



What you need to know about the bidi algorithm and inline markup

on this page: [definitions](#) - [visual vs logical](#) - [how the algorithm works](#) - [where help is needed](#) - [numbers are special](#) - [the bidi override](#)

This tutorial first describes some of the basic principles underlying how the Unicode bidirectional algorithm works. Then it looks at some of the more common scenarios where the bidi algorithm requires help through the addition of markup or control codes.

Although we try to take a markup independent view here, most of the examples will use XHTML, since it is widely recognizable. For advice relative to a specific markup language see the sidebar.

definitions

Bidirectional text is commonplace in right-to-left scripts such as Arabic, Hebrew, Syriac, and Thaana. Numerous different languages are written with these scripts.

Any embedded text from a left-to-right script and all numbers progress visually *left-to-right* within the general right-to-left visual flow. (Of course, the English text on this page also contains bidirectional text where it includes Arabic and Hebrew examples.)

We will use the term **bidi** to mean 'bidirectional'. We will also use **RTL** for 'right-to-left' and **LTR** for 'left-to-right'.

visual vs. logical order

Visual ordering of text was a common way of representing Hebrew in HTML on old user agents that didn't support the Unicode bidi algorithm. It still persists to a degree today, out of habit. Characters making up the text were stored in the source code in the same order you would see them displayed on screen when looking from left to right.

For example, take some Hebrew text with mixed directionality.

פעילות הבינאום, W3C

If you looked at the characters in memory, one by one, you would see the following for logically and visually ordered text. This also represents the typing order, so visual Hebrew text would have to be typed backwards.

Logical order	Visual order
05E4 פ HEBREW LETTER PE	0057 W LATIN CAPITAL LETTER W
05E2 ץ HEBREW LETTER AYIN	0033 3 DIGIT THREE
05D9 ך HEBREW LETTER YOD	0043 C LATIN CAPITAL LETTER C
05DC ל HEBREW LETTER LAMED	0020 SPACE
05D5 ם HEBREW LETTER VAV	002C , COMMA
05EA ת HEBREW LETTER TAV	05DD ם HEBREW LETTER FINAL MEM
0020 SPACE	05D5 ם HEBREW LETTER VAV
05D4 ה HEBREW LETTER HE	05D0 א HEBREW LETTER ALEF
05D1 ב HEBREW LETTER BET	05E0 נ HEBREW LETTER NUN
05D9 ך HEBREW LETTER YOD	05D9 ך HEBREW LETTER YOD
05E0 נ HEBREW LETTER NUN	05D1 ב HEBREW LETTER BET
05D0 א HEBREW LETTER ALEF	05D4 ה HEBREW LETTER HE
05D5 ם HEBREW LETTER VAV	0020 SPACE
05DD ם HEBREW LETTER FINAL MEM	05EA ת HEBREW LETTER TAV
002C , COMMA	05D5 ם HEBREW LETTER VAV
0020 SPACE	05DC ל HEBREW LETTER LAMED
0057 W LATIN CAPITAL LETTER W	05D9 ך HEBREW LETTER YOD
0033 3 DIGIT THREE	05E2 ץ HEBREW LETTER AYIN
0043 C LATIN CAPITAL LETTER C	05E4 פ HEBREW LETTER PE

Visual ordering isn't really seen much for Arabic. Since the Arabic letters are all joined up there was a stronger motivation on the part of Arabic implementers to enable the logical ordering approach.

Visual ordering of flowing text requires the author to disable any line wrapping and to explicitly add line-breaks and right-align text within block elements and table cells. Then the text needs to be served in an encoding that prevents the application of the Unicode bidi algorithm in later browsers. Here is an example written in HTML:

```
<table width="50%"><tr><td align="right" nowrap>
,INRIA מחליפה את שירותי הארחה באירופה מ-INRIA
W3C ממקמת בצרפת, ל-ERCIM. השינוי מאפשר ל-W3C
להעמיק את קשרי המחקר ברחבי אירופה, תוך שמירה
על הקשר ההיסטורי החזק עם INRIA, אחד ממייסדי
.ERCIM. השינוי יתבצע ב 1 לינואר 2003.
</td></tr></table>
```

(This is actually a fairly clean implementation. For example you can also find such things as right-aligned paragraphs with `<nobr>...</nobr>` tags around each line. If your window is too narrow, the beginning of each line disappears off the right side of the browser.)

