely that the files are identical.

27 It is possible to define a Universal Resource Name ([URN](#)), which is a kind of URI, for a 128-bit MD5
28 checksum:

```
29 MD5-URN = "urn:md5:" base64-number
```

30 Another checksum algorithm used in the DRP protocol is the SHA algorithm from the National
31 Institute of Standards and Technology. It is similar to the MD5 algorithm but its checksums are 160
32 bits long, and have different properties. See <http://csrc.nist.gov/fips/fip180-1.txt> for details on the
33 implementation of the SHA algorithm.

34 Here is how a URN is created using a 160-bit SHA checksum:

```
SHA-URN = "urn:sha:" base64-number
```

Both MD5 and SHA based URIs use base64 encoding to encode the checksum of the content. The base64 algorithm encodes each 3 bytes of the checksum in 4 characters containing 6 bits of information each. The details of the base64 encoding algorithm are part of the [MIME](#) specification.

Note that a base64 number can include '/' and '+' characters. Although these characters are treated specially in the [URN syntax](#) specification, it is not necessary to escape them when used in a checksum URN.

If a content identifier contains a URI that refers to a well known checksum algorithm such as MD5 or SHA, it is possible to verify the integrity of the content. If none of the URIs in the content identifier refer to a known checksum algorithm, the content identifier should be treated as an opaque string that can be used for addressing purposes only.

Here are some examples of valid content identifiers:

```
urn:md5:FNG4c6MJLdDEY1rcoGb4pQ==  
urn:md5:HUXZLQLMuI/KZ5KDcJPcOA==  
urn:sha:thvDyvhfIqlvFe+A9MYgxAfmlq5=  
http://www.acme.com/images/foo.gif,urn:md5:FNG4c6MJLdDEY1rcoGb4pQ==  
http://www.acme.com/Example/,urn:x-version:5  
ftp://www.acme.com/Hello%20World
```

In applications where the possibility of duplicates for a given URN is unacceptable, content identifiers can be generated with the required uniqueness by using multiple appropriate URIs. If for example version numbers are used, it is important to further qualify the content identifier in order to make it globally unique. This can be achieved by including the URL of the object to which the version number applies in the content identifier.

A content identifier is often embedded in an HTTP header, and therefore the URIs in the content identifier are not allowed to contain reserved characters such as spaces or comma characters. All reserved characters should be encoded as specified in the [URI](#) specification.

2.2 Index Format

To describe the exact state of a set of data files, the DRP protocol uses a data-structure called an *index*. An index not only describes the hierarchical structure of the files, but it also describes the version, size, and type of each file. An index is a snapshot of the state of a set of files at a particular moment in time.

An index is typically stored in memory as a tree data structure, but in order for clients and servers to communicate this information over HTTP, an index can be described using XML. The index DTD used by the DRP protocol is specified in [Appendix A](#). Using this DTD, here is an example of an index that describes a hierarchical set of files:

```
1
2   <?XML VERSION="1.0" RMD="NONE"?>
3
4   <index>
5     <file path="home.html"      size="12345"  id="urn:md5:PEFjWBDv/sd9a1S9BYuX0w==" />
6     <file path="layer1.js"      size="32112"  id="urn:md5:W25YCu3toJt3ZsDsHIZmpg==" />
7     <dir  path="images">
8
9     <file path="acme.gif"       size="4532"   id="urn:md5:+hbZN5XfU6QAJB1RF1/KSQ==" />
10    <file path="banner.gif"     size="10452"  id="urn:md5:tr3X+oN3r9kqvsiyDSSjg==" />
11
12    </dir>
13
14    <dir  path="java/classes">
15
16    <file path="Scroll.java"     size="14323"  id="urn:md5:xjBkgWouS6p6FTUMIkx/Zg==" />
17
18    <file path="gui.jar"        size="540321" id="urn:md5:tcUzw0DKut3SiTpmAsi8g==" />
19
20    </dir>
21  </index>
```

2.3 Index Retrieval

A DRP index is retrieved given a URL to the index. The index can be stored in any file and can be retrieved using a normal HTTP GET request. The mime type of the index file allows clients to treat the index as a special file which provides meta information about other files. The client can use the index to automatically download the files that are specified.

The index file can be an ordinary file on an HTTP server, but it can also be generated dynamically. Note that the index isn't necessarily generated from a file system, it could also represent hierarchical data from a different source such as a database.

Once the initial download is complete, a client can update the content by downloading a new version of the index, and comparing it against the previous version of the index. Because each file entry in the index has a content identifier, the client can determine which files have changed and so determine the minimal set of files that need to be downloaded in order to bring the client up-to-date.

Index Base URL

The index file is typically located in the root directory of the file hierarchy it describes. In that case the base URL of the files in the index is the same as the base URL of the index itself. The index file should not be listed in the index itself.

An index can also explicitly specify a different absolute or relative base URL using the **base** attribute of the **index** tag. This allows an index to describe files that are located in a different directory, or even on a different server.

For example, if the index is loaded from the following URL:

```


```

Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

Real-Time Litigation Alerts



Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

Advanced Docket Research



With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

Analytics At Your Fingertips



Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.