

# Package ‘xml2’

February 28, 2025

**Title** Parse XML

**Version** 1.3.7

**Description** Bindings to 'libxml2' for working with XML data using a simple, consistent interface based on 'XPath' expressions. Also supports XML schema validation; for 'XSLT' transformations see the 'xslt' package.

**License** MIT + file LICENSE

**URL** <https://xml2.r-lib.org>, <https://r-lib.r-universe.dev/xml2>

**BugReports** <https://github.com/r-lib/xml2/issues>

**Depends** R (>= 3.6.0)

**Imports** cli, methods, rlang (>= 1.1.0)

**Suggests** covr, curl, httr, knitr, magrittr, mockery, rmarkdown, testthat (>= 3.2.0), xslt

**VignetteBuilder** knitr

**Config/Needs/website** tidyverse/tidytemplate

**Encoding** UTF-8

**RoxygenNote** 7.2.3

**SystemRequirements** libxml2: libxml2-dev (deb), libxml2-devel (rpm)

**Collate** 'S4.R' 'as\_list.R' 'xml\_parse.R' 'as\_xml\_document.R' 'classes.R' 'format.R' 'import-standalone-obj-type.R' 'import-standalone-purrr.R' 'import-standalone-types-check.R' 'init.R' 'nodeset\_apply.R' 'paths.R' 'utils.R' 'xml2-package.R' 'xml\_attr.R' 'xml\_children.R' 'xml\_document.R' 'xml\_find.R' 'xml\_missing.R' 'xml\_modify.R' 'xml\_name.R' 'xml\_namespaces.R' 'xml\_node.R' 'xml\_nodeset.R' 'xml\_path.R' 'xml\_schema.R' 'xml\_serialize.R' 'xml\_structure.R' 'xml\_text.R' 'xml\_type.R' 'xml\_url.R' 'xml\_write.R' 'zzz.R'

**Config/testthat/edition** 3

**NeedsCompilation** yes

**Author** Hadley Wickham [aut],  
 Jim Hester [aut],  
 Jeroen Ooms [aut, cre],  
 Posit Software, PBC [cph, fnd],  
 R Foundation [ctb] (Copy of R-project homepage cached as example)

**Maintainer** Jeroen Ooms <jeroenooms@gmail.com>

**Repository** CRAN

**Date/Publication** 2025-02-28 16:10:01 UTC

## Contents

as_list . . . . .	3
as_xml_document . . . . .	4
download_xml . . . . .	4
read_xml . . . . .	6
url_absolute . . . . .	8
url_escape . . . . .	9
url_parse . . . . .	10
write_xml . . . . .	10
xml2_example . . . . .	11
xml_attr . . . . .	12
xml_cdata . . . . .	13
xml_children . . . . .	14
xml_comment . . . . .	15
xml_dtd . . . . .	16
xml_find_all . . . . .	17
xml_name . . . . .	19
xml_new_document . . . . .	20
xml_ns . . . . .	21
xml_ns_strip . . . . .	22
xml_path . . . . .	22
xml_replace . . . . .	23
xml_serialize . . . . .	24
xml_set_namespace . . . . .	25
xml_structure . . . . .	25
xml_text . . . . .	26
xml_type . . . . .	27
xml_url . . . . .	27
xml_validate . . . . .	28

<b>Index</b>	<b>29</b>
--------------	-----------

---

`as_list`*Coerce xml nodes to a list.*

---

## Description

This turns an XML document (or node or nodeset) into the equivalent R list. Note that this is `as_list()`, not `as.list()`: `lapply()` automatically calls `as.list()` on its inputs, so we can't override the default.

## Usage

```
as_list(x, ns = character(), ...)
```

## Arguments

<code>x</code>	A document, node, or node set.
<code>ns</code>	Optionally, a named vector giving prefix-url pairs, as produced by <code>xml_ns()</code> . If provided, all names will be explicitly qualified with the <code>ns</code> prefix, i.e. if the element <code>bar</code> is defined in namespace <code>foo</code> , it will be called <code>foo:bar</code> . (And similarly for attributes). Default namespaces must be given an explicit name. The <code>ns</code> is ignored when using <code>xml_name&lt;-()</code> and <code>xml_set_name()</code> .
<code>...</code>	Needed for compatibility with generic. Unused.

## Details

`as_list` currently only handles the four most common types of children that an element might have:

- Other elements, converted to lists.
- Attributes, stored as R attributes. Attributes that have special meanings in R (`class()`, `comment()`, `dim()`, `dimnames()`, `names()`, `row.names()` and `tsp()`) are escaped with `'`.
- Text, stored as a character vector.

## Examples

```
as_list(read_xml("<foo> a <b /><c><![CDATA[<d></d>]]></c></foo>"))
as_list(read_xml("<foo> <bar><baz /></bar> </foo>"))
as_list(read_xml("<foo id = 'a'></foo>"))
as_list(read_xml("<foo><bar id='a' /><bar id='b' /></foo>"))
```

---

as_xml_document	<i>Coerce a R list to xml nodes.</i>
-----------------	--------------------------------------

---

### Description

This turns an R list into the equivalent XML document. Not all R lists will produce valid XML, in particular there can only be one root node and all child nodes need to be named (or empty) lists. R attributes become XML attributes and R names become XML node names.

