

# HTML Reader C++ Class Library

**By gUrM33T**Posted : **30 Mar 2004**Updated : **30 Mar 2004**Views : **82,357**

A lightweight, fast, simple, and low-overhead C++ class library based on push model parsing.

79 votes for this Article. ☐ ☐ ☐ ☐ ☐Popularity: **8.73** Rating: **4.60** out of 5[Download Library Source \(Fully Documented\) - 24 Kb](#)

## Contents

- [Introduction](#)
- [Events-Based Parser](#)
- [Files](#)
- [Brief Description of Classes](#)
- [Usage](#)
- [More About Event Handling](#)
- [Class View](#)
- [License](#)

## Introduction

After failing to search for a class library that allows to read HTML text from either in-memory string buffers or physical disk files, I decided that there is a severe need to have a library like that. There are many parsers available for XML (eXtensible Markup Language), for instance, Simple API for XML (SAX), that allow you to parse XML simply by handling events that the reader generates as it parses specific symbols from the given XML document.

Inspired by the SAX parser for XML, I decided to develop an HTML Reader C++ Class Library myself from scratch that offers a simple, lightweight, fast, and the most important, a low-overhead solution to process an HTML document. Like SAX, I decided to develop an events-based parser, which raises events as it encounters various elements in the document. The advantage of an events-based parser is that the reader reads a section of an HTML document, generates an event, and then moves on to the next section. It uses less memory and is better for processing large documents.

## Events-Based Parser

An events-based parser uses the callback mechanism to report parsing events. These callbacks turn out to be protected virtual member functions that you will override. Events, such as the detection of an opening tag or the closing tag of an element, will trigger a call to the corresponding member function of your class. The application implements and registers an event handler with the reader. It is upto the application to put some code in the event handlers designed to achieve the objective of the application. Events-based parsers provide a simple, fast, and a lower-level access to the document being parsed.

Events-based parsers do not create an in-memory representation of the source document.

They simply parse the document and notify client applications about various elements they find along the way. What happens next is the responsibility of the client application. Events-based parsers don't cache information and have an enviably small memory footprint.

## Files

To use the HTML Reader Class Library in your MFC application project, you will need to add a number of files in your project:

Header File	Source File	Class
LiteHTMLReader.h	LiteHTMLReader.cpp	CLiteHTMLReader
LiteHTMLTag.h	-	CLiteHTMLTag
LiteHTMLAttributes.h	-	CLiteHTMLAttributes
LiteHTMLAttributes.h	LiteHTMLAttributes.cpp	CLiteHTMLElemAttr
LiteHTMLEntityResolver.h	LiteHTMLEntityResolver.cpp	CLiteHTMLEntityResolver

NOTE: `LiteHTMLCommon.h` must also be included in your project.

## Brief Description of Classes

- `CLiteHTMLReader` is the main class of our library that works in conjunction with other `CLiteHTML*` classes to parse the given HTML document. It contains methods (`Read` and `ReadFile`) to initiate the parsing process that can operate on either in-memory string buffers or a physical disk file. `CLiteHTMLReader` allows you to trap events that the reader generates as it finds various elements in the document such as the starting of a tag, ending of a tag, an HTML comment, etc. But to handles these events, your application must define a class that implements an interface `ILiteHTMLReaderEvents` declared in the `LiteHTMLReader.h` file.
- `CLiteHTMLTag` class, as its name implies, is related to the HTML tags. It deals with the parsing and storage of tag information from the given string such as the name of the tag and the attributes/properties of a tag. It provides a method (actually, all the above-specified classes provide a method named `parseFromStr`) that is called by the `CLiteHTMLReader` class' `Read` and `ReadFile` method as the document is being parsed. Typically, `CLiteHTMLTag` is not used directly by your application. As specified above, it works in conjunction with the reader helping in the parsing of HTML tags.
- The `CLiteHTMLElemAttr` and `CLiteHTMLAttributes` classes are inter-related as `CLiteHTMLAttributes` provides a collection-based mechanism to hold an array of `CLiteHTMLElemAttr` objects that are accessible either by the name of the attribute or a zero-based index value. As was the case with the `CLiteHTMLTag` class, these classes are also not typically used by your application directly.
- The last is the `CLiteHTMLEntityResolver` class that helps in resolving the entity references. Entity references are numeric or symbolic names for characters that may be included in an HTML document. They are useful for referring to rarely used characters, or those that authoring tools make it difficult or impossible to enter. Entity references begin with a "&" sign and end with a semi-colon (;). Some common examples are: `&lt;`; representing the `<` sign, `&gt;`; representing the `>` sign, etc.

