

Proposed Draft Technical Report ISO/IEC PDTR 19075-1	
Date: 2010-02-05	Reference number: ISO/JTC 1/SC 32N1969
Supersedes document – 32N1830	

THIS DOCUMENT IS STILL UNDER STUDY AND SUBJECT TO CHANGE. IT SHOULD NOT BE USED FOR REFERENCE PURPOSES.

ISO/IEC JTC 1/SC 32 Data Management and Interchange	Circulated to P- and O-members, and to technical committees and organizations in liaison for voting (P-members only) by:
Secretariat: USA (ANSI)	2010-05-05
	Please return all votes and comments in electronic form directly to the SC 32 Secretariat by the due date indicated.

ISO/IEC PDTR 19075-1:2010(E)
 Title: Information technology - Technical reports on ISO/IEC 9075 - Part 1: SQL XQuery regular expression support in SQL/Foundation
 Project: 1.32.09.01.01.00

Introductory note:
 initial committee stage draft; no disposition of comments; this text is sent to NBs for 3 month letter ballot. The ballot starts 2010-02-05.
 Medium: E
 No. of pages: 33

Dr. Timothy Schoechle, Secretary, ISO/IEC JTC 1/SC 32
 Farance Inc *, 3066 Sixth Street, Boulder, CO, United States of America
 Telephone: +1 303-443-5490; E-mail: Timothy@Schoechle.org
 available from the JTC 1/SC 32 WebSite <http://www.jtc1sc32.org/>
 *Farance Inc. administers the ISO/IEC JTC 1/SC 32 Secretariat on behalf of ANSI

WG3:KMG-015
DM32.2-2010-00031

ISO/IEC JTC 1/SC 32

Date: 2010-01-31

PDTR 19075-1:2011(E)

ISO/IEC JTC 1/SC 32/WG 3

The United States of America (ANSI)

Information technology — Database languages — SQL Technical Reports —

Part 1:
XQuery Regular Expression Support in SQL

Technologies de l'information — Langages de base de données — SQL rapports techniques —
Partie 1: le Support des Expressions Régulières d'XQuery en SQL

Document type: Technical Report
Document subtype: Proposed Draft Technical Report (PDTR)
Document stage: (2) PDTR under Consideration
Document language: English

Copyright notice

This ISO document is a working draft or a committee draft and is copyright-protected by ISO. While the reproduction of working drafts or committee drafts in any form for use by participants in the ISO standards development process is permitted without prior permission from ISO, neither this document nor any extract from it may be reproduced, stored or transmitted in any form for any other purpose without prior written permission from ISO.

Requests for permission to reproduce for the purpose of selling it should be addressed as shown below or to ISO's member body in the country of the requester.

*ANSI Customer Service Department
25 West 43rd Street, 4th Floor
New York, NY 10036
Tele: 1-212-642-4980
Fax: 1-212-302-1286
Email: storemanager@ansi.org
Web: www.ansi.org*

Reproduction for sales purposes may be subject to royalty payments or a licensing agreement.

Violators may be prosecuted.

Contents	Page
Foreword.....	v
Introduction.....	vi
1 Scope.....	1
2 Normative references.....	3
3 XQuery regular expressions.....	5
3.1 Matching a specific character.....	5
3.2 Metacharacters and escape sequences.....	6
3.3 Dot.....	7
3.4 Anchors.....	8
3.5 Line terminators.....	8
3.6 Bracket expressions.....	9
3.6.1 Listing characters.....	9
3.6.2 Matching a range.....	10
3.6.3 Negation.....	10
3.6.4 Character class subtraction.....	10
3.7 Alternation.....	10
3.8 Quantifiers.....	11
3.9 Locating a match.....	12
3.10 Capture and back-reference.....	13
3.11 Precedence.....	14
3.12 Modes.....	15
4 Operators using regular expressions.....	17
4.1 LIKE_REGEX.....	17
4.2 OCCURRENCES_REGEX.....	18
4.3 POSITION_REGEX.....	19
4.4 SUBSTRING_REGEX.....	21
4.5 TRANSLATE_REGEX.....	22
Index.....	25

(Blank page)

Foreword

TO BE SUPPLIED.