### Usage

```
as_xml_document(x, ...)
```

### Arguments

x	A document, node, or node set.
...	Needed for compatibility with generic. Unused.

### Examples

```
as_xml_document(list(x = list()))

# Nesting multiple nodes
as_xml_document(list(foo = list(bar = list(baz = list()))))

# attributes are stored as R attributes
as_xml_document(list(foo = structure(list(), id = "a")))
as_xml_document(list(foo = list(
  bar = structure(list(), id = "a"),
  bar = structure(list(), id = "b")
)))
```

---

download_xml	<i>Download a HTML or XML file</i>
--------------	------------------------------------

---

### Description

Libcurl implementation of `C_download` (the "internal" download method) with added support for https, ftps, gzip, etc. Default behavior is identical to `download.file()`, but request can be fully configured by passing a custom `curl::handle()`.

**Usage**

```
download_xml(  
  url,  
  file = basename(url),  
  quiet = TRUE,  
  mode = "wb",  
  handle = curl::new_handle()  
)
```

```
download_html(  
  url,  
  file = basename(url),  
  quiet = TRUE,  
  mode = "wb",  
  handle = curl::new_handle()  
)
```

**Arguments**

url	A character string naming the URL of a resource to be downloaded.
file	A character string with the name where the downloaded file is saved.
quiet	If TRUE, suppress status messages (if any), and the progress bar.
mode	A character string specifying the mode with which to write the file. Useful values are "w", "wb" (binary), "a" (append) and "ab".
handle	a curl handle object

**Details**

The main difference between `curl_download` and `curl_fetch_disk` is that `curl_download` checks the http status code before starting the download, and raises an error when status is non-successful. The behavior of `curl_fetch_disk` on the other hand is to proceed as normal and write the error page to disk in case of a non success response.

For a more advanced download interface which supports concurrent requests and resuming large files, have a look at the [multi\\_download](#) function.

**Value**

Path of downloaded file (invisibly).

**See Also**

[curl\\_download](#)

**Examples**

```
## Not run:  
download_html("http://tidyverse.org/index.html")
```

```
## End(Not run)
```

---

read_xml	<i>Read HTML or XML.</i>
----------	--------------------------

---

### Description

Read HTML or XML.

### Usage

```
read_xml(x, encoding = "", ..., as_html = FALSE, options = "NOBLANKS")

read_html(x, encoding = "", ..., options = c("RECOVER", "NOERROR", "NOBLANKS"))

## S3 method for class 'character'
read_xml(x, encoding = "", ..., as_html = FALSE, options = "NOBLANKS")

## S3 method for class 'raw'
read_xml(
  x,
  encoding = "",
  base_url = "",
  ...,
  as_html = FALSE,
  options = "NOBLANKS"
)

## S3 method for class 'connection'
read_xml(
  x,
  encoding = "",
  n = 64 * 1024,
  verbose = FALSE,
  ...,
  base_url = "",
  as_html = FALSE,
  options = "NOBLANKS"
)
```

### Arguments

x	A string, a connection, or a raw vector. A string can be either a path, a url or literal xml. Urls will be converted into connections either using <code>base::url</code> or, if installed, <code>curl::curl</code> . Local paths ending in <code>.gz</code> , <code>.bz2</code> , <code>.xz</code> , <code>.zip</code> will be automatically uncompressed. If a connection, the complete connection is read into a raw vector before being parsed.
---	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

encoding	Specify a default encoding for the document. Unless otherwise specified XML documents are assumed to be in UTF-8 or UTF-16. If the document is not UTF-8/16, and lacks an explicit encoding directive, this allows you to supply a default.
...	Additional arguments passed on to methods.
as_html	Optionally parse an xml file as if it's html.
options	Set parsing options for the libxml2 parser. Zero or more of <b>RECOVER</b> recover on errors <b>NOENT</b> substitute entities <b>DTDLOAD</b> load the external subset <b>DTDATTR</b> default DTD attributes <b>DTDVALID</b> validate with the DTD <b>NOERROR</b> suppress error reports <b>NOWARNING</b> suppress warning reports <b>PEDANTIC</b> pedantic error reporting <b>NOBLANKS</b> remove blank nodes <b>SAX1</b> use the SAX1 interface internally <b>XINCLUDE</b> Implement XInclude substitution <b>NONET</b> Forbid network access <b>NODICT</b> Do not reuse the context dictionary <b>NSCLEAN</b> remove redundant namespaces declarations <b>NOCDATA</b> merge CDATA as text nodes <b>NOXINCNODE</b> do not generate XINCLUDE START/END nodes <b>COMPACT</b> compact small text nodes; no modification of the tree allowed afterwards (will possibly crash if you try to modify the tree) <b>OLD10</b> parse using XML-1.0 before update 5 <b>NOBASEFIX</b> do not fixup XINCLUDE xml:base uris <b>HUGE</b> relax any hardcoded limit from the parser <b>OLDSAX</b> parse using SAX2 interface before 2.7.0 <b>IGNORE_ENC</b> ignore internal document encoding hint <b>BIG_LINES</b> Store big lines numbers in text PSVI field
base_url	When loading from a connection, raw vector or literal html/xml, this allows you to specify a base url for the document. Base urls are used to turn relative urls into absolute urls.
n	If file is a connection, the number of bytes to read per iteration. Defaults to 64kb.
verbose	When reading from a slow connection, this prints some output on every iteration so you know its working.