From the above discussion, one thing is clear that the `CLiteHTMLTag`, `CLiteHTMLAttributes`, and, `CLiteHTMLElemAttr` class provide a method named `parseFromStr` that is used by the `CLiteHTMLReader` to further delegate the parsing process

while reading an HTML document.

## Usage

OK, now let's come to the part of learning how to use this library in an MFC project:

1. The **first** step is pretty simple. All you have to do is to add all of the files (given in the [FILES](#) section above) in your project.
2. The **second** step, although optional, is to create a class that implements [ILiteHTMLReaderEvents](#) interface. [ILiteHTMLReaderEvents](#) is, in actual, an abstract class that acts as an interface which must be implemented by all those classes that need to handle events raised by the [CLiteHTMLReader](#) class. For example,

```
#include "stdafx.h"

#include "LiteHTMLReader.h"

class CEventHandler : public ILiteHTMLReaderEvents
{
private:
    void BeginParse(DWORD dwAppData, bool &bAbort);
    void StartTag(CLiteHTMLTag *pTag, DWORD dwAppData, bool &bAbort);
    void EndTag(CLiteHTMLTag *pTag, DWORD dwAppData, bool &bAbort);
    void Characters(const CString &rText, DWORD dwAppData, bool &bAbort);
    void Comment(const CString &rComment, DWORD dwAppData, bool &bAbort);
    void EndParse(DWORD dwAppData, bool bIsAborted);
};
```

You must have noticed that "optional" word I used above. The reason behind this is that if you do not provide your own implementation of the event handler(s), the [ILiteHTMLReaderEvents](#) class provides a default implementation that does nothing. To learn more about the [ILiteHTMLReaderEvents](#) interface, jump to the [ILiteHTMLReaderEvents Described](#) section of this article.

3. The **third** step is to create an instance of the [CLiteHTMLReader](#) class like this:

```
CLiteHTMLReader theReader;
```

4. Now we should call either [Read](#) or [ReadFile](#) method of the [CLiteHTMLReader](#) class???

NO! Our event handler implementation will not start receiving notifications until we register it with the reader by calling the [setEventHandler](#) method of the [CLiteHTMLReader](#) class. So, supposing that the name of our class that is implementing the [ILiteHTMLReaderEvents](#) interface is [CEventHandler](#), the **fourth** step is to create an instance of the [CEventHandler](#), and call [setEventHandler](#) by passing it the address of this instance variable.

```
CEventHandler theEventHandler;
theReader.setEventHandler(&theEventHandler);
```

Now, for all of you, who are thinking if it is possible to pass a [NULL](#) pointer to the [setEventHandler](#) method, the answer is YES and that too at any time you want. And not to mention, you can also change the event handler at any time by calling [setEventHandler](#) and passing the address of some other instance.

5. Now, the **fifth** and the final step is to call either `Read` or `ReadFile` method on the `CLiteHTMLReader` instance variable we created in step 3 by passing it the appropriate parameter i.e. if you decide to parse an in-memory string buffer, call the `Read` method and pass the address of the string you want to parse. In case, you need to parse an HTML document from a disk file, you can call another method `ReadFile` that is similar to `Read` but accepts a file handle (`HANDLE`) instead of a pointer to an array of characters. Take a look at the example:

```
TCHAR    strToParse[] = _T("<HTML>"
    "<HEAD>"
    "<TITLE>"
    "<!-- title goes here -->"
    "</TITLE>"
    "</HEAD>"
    "<BODY LEFTMARGIN=15px>This is a sample HTML document.</BODY>"
    "</HTML>");
theReader.Read(strToParse);
```

OR

```
CFile    fileToParse;
if (fileToParse.Open(_T("test.html"), CFile::modeRead))
{
    theReader.ReadFile(fileToParse.m_hFile);
    fileToParse.Close();
}
```

## More About Event Handling

The `ILiteHTMLReaderEvents` class presents an interface that must be implemented by all those classes that want to handle the notifications sent by the `CLiteHTMLReader` while parsing an HTML document. The order of events handled by the `ILiteHTMLReaderEvents` handler is determined by the order of information within the document being parsed. It's important to note that the interface includes a series of methods that the `CLiteHTMLReader` invokes during the parsing operation. The reader passes the appropriate information to the method's parameters. To perform some type of processing for a method, you simply add code to the method in your own `ILiteHTMLReaderEvents` implementation.