PDTR 19075-1:2011(E)

Introduction

TO BE SUPPLIED.

Information technology — Database languages — SQL Technical Reports —

Part 1:

XQuery Regular Expression Support in SQL**1 Scope**

This describes the regular expression support in SQL adopted from the regular expression syntax of [XQuery F&O], which is derived from Perl. This paper proposes five operators using this regular expression syntax:

- LIKE_REGEX predicate, to determine the existence of a match to a regular expression.
- OCCURRENCES_REGEX numeric function, to determine the number of matches to a regular expression.
- POSITION_REGEX function, to determine the position of a match.
- SUBSTRING_REGEX function, to extract a substring matching a regular expression.
- TRANSLATE_REGEX function, to perform replacements using a regular expression.

(Blank page)

2 Normative references

- [Foundation:2003] Jim Melton (ed), "ISO International Standard (IS) Database Language SQL - Part 2: SQL/Foundation", ISO/IEC 9075-2:2003
- [Foundation CD] Jim Melton (ed), "Committee Draft (CD) Database Language SQL - Part 2: SQL/Foundation", ISO/IEC JTC1/SC32 WG3:TXL-003 = ANSI INCITS H2-2004-005
- [Foundation WD] Jim Melton (ed), "Working Draft (WD) Database Language SQL - Part 2: SQL/Foundation", ISO/IEC JTC1/SC32 WG3:SIA-003 = ANSI INCITS H2-2003-420
- [MED CD] Jim Melton (ed), "Working Draft (WD) Database Language SQL - Part 9: SQL/MED", ISO/IEC JTC1/SC32 WG3:TXL-025 = ANSI INCITS H2-2005-021
- [SQL/XML WD] Jim Melton (ed.), "Working Draft (WD) XML-Related Specifications (SQL/XML)", ISO/IEC JTC1/SC32 WG3:TXL-010r1 = ANSI INCITS H2-2005-012r1
- [Friedl] Jeffrey E. F. Friedl, "Mastering Regular Expressions", O'Reilly & Associates, 1997
- [Perl 5.8 reg ex] Perl 5.8.0, Perl Manual Pages (perlre), Larry Wall, available at <http://www.perl-doc.com/perl5.8.0/pod/perlre.html>
- [POSIX reg ex] "(Open Group Technical Standard, Issue 6), Standard for Information Technology --- Portable Operating System Interface (POSIX) 2001", IEEE 1003.1-2001, ISBN 0-7381-3009-5, <http://www.opengroup.org/onlinepubs/7908799/xbd/re.html>
- [XML 1.0] "Extensible Markup Language (XML) 1.0 (Second edition)", W3C Recommendation, 6 October 2000, available at <http://www.w3.org/TR/2000/REC-xml-20001006>
- [XML 1.1] "Extensible Markup Language (XML) 1.1", W3C Candidate Recommendation, 15 October 2002, available at <http://www.w3.org/TR/2002/CRxml11-20021015>
- [XML Schema: Datatypes] "XML Schema Part 2: Datatypes", W3C Recommendation, 28 October 2004, available at <http://www.w3.org/TR/2004/REC-xmlschema-2-20041028>
- [XQuery F&O] "XQuery 1.0 and XPath 2.0 Functions and Operators", W3C "last call" working draft, 15 September 2005, available at <http://www.w3.org/TR/2005/WD-xpath-functions-20050915/>
- [Unicode18] Mark Davis, "Unicode technical standard #18, Unicode regular expressions", <http://www.unicode.org/reports/tr18/>

(Blank page)

3 XQuery regular expressions

XQuery regular expression syntax is specified in [XQuery F&O], section 7.6.1, “Regular expression syntax”. This paper references the XQuery specification, with two small modifications (required since character strings in an RDBMS are not necessarily normalized according to XML conventions). The following subsections provide an overview of this syntax.

The XQuery regular expression syntax is itself a modification of another regular expression syntax found in [XML Schema: Datatypes].

This section presents an overview of the capabilities of XQuery regular expression syntax. In the process, this section will illustrate some of the SQL operators. The SQL operators themselves are presented in the next section.