**Value**

An XML document. HTML is normalised to valid XML - this may not be exactly the same transformation performed by the browser, but it's a reasonable approximation.

### Setting the "user agent" header

When performing web scraping tasks it is both good practice — and often required — to set the **user agent** request header to a specific value. Sometimes this value is assigned to emulate a browser in order to have content render in a certain way (e.g. Mozilla/5.0 (Windows NT 5.1; rv:52.0) Gecko/20100101 Firefox/52.0 to emulate more recent Windows browsers). Most often, this value should be set to provide the web resource owner information on who you are and the intent of your actions like this Google scraping bot user agent identifier: Googlebot/2.1 (+http://www.google.com/bot.html).

You can set the HTTP user agent for URL-based requests using `httr::set_config()` and `httr::user_agent()`:

```
httr::set_config(httr::user_agent("me@example.com; +https://example.com/info.html"))
```

`httr::set_config()` changes the configuration globally, `httr::with_config()` can be used to change configuration temporarily.

### Examples

```
# Literal xml/html is useful for small examples
read_xml("<foo><bar /></foo>")
read_html("<html><title>Hi</title></html>")
read_html("<html><title>Hi")

# From a local path
read_html(system.file("extdata", "r-project.html", package = "xml2"))

## Not run:
# From a url
cd <- read_xml(xml2_example("cd_catalog.xml"))
me <- read_html("http://had.co.nz")

## End(Not run)
```

---

url\_absolute

*Convert between relative and absolute urls.*

---

### Description

Convert between relative and absolute urls.

### Usage

```
url_absolute(x, base)
```

```
url_relative(x, base)
```

### Arguments

`x` A character vector of urls relative to that base

`base` A string giving a base url.



**Value**

A character vector of urls

**See Also**

[xml\\_url](#) to retrieve the URL associated with a document

**Examples**

```
url_absolute(c(".", "..", "/", "/x"), "http://hadley.nz/a/b/c/d")

url_relative("http://hadley.nz/a/c", "http://hadley.nz")
url_relative("http://hadley.nz/a/c", "http://hadley.nz/")
url_relative("http://hadley.nz/a/c", "http://hadley.nz/a/b")
url_relative("http://hadley.nz/a/c", "http://hadley.nz/a/b/")
```

---

url\_escape

*Escape and unescape urls.*

---

**Description**

Escape and unescape urls.

**Usage**

```
url_escape(x, reserved = "")
url_unescape(x)
```

**Arguments**

x                    A character vector of urls.  
reserved            A string containing additional characters to avoid escaping.

**Examples**

```
url_escape("a b c")
url_escape("a b c", "")

url_unescape("a%20b%2fc")
url_unescape("%C2%B5")
```

---

url_parse	<i>Parse a url into its component pieces.</i>
-----------	-----------------------------------------------

---

**Description**

Parse a url into its component pieces.

**Usage**

```
url_parse(x)
```

**Arguments**

x                    A character vector of urls.

**Value**

A dataframe with one row for each element of x and columns: scheme, server, port, user, path, query, fragment.

**Examples**

```
url_parse("http://had.co.nz/")
url_parse("http://had.co.nz:1234/")
url_parse("http://had.co.nz:1234/?a=1&b=2")
url_parse("http://had.co.nz:1234/?a=1&b=2#def")
```

---

write_xml	<i>Write XML or HTML to disk.</i>
-----------	-----------------------------------

---

**Description**

This writes out both XML and normalised HTML. The default behavior will output the same format which was read. If you want to force output pass option = "as\_xml" or option = "as\_html" respectively.

**Usage**

```
write_xml(x, file, ...)
```

```
## S3 method for class 'xml_document'
write_xml(x, file, ..., options = "format", encoding = "UTF-8")
```

```
write_html(x, file, ...)
```

```
## S3 method for class 'xml_document'
write_html(x, file, ..., options = "format", encoding = "UTF-8")
```

**Arguments**

x	A document or node to write to disk. It's not possible to save nodesets containing more than one node.
file	Path to file or connection to write to.
...	additional arguments passed to methods.
options	default: 'format'. Zero or more of <b>format</b> Format output <b>no_declaration</b> Drop the XML declaration <b>no_empty_tags</b> Remove empty tags <b>no_xhtml</b> Disable XHTML1 rules <b>require_xhtml</b> Force XHTML rules <b>as_xml</b> Force XML output <b>as_html</b> Force HTML output <b>format_whitespace</b> Format with non-significant whitespace
encoding	The character encoding to use in the document. The default encoding is 'UTF-8'. Available encodings are specified at <a href="http://xmlsoft.org/html/libxml-encoding.html#xmlCharEncoding">http://xmlsoft.org/html/libxml-encoding.html#xmlCharEncoding</a> .

**Examples**

```
h <- read_html("<p>Hi!</p>")

tmp <- tempfile(fileext = ".xml")
write_xml(h, tmp, options = "format")
readLines(tmp)

# write formatted HTML output
write_html(h, tmp, options = "format")
readLines(tmp)
```

---

xml2\_example

*Get path to a xml2 example*


---

**Description**

xml2 comes bundled with a number of sample files in its 'inst/extdata' directory. This function makes them easy to access.

