

Digest Values for DOM (DOMHASH)

Status of this Memo

This memo provides information for the Internet community. It does not specify an Internet standard of any kind. Distribution of this memo is unlimited.

Copyright Notice

Copyright (C) The Internet Society (2000). All Rights Reserved.

Abstract

This memo defines a clear and unambiguous definition of digest (hash) values of the XML objects regardless of the surface string variation of XML. This definition can be used for XML digital signature as well efficient replication of XML objects.

Table of Contents

1. Introduction.....	2
2. Digest Calculation.....	3
2.1. Overview.....	3
2.2. Namespace Considerations.....	4
2.3. Definition with Code Fragments.....	5
2.3.1. Text Nodes.....	5
2.3.2. Processing Instruction Nodes.....	6
2.3.3. Attr Nodes.....	6
2.3.4. Element Nodes.....	7
2.3.5. Document Nodes.....	9
3. Discussion.....	9
4. Security Considerations.....	9
References.....	10
Authors' Addresses.....	10
Full Copyright Statement.....	11

1. Introduction

The purpose of this document is to give a clear and unambiguous definition of digest (hash) values of the XML objects [XML]. Two subtrees are considered identical if their hash values are the same, and different if their hash values are different.

There are at least two usage scenarios of DOMHASH. One is as a basis for digital signatures for XML. Digital signature algorithms normally require hashing a signed content before signing. DOMHASH provides a concrete definition of the hash value calculation.

The other is to use DOMHASH when synchronizing two DOM structures [DOM]. Suppose that a server program generates a DOM structure which is to be rendered by clients. If the server makes frequent small changes on a large DOM tree, it is desirable that only the modified parts are sent over to the client. A client can initiate a request by sending the root hash value of the structure in the cache memory. If it matches with the root hash value of the current server structure, nothing needs be sent. If not, then the server compares the client hash with the older versions in the server's cache. If it finds one that matches the client's version of the structure, then it locates differences with the current version by recursively comparing the hash values of each node. This way, the client can receive only an updated portion of a large structure without requesting the whole thing.

One way of defining digest values is to take a surface string as the input for a digest algorithm. However, this approach has several drawbacks. The same internal DOM structure may be represented in many different ways as surface strings even if they strictly conform to the XML specification. Treatment of white spaces, selection of character encodings, entity references (i.e., use of ampersands), and so on have impact on the generation of a surface string. If the implementations of surface string generation are different, the hash values would be different, resulting in unvalidatable digital signatures and unsuccessful detection of identical DOM structures. Therefore, it is desirable that digest of DOM is defined in the DOM terms -- that is, as an unambiguous algorithm operating on a DOM tree. This is the approach we take in this specification.

Introduction of namespace is another source of variation of surface string because different namespace prefixes can be used for representing the same namespace URI [URI]. In the following example, the namespace prefix "edi" is bound to the URI "http://ecommerce.org/schema" but this prefix can be arbitrarily chosen without changing the logical contents as shown in the second example.

```
<?xml version="1.0"?>
<root xmlns:edi='http://ecommerce.org/schema'>
  <edi:order>
    :
  </edi:order>
</root>
```

```
<?xml version="1.0"?>
<root xmlns:ec='http://ecommerce.org/schema'>
  <ec:order>
    :
  </ec:order>
</root>
```

The DOMHASH defined in this document is designed so that the choice of the namespace prefix does not affect the digest value. In the above example, both the "root" elements will get the same digest value.

2. Digest Calculation

2.1. Overview

Hash values are defined on the DOM type Node. We consider the following five node types that are used for representing a DOM document structure:

- Text
- ProcessingInstruction
- Attr
- Element
- Document

Comment nodes and Document Type Definitions (DTDs) do not participate in the digest value calculation. This is because DOM does not require a conformant processor to create data structures for these. DOMHASH is designed so that it can be computed with any XML processor conformant to the DOM or SAX [SAX] specification.

Nodes with the node type EntityReference must be expanded prior to digest calculation.

The digest values are defined recursively on each level of the DOM tree so that only a relevant part needs to be recalculated when a small portion of the tree is changed.