The result is very fragile code that is difficult to maintain. For example, apart from the difficulty of typing the Hebrew backwards, if you wanted to add a few words on the second line of this paragraph, you would have to readjust all the following line breaks. You would also have to add and maintain separate spans of link or emphasis markup for any marked up text that wrapped onto another line.

Note: Visual ordering can also cause problems at a higher level. For example, it requires the order of table columns to be manually reversed when translating into another language. Line breaks will also need to be manually re-flowed if the page geometry changes. Etc.

Logical ordering is a much better approach. In this approach text is stored in memory in the order

in which it is normally typed (and usually pronounced). The Unicode bidirectional algorithm then produces the reordering required for correct visual display.

This makes it almost trivial to create long paragraphs of flowing text that automatically wrap to the width of the block element. It also makes it much easier to use such things as screen readers.

The bidi algorithm works on logically ordered text. If you prefer to use visual ordering there is no point reading further (except that you could make your life a whole lot easier).

how the bidi algorithm works

Here we introduce some important basic concepts. If it seems boring, stick with it because without understanding this stuff properly you'll get lost when you need to write marked up bidi text.

Directional typing of characters

We already know that a sequence of Latin characters is rendered (ie. displayed) one after the other from left to right (we can see that on this page). On the other hand, the bidi algorithm will render a sequence of strongly typed RTL characters one after the other from right to left.

This works because each character in Unicode has an associated directional property. Most letters are **strongly typed** as LTR. Letters from bidirectional scripts are strongly typed as RTL.

When text with different directionality is mixed inline, the bidi algorithm renders each sequence of characters with the same directionality as a separate **directional run**. So in the following example there are three directional runs:

bahrain مصر kuwait

Directional runs are ordered according to the overall context. In the example above, which has an overall context of LTR, you would read 'bahrain', then 'مصر' (RTL), then 'kuwait'. Note that you don't need any markup or styling to make this happen.

Directional context

The result of the bidirectional algorithm will depend on the overall directional context of the paragraph, block or page in which it is applied.

In XHTML that context would normally be expressed by adding nothing to the `html` tag, in which case the default directionality is LTR, or adding `dir="rtl"` to the `html` tag, in which case all the elements in the document will inherit a context of RTL. The context may be changed later by using the `dir` attribute on the relevant block element.

Neutral characters

Spaces and punctuation do not have LTR and RTL forms in Unicode, because they may be used in either type of script. They are therefore classed as **neutral**.

This is where things begin to get interesting. When the bidi algorithm encounters characters with neutral directional properties (such as spaces and punctuation) it works out how to handle them by looking at the surrounding characters.

4

A neutral character between two strongly typed RTL characters will be treated as a RTL character itself, and will have the effect of extending the directional run. This is why the three arabic words in this LTR phrase are read from right to left. (The first Arabic word you read is مفتاح then معايير then الويب.)

The title is مفتاح معايير الويب in Arabic.

Note that you still don't need any markup or styling for this. There are still only three directional runs here.

The really interesting part comes when a space or punctuation falls between two strongly typed characters with *different* directionality. In such a case they will be treated as if they have the directionality of the *overall paragraph or context*. Even if there are several neutral characters between the two different strongly typed characters, they will all be treated in the same way.

This has the effect of creating a boundary between directional runs.

where the algorithm needs help

The bidi algorithm will handle text perfectly well in most situations, and typically no special markup or other device is needed other than to set the overall direction for the document. You would be very lucky, however, if you got off that easily all the time.

Neutrals that get misplaced

Let's type some punctuation at the end of the Arabic phrase in the last example. By default we will see the following:

The title is "مفتاح معايير الويب!" in Arabic.

The quotation marks are OK, but the exclamation mark is in the wrong position. It should appear at the end of the Arabic text, ie. to the left, like this:

The title is "مفتاح معايير الويب!" in Arabic.

Given our understanding of the bidi algorithm we can easily understand why this happened. Because the exclamation mark was typed in between the last RTL letter 'ب' (on the left) and the LTR letter 'i' (of the word 'in') its directionality is determined by the overall context of the paragraph (here LTR). Note that it makes no difference that there are actually two punctuation characters and a space in this position - they are all neutrals and so are all affected the same way. Because the exclamation mark is seen as LTR it joins the directional run that includes the text 'in Arabic'.