**Usage**

```
xml2_example(path = NULL)
```

**Arguments**

path	Name of file. If NULL, the example files will be listed.
------	----------------------------------------------------------

---

xml_attr	<i>Retrieve an attribute.</i>
----------	-------------------------------

---

### Description

xml\_attrs() retrieves all attributes values as a named character vector, xml\_attrs() <- or xml\_set\_attrs() sets all attribute values. xml\_attr() retrieves the value of single attribute and xml\_attr() <- or xml\_set\_attr() modifies its value. If the attribute doesn't exist, it will return default, which defaults to NA. xml\_has\_attr() tests if an attribute is present.

### Usage

```
xml_attr(x, attr, ns = character(), default = NA_character_)
```

```
xml_has_attr(x, attr, ns = character())
```

```
xml_attrs(x, ns = character())
```

```
xml_attr(x, attr, ns = character()) <- value
```

```
xml_set_attr(x, attr, value, ns = character())
```

```
xml_attrs(x, ns = character()) <- value
```

```
xml_set_attrs(x, value, ns = character())
```

### Arguments

x	A document, node, or node set.
attr	Name of attribute to extract.
ns	Optionally, a named vector giving prefix-url pairs, as produced by <a href="#">xml_ns()</a> . If provided, all names will be explicitly qualified with the ns prefix, i.e. if the element bar is defined in namespace foo, it will be called foo:bar. (And similarly for attributes). Default namespaces must be given an explicit name. The ns is ignored when using <a href="#">xml_name&lt;-()</a> and <a href="#">xml_set_name()</a> .
default	Default value to use when attribute is not present.
value	character vector of new value.

### Value

xml\_attr() returns a character vector. NA is used to represent of attributes that aren't defined.

xml\_has\_attr() returns a logical vector.

xml\_attrs() returns a named character vector if x x is single node, or a list of character vectors if given a nodeset

**Examples**

```

x <- read_xml("<root id='1'><child id='a' /><child id='b' d='b' /></root>")
xml_attr(x, "id")
xml_attr(x, "apple")
xml_attrs(x)

kids <- xml_children(x)
kids
xml_attr(kids, "id")
xml_has_attr(kids, "id")
xml_attrs(kids)

# Missing attributes give missing values
xml_attr(xml_children(x), "d")
xml_has_attr(xml_children(x), "d")

# If the document has a namespace, use the ns argument and
# qualified attribute names
x <- read_xml('
<root xmlns:b="http://bar.com" xmlns:f="http://foo.com">
  <doc b:id="b" f:id="f" id="" />
</root>
')
doc <- xml_children(x)[[1]]
ns <- xml_ns(x)

xml_attrs(doc)
xml_attrs(doc, ns)

# If you don't supply a ns spec, you get the first matching attribute
xml_attr(doc, "id")
xml_attr(doc, "b:id", ns)
xml_attr(doc, "id", ns)

# Can set a single attribute with `xml_attr() <-` or `xml_set_attr()`
xml_attr(doc, "id") <- "one"
xml_set_attr(doc, "id", "two")

# Or set multiple attributes with `xml_attrs()` or `xml_set_attrs()`
xml_attrs(doc) <- c("b:id" = "one", "f:id" = "two", "id" = "three")
xml_set_attrs(doc, c("b:id" = "one", "f:id" = "two", "id" = "three"))

```

---

`xml_cdata`*Construct a cdata node*

---

**Description**

Construct a cdata node

**Usage**

```
xml_cdata(content)
```

**Arguments**

content            The CDATA content, does not include </p>
</div>
<div data-bbox="159 238 234 254" data-label="Section-Header"><b>Examples</b></div>
<div data-bbox="187 262 450 303" data-label="Text">
<pre>x &lt;- xml\_new\_root("root")
xml\_add\_child(x, xml\_cdata("&lt;d/&gt;"))
as.character(x)</pre>
</div>
<div data-bbox="175 347 577 363" data-label="Text">
<hr/>
<table>
<tr>
<td>xml\_children</td>
<td><i>Navigate around the family tree.</i></td>
</tr>
</table>
<hr/>
</div>
<div data-bbox="159 400 246 416" data-label="Section-Header"><b>Description</b></div>
<div data-bbox="187 424 840 485" data-label="Text">
<p>xml\_children returns only elements, xml\_contents returns all nodes. xml\_length returns the number of children. xml\_parent returns the parent node, xml\_parents returns all parents up to the root. xml\_siblings returns all nodes at the same level. xml\_child makes it easy to specify a specific child to return.</p>
</div>
<div data-bbox="159 503 207 519" data-label="Section-Header"><b>Usage</b></div>
<div data-bbox="187 527 519 753" data-label="Text">
<pre>xml\_children(x)

xml\_child(x, search = 1, ns = xml\_ns(x))

xml\_contents(x)

xml\_parents(x)

xml\_siblings(x)

xml\_parent(x)

xml\_length(x, only\_elements = TRUE)

xml\_root(x)</pre>
</div>
<div data-bbox="159 771 243 787" data-label="Section-Header"><b>Arguments</b></div>
<div data-bbox="187 795 840 860" data-label="Text">
<p>x                    A document, node, or node set.</p>
<p>search              For xml\_child, either the child number to return (by position), or the name of the child node to return. If there are multiple child nodes with the same name, the first will be returned</p>
</div>

- `ns`                    Optionally, a named vector giving prefix-url pairs, as produced by `xml_ns()`. If provided, all names will be explicitly qualified with the `ns` prefix, i.e. if the element `bar` is defined in namespace `foo`, it will be called `foo:bar`. (And similarly for attributes). Default namespaces must be given an explicit name. The `ns` is ignored when using `xml_name<-()` and `xml_set_name()`.
- `only_elements`        For `xml_length`, should it count all children, or just children that are elements (the default)?