Below, we give the precise definitions of digest for these types. We describe the format of the data to be supplied to a hash algorithm using a figure and a simple description, followed by a Java code fragment using the DOM API and the JDK 1.1 Platform Core API only. Therefore, the semantics should be unambiguous.

As the rule of thumb, all strings are to be in UTF-16BE [UTF16]. If there is a sequence of text nodes without any element nodes in between, these text nodes are merged into one by concatenating them. A zero-length text node is always ignored.

Note that validating and non-validating XML processors may generate different DOM trees from the same XML document, due to attribute normalization and default attributes. If DOMHASH is to be used for testing logical equivalence between two XML documents (as opposed to DOM trees), it may be necessary to normalize attributes and supply default attributes prior to DOMHASH calculation.

Some legacy character encodings (such as ISO-2022-JP) have certain ambiguity in translating into Unicode. This is again dependent on XML processors. Treatment of such processor dependencies is out of scope of this document.

2.2. Namespace Considerations

To avoid the dependence on the namespace prefix, we use "expanded names" to do digest calculation. If an element name or an attribute name is qualified either by an explicit namespace prefix or by a default namespace, the name's LocalPart is prepended by the URI of the namespace (the namespace name as defined in the Namespace specification [NAM]) and a colon before digest calculation. In the following example, the default qualified name "order" is expanded into "http://ecommerce.org/schema:order" while the explicit qualified name "book:title" is expanded into "urn:loc.gov:books:title" before digest calculation.

```
<?xml version="1.0"?>

<root xmlns='http://ecommerce.org/schema'
      xmlns:book='urn:loc.gov:books'>
  <order>
    <book:title> ... </book:title>
    :
  </order>
</root>
```

We define an expanded name (either for element or attribute) as follows:

If a name is not qualified, the expanded name is the name itself.

If a name is qualified with the prefix "xmlns", the expanded name is undefined.

If a name is qualified either by default or by an explicit namespace prefix, the expanded name is URI bound to the namespace + ":" + LocalPart

In the following example code, we assume that the `getExpandedName()` method (which returns the expanded name as defined above) is defined in both `Element` and `Attr` interfaces of DOM.

Note that the digest values are not defined on namespace declarations. In other words, the digest value is not defined for an attribute when

- the attribute name is "xmlns", or
- the namespace prefix is "xmlns".

In the above example, the two attributes which are namespace declarations do not have digest values and therefore will not participate in the calculation of the digest value of the "root" element.

2.3. Definition with Code Fragments

The code fragments in the definitions below assume that they are in implementation classes of `Node`. Therefore, a methods call without an explicit object reference is for the `Node` itself. For example, `getData()` returns the text data of the current node if it is a `Text` node. The parameter `digestAlgorithm` is to be replaced by an identifier of the digest algorithm, such as "MD5" [MD5] and "SHA-1" [SHA].

The computation should begin with a four byte integer that represents the type of the node, such as `TEXT_NODE` or `ELEMENT_NODE`.

2.3.1. Text Nodes

The hash value of a `Text` node is computed on the four byte header followed by the UTF-16BE encoded text string.

- `TEXT_NODE` (3) in 32 bit network-byte-ordered integer
- Text data in UTF-16BE stream (variable length)

```

public byte[] getDigest(String digestAlgorithm) {
    MessageDigest md = MessageDigest.getInstance(digestAlgorithm);
    md.update((byte)0);
    md.update((byte)0);
    md.update((byte)0);
    md.update((byte)3);
    md.update(getData().getBytes("UnicodeBigUnmarked"));
    return md.digest();
}

```

Here, MessageDigest is in the package java.security.*, one of the built-in packages of JDK 1.1.

2.3.2. ProcessingInstruction Nodes

A ProcessingInstruction (PI) node has two components: the target and the data. Accordingly, the hash is computed on the concatenation of both, separated by 'x0000'. PI data is from the first non white space character after the target to the character immediately preceding the "?>".