So how to get the punctuation in the right place? In situations such as this you would typically have surrounded the Arabic quotation with markup - either to denote it as a quotation, or to simply apply language information. In this case, there is a simple remedy. Apply an attribute to the tag to change the directional flow within to RTL.

Here is how that might look in XHTML.

```
The title is "<span dir="rtl" lang="ar" xml:lang="ar">مفتاح معايير الويب!</span>"  
in Arabic.
```

Note carefully how the span tag falls *inside* the quote marks - these are part of the surrounding English text.

Another possibility would be to type an invisible, strongly-typed RTL character after the exclamation mark. That way the exclamation mark would be interpreted as RTL and join the Arabic directional run.

It just so happens that there is such a character - the Unicode character U+200F, called the RIGHT-TO-LEFT MARK (RLM). There is a similar character, U+200E, called the LEFT-TO-RIGHT MARK (LRM).

Because the character is invisible you may prefer to actually type in a numeric character reference (`‏`) or, if available, a character entity (such as `&r1m;` in XHTML). In the following example `&r1m;` has been added after the exclamation mark and the result looks fine:

The title is "مفتاح معايير الويب!" in Arabic.

Neutrals with an identity crisis

In our next example the list order is incorrect because the first two Arabic words should be reversed and the intervening comma, which is part of the English text, should appear immediately to the right of the first word.

The names of these states in Arabic are الكويت, مصر, البحرين respectively.

What was wanted was:

The names of these states in Arabic are مصر, البحرين and الكويت respectively.

The reason for the failure is that, with a strongly typed right-to-left (RTL) character on either side, the bidirectional algorithm sees the neutral comma as part of the Arabic text. In fact it is part of the English text, and should mark the boundary of two directional runs in Arabic.

Whereas in the previous section the neutral character thought it was part of the overall context, but wasn't, in this section the neutral character thinks it is part of the directional run, when it is really part of the overall context! No-one said life was simple...

An elegant solution is to use another invisible Unicode character, this time the LEFT-TO-RIGHT MARK, next to the comma. This puts our neutral punctuation between strongly typed RTL and LTR characters and forces it to take on the directionality of the overall context, which is the English LTR flow. This breaks the Arabic words into two separate directional runs, which are ordered LTR in accordance with the prevailing direction of the paragraph.

Again, you may prefer to use an NCR (`‎`) or a character entity (such as `&l1r;`) if available.

Putting markup around the comma is a bit like cracking an egg with a hammer in this case.

Once again, together

The examples we have used so far have been English and LTR based. The same principles apply for RTL text in languages such as Hebrew and Arabic. Lets see one more example.

Here's what we want to see, in a paragraph with its base direction set to RTL.

.ERCIM - מעביר את שירותי הארחה באירופה ל (World Wide Web Consortium) W3C

Unfortunately, on its own the bidirectional algorithm creates a real mess of this.

.ERCIM - ל - W3C (World Wide Web Consortium) מעביר את שירותי הארחה באירופה

It looks like this is going to be incredibly complicated to sort out, but actually the solution is trivial. Just insert an RLM after 'W3C' and you're done. It's really that simple!

If you're not convinced, here's the explanation. The RLM after 'W3C' makes this piece of LTR text a separate directional flow from the text in parentheses that follows it. Remember that flows will be ordered right-to-left because this is the overall paragraph direction. Since the 'W3C' text was typed in first, that now appears farthest to the right. The parenthesis is now between strongly typed RTL and LTR characters, and so also takes on the directionality of the overall paragraph, RTL. So that comes next. Then comes the unbroken LTR directional flow which is the stuff inside the parentheses.

Note: To better understand the sequence of characters in the examples above it may help to move your cursor through the source text in an editor, rather than just staring at the screen.

(The changed glyph shape of the right-most parenthesis comes automatically. The glyph used for these 'mirrored characters' changes according to its directionality. It's still the same character.)

Nesting directional runs

The Unicode bidi algorithm and the directional markers work quite well when there is only a single level of mixed text. If you have a situation where there are two or more nested levels of directional text you will need a different solution. Here is an example of incorrectly ordered text:

The title says "פעילות הבינאום, W3C" in Hebrew.