The common parameters received by all of the methods defined in `ILiteHTMLReaderEvents` class, except `EndParse` include:

- `dwAppData`: A 32-bit application-specific data.
- `bAbort`: You can set this parameter to either `true` or `false` according to your application's needs to specify whether the reader should continue parsing rest of the data in the buffer or aborts immediately after the current event handler completes processing.

The `EndParse` method receives `bIsAborted` parameter instead of the `bAbort` that signifies if `EndParse` has occurred because of the normal parsing termination.

Along with the above-specified parameters, all of the methods except `BeginParse` and `EndParse`, receive some extra information (specific to the event) which is retrieved by the reader while parsing the HTML document. For instance, when an HTML tag (either opening or closing) is parsed, the `StartTag` or `EndTag` methods receive a pointer to a `CLiteHTMLTag` that contains the name of the tag and the attributes (if any) of the tag. Attribute information is retrieved only if the tag parsed is an opening tag as closing tags cannot contain any

attribute/value pairs. If there is no attribute information associated with a `CLiteHTMLTag`, the pointer variable contains `NULL`. So it is obvious that `EndTag` method always receives a `NULL` pointer. It is the responsibility of an application (and a good programming practice) to check for `NULL` pointer before using it.

Similarly, the `Comment` and `Characters` method of the class receives a reference to a `CString` containing the extracted text. The `Comment` method receives `rComment` parameter containing the comment text excluding the delimiters i.e. without `<!--` and `-->`. The `Characters` method receives a `rText` parameter that signifies either the contents of an element or some text that could not be parsed by the reader.

## Class View

### CLiteHTMLReader Class Members

Member	Description
<code>CLiteHTMLReader()</code>	Constructs a <code>CLiteHTMLReader</code> object.
<code>EventMaskEnum setEventMask(DWORD);</code>	Sets a new event mask.
<code>EventMaskEnum setEventMask(DWORD, DWORD);</code>	Changes the current event mask by adding and/or removing flags.
<code>EventMaskEnum getEventMask(void) const;</code>	Returns the event mask previously set by a call to <code>setEventMask</code> .
<code>DWORD setAppData(DWORD);</code>	Sets application-specific data to be passed to event handlers.
<code>DWORD getAppData(void) const;</code>	Returns app-specific data previously set by a call to <code>setAppData</code> .
<code>ILiteHTMLReaderEvents* setEventHandler(ILiteHTMLReaderEvents*);</code>	Registers an event handler with the reader.
<code>ILiteHTMLReaderEvents* getEventHandler(void) const;</code>	Returns the currently associated event handler.
<code>UINT Read(LPCTSTR);</code>	Parses an HTML document from the specified string.
<code>UINT Read(HANDLE);</code>	Parses an HTML document from a file given its <code>HANDLE</code> .

### CLiteHTMLTag Class Members

Member	Description
<code>CLiteHTMLTag()</code>	Constructs a <code>CLiteHTMLTag</code> object.
<code>CLiteHTMLTag(CLiteHTMLTag&amp;, bool)</code>	Constructs a <code>CLiteHTMLTag</code> object from an existing instance. The first parameter is the reference to a source <code>CLiteHTMLTag</code> , and the second

	parameter determines whether to make a copy or to take ownership of the encapsulated <code>CLiteHTMLAttributes</code> pointer.
<code>~CLiteHTMLTag()</code>	Destroys a <code>CLiteHTMLTag</code> object.
<code>CString getTagName(void) const;</code>	Returns the name of the tag.
<code>const CLiteHTMLAttributes* getAttributes(void) const;</code>	Returns a pointer to an attribute collection associated with this <code>CLiteHTMLTag</code> .
<code>UINT parseFromStr(LPCTSTR, bool&amp;, bool&amp;, bool);</code>	Parses an HTML tag from the string specified by the first parameter. The second and third parameter receive a boolean true/false indicating that the tag parsed is an opening and/or closing tag, respectively. The fourth parameter specifies whether to parse tag's attributes also.

### CLiteHTMLAttributes Class Members

Member	Description
<code>CLiteHTMLAttributes()</code>	Constructs a <code>CLiteHTMLAttributes</code> object.
<code>CLiteHTMLAttributes(CLiteHTMLAttributes&amp;, bool)</code>	Constructs a <code>CLiteHTMLAttributes</code> object from an existing instance. The first parameter is the reference to a source <code>CLiteHTMLAttributes</code> , and the second parameter determines whether to make a copy or to take ownership of the encapsulated pointer.
<code>~CLiteHTMLAttributes()</code>	Destroys a <code>CLiteHTMLAttributes</code> object.
<code>int getCount() const;</code>	Returns the count of <code>CLiteHTMLElemAttr</code> items.
<code>int getIndexFromName(LPCTSTR) const;</code>	Looks up the index of an attribute given its name.
<code>CLiteHTMLElemAttr operator[](int) const;</code>	Returns a <code>CLiteHTMLElemAttr</code> object given an attribute's index.
<code>CLiteHTMLElemAttr getAttribute(int) const;</code>	Returns a <code>CLiteHTMLElemAttr</code> object given an attribute's index.