**Value**

A node or nodeset (possibly empty). Results are always de-duplicated.

**Examples**

```
x <- read_xml("<foo> <bar><boo /></bar> <baz/> </foo>")
xml_children(x)
xml_children(xml_children(x))
xml_siblings(xml_children(x)[[1]])

# Note the each unique node only appears once in the output
xml_parent(xml_children(x))

# Mixed content
x <- read_xml("<foo> a <b/> c <d>e</d> f</foo>")
# Children gets the elements, contents gets all node types
xml_children(x)
xml_contents(x)

xml_length(x)
xml_length(x, only_elements = FALSE)

# xml_child makes it easier to select specific children
xml_child(x)
xml_child(x, 2)
xml_child(x, "baz")
```

---

xml\_comment

*Construct a comment node*

---

**Description**

Construct a comment node

**Usage**

```
xml_comment(content)
```

**Arguments**

`content`                The comment content

**Examples**

```
x <- xml_new_document()
r <- xml_add_child(x, "root")
xml_add_child(r, xml_comment("Hello!"))
as.character(x)
```

---

xml\_dtd

*Construct a document type definition*


---

**Description**

This is used to create simple document type definitions. If you need to create a more complicated definition with internal subsets it is recommended to parse a string directly with `read_xml()`.

**Usage**

```
xml_dtd(name = "", external_id = "", system_id = "")
```

**Arguments**

name	The name of the declaration
external_id	The external ID of the declaration
system_id	The system ID of the declaration

**Examples**

```
r <- xml_new_root(
  xml_dtd(
    "html",
    "-//W3C//DTD XHTML 1.0 Transitional//EN",
    "http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd"
  )
)

# Use read_xml directly for more complicated DTD
d <- read_xml(
  '<!DOCTYPE doc [
<!ELEMENT doc (#PCDATA)>
<!ENTITY foo " test ">
]>
<doc>This is a valid document &foo; !</doc>'
)
```



---

xml_find_all	<i>Find nodes that match an xpath expression.</i>
--------------	---------------------------------------------------

---

## Description

Xpath is like regular expressions for trees - it's worth learning if you're trying to extract nodes from arbitrary locations in a document. Use `xml_find_all` to find all matches - if there's no match you'll get an empty result. Use `xml_find_first` to find a specific match - if there's no match you'll get an `xml_missing` node.

## Usage

```
xml_find_all(x, xpath, ns = xml_ns(x), ...)

## S3 method for class 'xml_nodeset'
xml_find_all(x, xpath, ns = xml_ns(x), flatten = TRUE, ...)

xml_find_first(x, xpath, ns = xml_ns(x))

xml_find_num(x, xpath, ns = xml_ns(x))

xml_find_int(x, xpath, ns = xml_ns(x))

xml_find_chr(x, xpath, ns = xml_ns(x))

xml_find_lgl(x, xpath, ns = xml_ns(x))
```

## Arguments

<code>x</code>	A document, node, or node set.
<code>xpath</code>	A string containing an xpath (1.0) expression.
<code>ns</code>	Optionally, a named vector giving prefix-url pairs, as produced by <code>xml_ns()</code> . If provided, all names will be explicitly qualified with the ns prefix, i.e. if the element bar is defined in namespace foo, it will be called <code>foo:bar</code> . (And similarly for attributes). Default namespaces must be given an explicit name. The ns is ignored when using <code>xml_name&lt;-( )</code> and <code>xml_set_name()</code> .
<code>...</code>	Further arguments passed to or from other methods.
<code>flatten</code>	A logical indicating whether to return a single, flattened nodeset or a list of nodesets.

## Value

`xml_find_all` returns a nodeset if applied to a node, and a nodeset or a list of nodesets if applied to a nodeset. If there are no matches, the nodeset(s) will be empty. Within each nodeset, the result will always be unique; repeated nodes are automatically de-duplicated.

xml\_find\_first returns a node if applied to a node, and a nodeset if applied to a nodeset. The output is *always* the same size as the input. If there are no matches, xml\_find\_first will return a missing node; if there are multiple matches, it will return the first only.

xml\_find\_num, xml\_find\_chr, xml\_find\_lgl return numeric, character and logical results respectively.

### Deprecated functions

xml\_find\_one() has been deprecated. Instead use xml\_find\_first().