The order of the two Hebrew words is correct, but the text 'W3C' should appear on the left hand side of the quotation and the comma should appear between the Hebrew text and 'W3C'. In other words, the desired result is:

The title says "W3C ,פעילות הבינאום" in Hebrew.

The problem arises because the directional flows are being ordered according to the LTR context of the paragraph. Inside the Hebrew quotation, however, the correct default ordering should be RTL.

To resolve this problem we need to open a new **embedding level**. In XHTML this would be done by enclosing the quotation in markup and assigning it a directionality of RTL using the `dir` attribute.

The title says "`W3C ,פעילות הבינאום`" in Hebrew.

In markup languages other than XHTML/HTML you may find a similar attribute to which you can apply styling to achieve the correct effect. If you don't have such an attribute you may have to resort to individually styling the appropriate inline markup, but it would probably be better to lobby your markup developer to provide you with one.

There are Unicode control characters you could use to achieve the same result, but because they create states with invisible boundaries this is [not recommended](#).

numbers are a little special

Numbers in RTL scripts run left-to-right within the right-to-left flow, but they are handled a little

differently than words by the bidi algorithm. They are said to have **weak directionality**. The following two examples illustrate this difference. Comparison of the two cases shows the words on either side of the fourth item in the sequence have been reversed. The only difference between the two sentences in memory is the use of '1234' versus 'four'.

one two ثلاثة four خمسة

one two خمسة 1234 ثلاثة

In the first example the letters in the word 'four' are strongly typed and therefore break the two Arabic words into separate directional runs, ordered left to right as per the paragraph context.

In the second example, the weakly typed number '1234' is seen as part of the Arabic text, so the two Arabic words are treated as part of the same directional run - even though the sequence of digits runs LTR on screen.

Note: In some versions of Netscape and Mozilla browsers, the fourth item in the second row will be rendered as ١٢٣٤ rather than 1234. It means the same. Note that Mozilla 1.4 reinstates European digits as the default.

This only happens in RTL text.

Sounds complicated? Don't worry, usually the bidi algorithm will just take care of things for you. I really only included this section for those who notice the difference and wonder what's going on.

Note also that alongside a number certain otherwise neutral characters, such as currency symbols, will be treated as part of the number rather than a neutral .

overriding the algorithm

There may be occasions where you don't want the bidi algorithm to do its reordering work at all. In these cases you need some additional markup to surround the text you want left unordered.

In XHTML 1.0 this is achieved using the inline `bdo` element. In XHTML 2 it will probably be implemented as a value of `rl0` or `lr0` on the `direction` attribute, enabling it to be applied to any element. Again, there are Unicode control characters you could use to achieve the same result, but because they create states with invisible boundaries this is **not recommended**.

Examples that show the characters as ordered in memory use the `bdo` tag to achieve that effect. For example, to show the underlying sequence of characters for:

פעילות הבינאום, W3C

we would use the following markup in XHTML 1.0:

```
<p><bdo dir="ltr">W3C , פעילות הבינאום</bdo></p>
```

The result, drawn left-to-right, would be:

W3C, מואניבה תוליעפ

further reading

- [Bidi techniques in HTML](http://www.w3.org/International/geo/html-tech/tech-bidi.html) (in devt)
<http://www.w3.org/International/geo/html-tech/tech-bidi.html>

- [The Unicode Bidirectional Algorithm](http://www.unicode.org/reports/tr9/) *http://www.unicode.org/reports/tr9/*
- [Non-Latin scripts tutorial](http://people.w3.org/rishida/scripts/tutorial/) *http://people.w3.org/rishida/scripts/tutorial/*
- [Other W3C I18N resources relating to Bidirectional text](http://www.w3.org/International/resource-index.html#bidi)
http://www.w3.org/International/resource-index.html#bidi

acknowledgements

Many thanks to Martin Dürst for many useful comments on an early draft of this document.

Author: Richard Ishida, W3C.

Content created 29 September, 2003. Last update 2004-09-20 15:48 GMT

Copyright © 2004 W3C[®] (MIT, ERCIM, Keio), All Rights Reserved. W3C [liability](#), [trademark](#), [document use](#) and [software licensing](#) rules apply. Your interactions with this site are in accordance with our [public](#) and [Member](#) privacy statements.