- PROCESSING_INSTRUCTION_NODE (7) in 32 bit network-byte-ordered integer
- PI target in UTF-16BE stream (variable length)
- 0x00 0x00
- PI data in UTF-16BE stream (variable length)

```

public byte[] getDigest(String digestAlgorithm) {
    MessageDigest md = MessageDigest.getInstance(digestAlgorithm);
    md.update((byte)0);
    md.update((byte)0);
    md.update((byte)0);
    md.update((byte)7);
    md.update(getName().getBytes("UnicodeBigUnmarked"));
    md.update((byte)0);
    md.update((byte)0);
    md.update(getData().getBytes("UnicodeBigUnmarked"));
    return md.digest();
}

```

2.3.3. Attr Nodes

The digest value of Attr nodes are defined similarly to PI nodes, except that we need a separator between the expanded attribute name and the attribute value. The '0x0000' value in UTF-16BE is allowed nowhere in an XML document, so it can serve as an unambiguous separator. The expanded name must be used as the attribute name because it may be qualified. Note that if the attribute is a

namespace declaration (either the attribute name is "xmlns" or its prefix is "xmlns"), the digest value is undefined and the getDigest() method should return null.

- ATTRIBUTE_NODE (2) in 32 bit network-byte-ordered integer
- Expanded attribute name in UTF-16BE stream (variable length)
- 0x00 0x00
- Attribute value in UTF-16BE stream (variable length)

```
public byte[] getDigest(String digestAlgorithm) {
    if (getNodeName().equals("xmlns")
        || getNodeName().startsWith("xmlns:"))
        return null;
    MessageDigest md = MessageDigest.getInstance(digestAlgorithm);
    md.update((byte)0);
    md.update((byte)0);
    md.update((byte)0);
    md.update((byte)2);
    md.update(getExpandedName().getBytes("UnicodeBigUnmarked"));
    md.update((byte)0);
    md.update((byte)0);
    md.update(getValue().getBytes("UnicodeBigUnmarked"));
    return md.digest();
}
```

2.3.4. Element Nodes

Element nodes are the most complex because they consist of other nodes recursively. Hash values of these component nodes are used to calculate the node's digest so that we can save computation when the structure is partially changed.

First, all the attributes except for namespace declarations must be collected. This list is sorted lexicographically by the expanded attribute names (based on Unicode character code points). When no surrogate characters are involved, this is the same as sorting in ascending order in terms of the UTF-16BE encoded expanded attribute names, using the string comparison operator `String.compareTo()` in Java.

- ELEMENT_NODE (1) in 32 bit network-byte-ordered integer
- Expanded element name in UTF-16BE stream (variable length)
- 0x00 0x00
- A number of non-namespace-declaration attributes in 32 bit network-byte-ordered unsigned integer
- Sequence of digest values of non-namespace-declaration attributes, sorted lexicographically by expanded attribute names
- A number of child nodes (except for Comment nodes) in 32bit

```

network-byte-ordered unsigned integer
- Sequence of digest values of each child node except for Comment
  nodes (variable length) (A sequence of child texts is merged to one
  text. A zero-length text and Comment nodes are not counted as
  child)

public byte[] getDigest(String digestAlgorithm) {
    MessageDigest md = MessageDigest.getInstance(digestAlgorithm);
    ByteArrayOutputStream baos = new ByteArrayOutputStream();
    DataOutputStream dos = new DataOutputStream(baos);
    dos.writeInt(ELEMENT_NODE); // This is stored in network byte order
    dos.write(getExpandedName().getBytes("UnicodeBigUnmarked"));
    dos.write((byte)0);
    dos.write((byte)0);
    // Collect all attributes except for namespace declarations
    NamedNodeMap nnm = this.getAttributes();
    int len = nnm.getLength()
        // Find "xmlns" or "xmlns:foo" in nnm and omit it.
    ...
    dos.writeInt(len); // This is sorted in the network byte order
    // Sort attributes lexicographically by expanded attribute
    // names.
    ...
    // Assume that 'Attr[] aattr' has sorted Attribute instances.
    for (int i = 0; i < len; i++)
        dos.write(aattr[i].getDigest(digestAlgorithm));
    Node n = this.getFirstChild();
    // Assume that adjoining Texts are merged,
    // there is no 0-length Text, and
    // comment nodes are removed.
    len = this.getChildNodes().getLength();
    dos.writeInt(len); // This is stored in the network byte order
    while (n != null) {
        dos.write(n.getDigest(digestAlgorithm));
        n = n.getNextSibling();
    }
    dos.close();
    md.update(baos.toByteArray());
    return md.digest();
}