### See Also

[xml\\_ns\\_strip\(\)](#) to remove the default namespaces

### Examples

```
x <- read_xml("<foo><bar><baz/></bar><baz/></foo>")
xml_find_all(x, ".//baz")
xml_path(xml_find_all(x, ".//baz"))

# Note the difference between .// and //
# // finds anywhere in the document (ignoring the current node)
# .// finds anywhere beneath the current node
(bar <- xml_find_all(x, ".//bar"))
xml_find_all(bar, ".//baz")
xml_find_all(bar, "//baz")

# Find all vs find one -----
x <- read_xml("<body>
  <p>Some <b>text</b>.</p>
  <p>Some <b>other</b> <b>text</b>.</p>
  <p>No bold here!</p>
</body>")
para <- xml_find_all(x, ".//p")

# By default, if you apply xml_find_all to a nodeset, it finds all matches,
# de-duplicates them, and returns as a single nodeset. This means you
# never know how many results you'll get
xml_find_all(para, ".//b")

# If you set flatten to FALSE, though, xml_find_all will return a list of
# nodesets, where each nodeset contains the matches for the corresponding
# node in the original nodeset.
xml_find_all(para, ".//b", flatten = FALSE)

# xml_find_first only returns the first match per input node. If there are 0
# matches it will return a missing node
xml_find_first(para, ".//b")
xml_text(xml_find_first(para, ".//b"))

# Namespaces -----
# If the document uses namespaces, you'll need use xml_ns to form
```

```
# a unique mapping between full namespace url and a short prefix
x <- read_xml('
<root xmlns:f = "http://foo.com" xmlns:g = "http://bar.com">
  <f:doc><g:baz /></f:doc>
  <f:doc><g:baz /></f:doc>
</root>
')
xml_find_all(x, "../f:doc")
xml_find_all(x, "../f:doc", xml_ns(x))
```

---

xml_name	<i>The (tag) name of an xml element.</i>
----------	------------------------------------------

---

### Description

The (tag) name of an xml element.  
 Modify the (tag) name of an element

### Usage

```
xml_name(x, ns = character())

xml_name(x, ns = character()) <- value

xml_set_name(x, value, ns = character())
```

### Arguments

x	A document, node, or node set.
ns	Optionally, a named vector giving prefix-url pairs, as produced by <a href="#">xml_ns()</a> . If provided, all names will be explicitly qualified with the ns prefix, i.e. if the element bar is defined in namespace foo, it will be called foo:bar. (And similarly for attributes). Default namespaces must be given an explicit name. The ns is ignored when using <a href="#">xml_name&lt;-()</a> and <a href="#">xml_set_name()</a> .
value	a character vector with replacement name.

### Value

A character vector.

### Examples

```
x <- read_xml("<bar>123</bar>")
xml_name(x)

y <- read_xml("<bar><baz>1</baz>abc<foo /></bar>")
z <- xml_children(y)
xml_name(xml_children(y))
```

---

xml_new_document	<i>Create a new document, possibly with a root node</i>
------------------	---------------------------------------------------------

---

### Description

xml\_new\_document creates only a new document without a root node. In most cases you should instead use xml\_new\_root, which creates a new document and assigns the root node in one step.

### Usage

```
xml_new_document(version = "1.0", encoding = "UTF-8")

xml_new_root(
  .value,
  ...,
  .copy = inherits(.value, "xml_node"),
  .version = "1.0",
  .encoding = "UTF-8"
)
```

### Arguments

version	The version number of the document.
encoding	The character encoding to use in the document. The default encoding is 'UTF-8'. Available encodings are specified at <a href="http://xmlsoft.org/html/libxml-encoding.html#xmlCharEncoding">http://xmlsoft.org/html/libxml-encoding.html#xmlCharEncoding</a> .
.value	node to insert.
...	If named attributes or namespaces to set on the node, if unnamed text to assign to the node.
.copy	whether to copy the .value before replacing. If this is FALSE then the node will be moved from it's current location.
.version	The version number of the document, passed to xml_new_document(version).
.encoding	The encoding of the document, passed to xml_new_document(encoding).

### Value

A xml\_document object.

---

xml_ns	<i>XML namespaces.</i>
--------	------------------------

---

### Description

xml\_ns extracts all namespaces from a document, matching each unique namespace url with the prefix it was first associated with. Default namespaces are named d1, d2 etc. Use xml\_ns\_rename to change the prefixes. Once you have a namespace object, you can pass it to other functions to work with fully qualified names instead of local names.

### Usage

```
xml_ns(x)

xml_ns_rename(old, ...)
```

### Arguments

x	A document, node, or node set.
old, ...	An existing xml_namespace object followed by name-value (old prefix-new prefix) pairs to replace.

### Value

A character vector with class xml\_namespace so the default display is a little nicer.

### Examples

```
x <- read_xml('
<root>
  <doc1 xmlns = "http://foo.com"><baz /></doc1>
  <doc2 xmlns = "http://bar.com"><baz /></doc2>
</root>
')
xml_ns(x)

# When there are default namespaces, it's a good idea to rename
# them to give informative names:
ns <- xml_ns_rename(xml_ns(x), d1 = "foo", d2 = "bar")
ns

# Now we can pass ns to other xml function to use fully qualified names
baz <- xml_children(xml_children(x))
xml_name(baz)
xml_name(baz, ns)

xml_find_all(x, "//baz")
xml_find_all(x, "//foo:baz", ns)
```

```
str(as_list(x))
str(as_list(x, ns))
```

---

xml\_ns\_strip                      *Strip the default namespaces from a document*

---

### Description

Strip the default namespaces from a document

### Usage

```
xml_ns_strip(x)
```

### Arguments

x                      A document, node, or node set.