<code>CLiteHTMLElemAttr operator[] (LPCTSTR) const;</code>	Returns a <code>CLiteHTMLElemAttr</code> object given an attribute name.
<code>CLiteHTMLElemAttr getAttribute(LPCTSTR) const;</code>	Returns a <code>CLiteHTMLElemAttr</code> object given an attribute name.
<code>CString getName(int) const;</code>	Returns the name of an attribute given its index.
<code>CString getValue(int) const;</code>	Returns the value of an attribute given its index.
<code>CString getValueFromName(LPCTSTR) const;</code>	Returns the value of an attribute given its name.
<code>CLiteHTMLElemAttr* addAttribute(LPCTSTR, LPCTSTR);</code>	Adds a new <code>CLiteHTMLElemAttr</code> item to the collection.
<code>bool removeAttribute(int);</code>	Removes an <code>CLiteHTMLElemAttr</code> item from the collection.
<code>bool removeAll(void);</code>	Removes all <code>CLiteHTMLElemAttr</code> items from the collection.
<code>UINT parseFromStr(LPCTSTR);</code>	Parses attribute/value pairs from the given string.

## CLiteHTMLElemAttr Class Members

Member	Description
<code>CString getName(void) const;</code>	Returns the name of an <code>CLiteHTMLElemAttr</code> .
<code>CString getValue(void) const;</code>	Returns the value of an <code>CLiteHTMLElemAttr</code> .
<code>bool isColorValue(void) const;</code>	Determines if the attribute value contains a color reference.
<code>bool isNamedColorValue(void) const;</code>	Determines if the attribute value is a named color value.
<code>bool isSysColorValue(void) const;</code>	Determines if the attribute value is a named system color value.
<code>bool isHexColorValue(void) const;</code>	Determines if the attribute value is a color value in hexadecimal format.
<code>bool isPercentValue(void) const;</code>	Checks to see if the attribute contains a percent value.
<code>COLORREF getColorValue(void) const;</code>	Returns the color value of the attribute.
<code>CString getColorHexValue(void) const;</code>	Returns the RGB value of the attribute in hexadecimal format.
<code>unsigned short getPercentValue() const;</code>	Returns a percent value of the attribute.

<code>short</code>	Returns a length value of the attribute.
<code>getLengthValue(LengthUnitsEnum&amp;) const;</code>	
<code>operator bool() const;</code>	Converts attribute value to <code>bool</code> .
<code>operator BYTE() const;</code>	Converts attribute value to <code>BYTE</code> ( <code>unsigned char</code> ).
<code>operator double() const;</code>	Converts attribute value to <code>double</code> .
<code>operator short() const;</code>	Converts attribute value to <code>signed short int</code> .
<code>operator LPCTSTR() const;</code>	Returns the value of the attribute.
<code>UINT parseFromStr(LPCTSTR);</code>	Parses an attribute/value pair from the given string.

## License

*This code may be used in compiled form in any way you desire (including commercial use). The code may be redistributed unmodified by any means **providing** it is not sold for profit without the authors written consent, and providing that this notice and the authors name and all copyright notices remains intact. However, this file and the accompanying source code **may not** be hosted on a website or bulletin board without the authors written permission.*

***This software is provided "AS IS" without express or implied warranty. The author accepts no liability for any damage/loss of business that this product may cause. Use it at your own risk!***

## About the Author

**gUrM33T**



Occupation: Web Developer

Location:  Canada

## Discussions and Feedback

 **62 messages** have been posted for this article. Visit <http://www.codeproject.com/KB/library/GomzyHTMLReader.aspx> to post and view comments on this article, or click [here](#) to get a print view with messages.

[PermaLink](#) | [Privacy](#) | [Terms of Use](#)  
Last Updated: 30 Mar 2004  
Editor: [Nishant Sivakumar](#)

Copyright 2004 by [gUrM33T](#)  
Everything else Copyright © [CodeProject](#), 1999-2007  
Web08 | [Advertise on the Code Project](#)