```

2.3.5. Document Nodes

A Document node may have PI nodes before and after the root Element node. The digest value of a Document node is computed based on the sequence of the digest values of the pre-root PI nodes, the root Element node, and the post-root PI nodes in this order. Comment nodes and DocumentType nodes, if any, are ignored.

- DOCUMENT_NODE (9) in 32 bit network-byte-ordered integer
- A number of child nodes (except for Comment and DocumentType nodes) in 32bit network-byte-ordered unsigned integer
- Sequence of digest values of each child node except for Comment and DocumentType nodes (variable length)

```
public byte[] getDigest(String digestAlgorithm) {
    MessageDigest md = MessageDigest.getInstance(digestAlgorithm);
    ByteArrayOutputStream baos = new ByteArrayOutputStream();
    DataOutputStream dos = new DataOutputStream(baos);
    dos.writeInt(DOCUMENT_NODE); //This is stored in network byte order

    // Assume that Comment and DocumentType nodes are removed and this
    // node has only an Element node and PI nodes.
    len = this.getChildNodes().getLength();
    dos.writeInt(len); // This is stored in the network byte order
    Node n = this.getFirstChild();
    while (n != null) {
        dos.write(n.getDigest(digestAlgorithm));
        n = n.getNextSibling();
    }
    dos.close();
    md.update(baos.toByteArray());
    return md.digest();
}
```

3. Discussion

The definition described above can be efficiently implemented with any XML processor that is conformant to either DOM and SAX specification. Reference implementations are available on request.

4. Security Considerations

DOMHASH is expected to be used as the basis for digital signatures and other security and integrity uses. It's appropriateness for such uses depends on the security of the hash algorithm used and inclusion of the fundamental characteristics it is desired to check in parts of the DOM model incorporated in the digest by DOMHASH.

References

- [DOM] "Document Object Model (DOM), Level 1 Specification", October 1998, <http://www.w3.org/TR/REC-DOM-Level-1/>
- [MD5] Rivest, R., "The MD5 Message-Digest Algorithm", RFC 1321, April 1992.
- [NAM] Tim Bray, Dave Hollander, Andrew Layman, "Namespaces in XML", <http://www.w3.org/TR/1999/REC-xml-names-19990114>.
- [SAX] David Megginson, "SAX 1.0: The Simple API for XML", <http://www.megginson.com/SAX/>, May 1998.
- [SHA] (US) National Institute of Standards and Technology, "Federal Information Processing Standards Publication 180-1: Secure Hash Standard", 17 April 1995.
- [URI] Berners-Lee, T., Fielding, R. and L. Masinter, "Uniform Resource Identifiers (URI): Generic Syntax", RFC 2396, August 1998.
- [UTF16] Hoffman, P., Yergeau, F., "UTF-16, an encoding of ISO 10646", RFC 2781, February 2000.
- [XML] Tim Bray, Jean Paoli, C. M. Sperber-McQueen, "Extensible Markup Language (XML) 1.0", <http://www.w3.org/TR/1998/REC-xml-19980210>

Authors' Addresses

Hiroshi Maruyama,
IBM Research, Tokyo Research Laboratory

EMail: maruyama@jp.ibm.com

Kent Tamura,
IBM Research, Tokyo Research Laboratory

EMail: kent@trl.ibm.co.jp

Naohiko Uramoto,
IBM Research, Tokyo Research Laboratory

EMail: uramoto@jp.ibm.com

Full Copyright Statement

Copyright (C) The Internet Society (2000). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Acknowledgment

Funding for the RFC Editor function is currently provided by the Internet Society.