### Examples

```
x <- read_xml(
  "<foo xmlns = 'http://foo.com'>
  <baz/>
  <bar xmlns = 'http://bar.com'>
  <baz/>
  </bar>
  </foo>"
)
# Need to specify the default namespaces to find the baz nodes
xml_find_all(x, "//d1:baz")
xml_find_all(x, "//d2:baz")

# After stripping the default namespaces you can find both baz nodes directly
xml_ns_strip(x)
xml_find_all(x, "//baz")
```

---

xml\_path                      *Retrieve the xpath to a node*

---

### Description

This is useful when you want to figure out where nodes matching an xpath expression live in a document.

### Usage

```
xml_path(x)
```

**Arguments**

`x` A document, node, or node set.

**Value**

A character vector.

**Examples**

```
x <- read_xml("<foo><bar><baz /></bar><baz /></foo>")
xml_path(xml_find_all(x, ".//baz"))
```

---

xml\_replace

*Modify a tree by inserting, replacing or removing nodes*

---

**Description**

`xml_add_sibling()` and `xml_add_child()` are used to insert a node as a sibling or a child. `xml_add_parent()` adds a new parent in between the input node and the current parent. `xml_replace()` replaces an existing node with a new node. `xml_remove()` removes a node from the tree.

**Usage**

```
xml_replace(.x, .value, ..., .copy = TRUE)
```

```
xml_add_sibling(.x, .value, ..., .where = c("after", "before"), .copy = TRUE)
```

```
xml_add_child(.x, .value, ..., .where = length(xml_children(.x)), .copy = TRUE)
```

```
xml_add_parent(.x, .value, ...)
```

```
xml_remove(.x, free = FALSE)
```

**Arguments**

`.x` a document, node or nodeset.

`.value` node to insert.

`...` If named attributes or namespaces to set on the node, if unnamed text to assign to the node.

`.copy` whether to copy the `.value` before replacing. If this is `FALSE` then the node will be moved from it's current location.

`.where` to add the new node, for `xml_add_child` the position after which to add, use `0` for the first child. For `xml_add_sibling` either `"before"` or `"after"` indicating if the new node should be before or after `.x`.

`free` When removing the node also free the memory used for that node. Note if you use this option you cannot use any existing objects pointing to the node or its children, it is likely to crash R or return garbage.

## Details

Care needs to be taken when using `xml_remove()`,

---

<code>xml_serialize</code>	<i>Serializing XML objects to connections.</i>
----------------------------	------------------------------------------------

---

## Description

Serializing XML objects to connections.

## Usage

```
xml_serialize(object, connection, ...)
```

```
xml_unserialize(connection, ...)
```

## Arguments

<code>object</code>	R object to serialize.
<code>connection</code>	an open <a href="#">connection</a> or (for <code>serialize</code> ) <code>NULL</code> or (for <code>unserialize</code> ) a raw vector (see ‘Details’).
<code>...</code>	Additional arguments passed to <a href="#">read_xml()</a> .

## Value

For `serialize`, `NULL` unless `connection = NULL`, when the result is returned in a raw vector.

For `unserialize` an R object.

## Examples

```
library(xml2)
x <- read_xml("<a>
  <b><c>123</c></b>
  <b><c>456</c></b>
</a>")

b <- xml_find_all(x, "//b")
out <- xml_serialize(b, NULL)
xml_unserialize(out)
```



---

xml_set_namespace	<i>Set the node's namespace</i>
-------------------	---------------------------------

---

**Description**

The namespace to be set must be already defined in one of the node's ancestors.

**Usage**

```
xml_set_namespace(.x, prefix = "", uri = "")
```

**Arguments**

.x	a node
prefix	The namespace prefix to use
uri	The namespace URI to use

**Value**

the node (invisibly)

---

xml_structure	<i>Show the structure of an html/xml document.</i>
---------------	----------------------------------------------------

---

**Description**

Show the structure of an html/xml document without displaying any of the values. This is useful if you want to get a high level view of the way a document is organised. Compared to `xml_structure`, `html_structure` prints the id and class attributes.

**Usage**

```
xml_structure(x, indent = 2, file = "")
```

```
html_structure(x, indent = 2, file = "")
```

**Arguments**

x	HTML/XML document (or part there of)
indent	Number of spaces to indent
file	A <a href="#">connection</a> , or a character string naming the file to print to. If "" (the default), <code>cat</code> prints to the standard output connection, the console unless redirected by <a href="#">sink</a> . If it is " cmd", the output is piped to the command given by 'cmd', by opening a pipe connection.

**Examples**

```
xml_structure(read_xml("<a><b><c/><c/></b><d/></a>"))

rproj <- read_html(system.file("extdata", "r-project.html", package = "xml2"))
xml_structure(rproj)
xml_structure(xml_find_all(rproj, ".//p"))

h <- read_html("<body><p id = 'a'></p><p class = 'c d'></p></body>")
html_structure(h)
```

---

xml\_text

*Extract or modify the text*


---

**Description**

xml\_text returns a character vector, xml\_double returns a numeric vector, xml\_integer returns an integer vector.

**Usage**

```
xml_text(x, trim = FALSE)

xml_text(x) <- value

xml_set_text(x, value)

xml_double(x)

xml_integer(x)
```

**Arguments**

x	A document, node, or node set.
trim	If TRUE will trim leading and trailing spaces.
value	character vector with replacement text.

**Value**

A character vector, the same length as x.

**Examples**

```
x <- read_xml("<p>This is some text. This is <b>bold!</b></p>")
xml_text(x)
xml_text(xml_children(x))

x <- read_xml("<x>This is some text. <x>This is some nested text.</x></x>")
xml_text(x)
```

```
xml_text(xml_find_all(x, "//x"))

x <- read_xml("<p> Some text </p>")
xml_text(x, trim = TRUE)

# xml_double() and xml_integer() are useful for extracting numeric attributes
x <- read_xml("<plot><point x='1' y='2' /><point x='2' y='1' /></plot>")
xml_integer(xml_find_all(x, "//@x"))
```

---

xml_type	<i>Determine the type of a node.</i>
----------	--------------------------------------

---

**Description**

Determine the type of a node.

**Usage**

```
xml_type(x)
```

**Arguments**

x                    A document, node, or node set.

**Examples**

```
x <- read_xml("<foo> a <b /> <![CDATA[ blah]]></foo>")
xml_type(x)
xml_type(xml_contents(x))
```

---

xml_url	<i>The URL of an XML document</i>
---------	-----------------------------------

---

**Description**

This is useful for interpreting relative urls with [url\\_relative\(\)](#).

**Usage**

```
xml_url(x)
```

**Arguments**

x                    A node or document.

**Value**

A character vector of length 1. Returns NA if the name is not set.

**Examples**

```
catalog <- read_xml(xml2_example("cd_catalog.xml"))
xml_url(catalog)

x <- read_xml("<foo/>")
xml_url(x)
```

---

xml_validate	<i>Validate XML schema</i>
--------------	----------------------------

---

**Description**

Validate an XML document against an XML 1.0 schema.

**Usage**

```
xml_validate(x, schema)
```

**Arguments**

x	A document, node, or node set.
schema	an XML document containing the schema

**Value**

TRUE or FALSE

**Examples**

```
# Example from https://msdn.microsoft.com/en-us/library/ms256129(v=vs.110).aspx
doc <- read_xml(system.file("extdata/order-doc.xml", package = "xml2"))
schema <- read_xml(system.file("extdata/order-schema.xml", package = "xml2"))
xml_validate(doc, schema)
```

# Index

as\_list, 3  
as\_xml\_document, 4

class(), 3  
comment(), 3  
connection, 24, 25  
curl::handle(), 4  
curl\_download, 5

dim(), 3  
dimnames(), 3  
download.file(), 4  
download\_html (download\_xml), 4  
download\_xml, 4

html\_structure (xml\_structure), 25  
htr::set\_config(), 8  
htr::user\_agent(), 8  
htr::with\_config(), 8

multi\_download, 5

names(), 3

read\_html (read\_xml), 6  
read\_xml, 6  
read\_xml(), 24  
row.names(), 3

sink, 25

tsp(), 3

url\_absolute, 8  
url\_escape, 9  
url\_parse, 10  
url\_relative (url\_absolute), 8  
url\_relative(), 27  
url\_unescape (url\_escape), 9

write\_html (write\_xml), 10

write\_xml, 10

xml2\_example, 11  
xml\_add\_child (xml\_replace), 23  
xml\_add\_parent (xml\_replace), 23  
xml\_add\_sibling (xml\_replace), 23  
xml\_attr, 12  
xml\_attr<- (xml\_attr), 12  
xml\_attrs (xml\_attr), 12  
xml\_attrs<- (xml\_attr), 12  
xml\_cdata, 13  
xml\_child (xml\_children), 14  
xml\_children, 14  
xml\_comment, 15  
xml\_contents (xml\_children), 14  
xml\_double (xml\_text), 26  
xml\_dtd, 16  
xml\_find\_all, 17  
xml\_find\_chr (xml\_find\_all), 17  
xml\_find\_first (xml\_find\_all), 17  
xml\_find\_int (xml\_find\_all), 17  
xml\_find\_lgl (xml\_find\_all), 17  
xml\_find\_num (xml\_find\_all), 17  
xml\_find\_one (xml\_find\_all), 17  
xml\_has\_attr (xml\_attr), 12  
xml\_integer (xml\_text), 26  
xml\_length (xml\_children), 14  
xml\_name, 19  
xml\_name<- (xml\_name), 19  
xml\_new\_document, 20  
xml\_new\_root (xml\_new\_document), 20  
xml\_ns, 21  
xml\_ns(), 3, 12, 15, 17, 19  
xml\_ns\_rename (xml\_ns), 21  
xml\_ns\_strip, 22  
xml\_ns\_strip(), 18  
xml\_parent (xml\_children), 14  
xml\_parents (xml\_children), 14  
xml\_path, 22  
xml\_remove (xml\_replace), 23

`xml_replace`, [23](#)  
`xml_root(xml_children)`, [14](#)  
`xml_serialize`, [24](#)  
`xml_set_attr(xml_attr)`, [12](#)  
`xml_set_attrs(xml_attr)`, [12](#)  
`xml_set_name(xml_name)`, [19](#)  
`xml_set_name()`, [3](#), [12](#), [15](#), [17](#), [19](#)  
`xml_set_namespace`, [25](#)  
`xml_set_text(xml_text)`, [26](#)  
`xml_siblings(xml_children)`, [14](#)  
`xml_structure`, [25](#)  
`xml_text`, [26](#)  
`xml_text<-(xml_text)`, [26](#)  
`xml_type`, [27](#)  
`xml_unserialize(xml_serialize)`, [24](#)  
`xml_url`, [9](#), [27](#)  
`xml_validate`, [28](#)