The following discussion does not cover every aspect of XQuery regular expressions; for this, [XQuery F&O] is the reference (though hardly a tutorial; try [Friedl] for a detailed treatment of regular expressions).

3.1 Matching a specific character

Perhaps the most elementary pattern matching requirement is the ability to match a single character or string. For most characters, this is done by simply writing the character in the regular expression. For example, suppose you want to know if a string *S* contains the letters “xyz”. This could be done with the following predicate:

```
S LIKE_REGEX 'xyz'
```

Note that the SQL LIKE predicate would require an exact match for “xyz”. However, the convention with regular expressions is that *S* need only contain a substring that is “xyz”. For example, all of the following values of *S* would yield *True* for the predicate above:

```
xyz
abcxyz123
1 xyz 2 xyz 3 xyz
```

Note that in the last example, there are actually three occurrences of the regular expression “xyz” within the tested value. The user may wish to know the number of occurrences of a match. This can be done with OCCURRENCES_REGEX. For example:

```
OCCURRENCES_REGEX ('xyz' IN '1 xyz 2 xyz 3 xyz') = 3
```

The user might also wish to know the position of a specific match. This can be done using POSITION_REGEX. For example, to learn the starting character position of the second occurrence,

```
POSITION_REGEX ('xyz' IN '1 xyz 2 xyz 3 xyz' OCCURRENCE 2 ) = 9
```

It is also possible to ask for the character position of the first character after the match. For example:

PDTR 19075-1:2011(E)

3.1 Matching a specific character

```
POSITION_REGEX ( AFTER 'xyz' IN '1 xyz 2 xyz 3 xyz' OCCURRENCE 2 ) = 12
```

If AFTER is used and the last character of the subject string is consumed, then the result is the length of the string plus 1 (one):

```
POSITION_REGEX ( AFTER 'xyz' IN 'xyz' ) = 4
```

3.2 Metacharacters and escape sequences

As mentioned, most characters can be matched by simply writing the character in the regular expression. However, certain characters are reserved as *metacharacters*. The complete list of metacharacters is:

```
. \ ? * + { } ( ) | [ ] ^ $
```

The use of each of these metacharacters will be explained later. If you want to match a metacharacter, then you need to use an *escape sequence*, consisting of a backslash (“/”) followed by the metacharacter. For example, to test whether a string contains a dollar sign, you could write

```
S LIKE_REGEX '\$'
```

In particular, the escape sequence representing a backslash is two consecutive backslashes. There are various other defined escape sequences, matching either a single character, or any of a group of characters. The *single character escape sequences* are:

<code>\n</code>	newline (U+000A)
<code>\r</code>	return (U+000D)
<code>\t</code>	tab (U+0009)
<code>\-</code>	minus sign ('-')

The so-called *category escapes* are exemplified by “`\p{L}`” or “`\p{Lu}`”. A category escape begins with “`\p{`” followed by one uppercase letter, optionally a lowercase letter, and then the closing brace. In these example, “`\p{L}`” matches any letter (as defined by Unicode) and “`\p{Lu}`” matches any uppercase letter. Some interesting category escapes are listed below:

<code>\p{L}</code>	Any letter.
<code>\p{Lu}</code>	Any uppercase letter.
<code>\p{Ll}</code>	Any lowercase letter.
<code>\p{Nd}</code>	Any decimal digit.
<code>\p{P}</code>	Any punctuation mark.
<code>\p{Z}</code>	Any separator (space, line, paragraph, <i>etc.</i>).

The complete list of category escapes is found in [XML Schema: Datatypes], section F.1.1, “Character class escapes”.

There are also *complementary category escapes*, which are exemplified by “`\P{L}`” or “`\P{Lu}`”. A complementary category escape matches any character that would not be matched by the corresponding category

escape. The difference is that the (positive) character escape is written with a lowercase “p” whereas the complementary character escape is written with an uppercase “P”.

The so-called *block escapes* match any character in a block of Unicode, that is, a predefined consecutive range of code points. For example, “\p{IsBasicLatin}” matches the ASCII character set. There are also *complementary block escapes*, such as “\P{IsBasicLatin}”, which matches any single character that is not an ASCII character.

Finally, there are the following *multi-character escape sequences*:

- \s As defined by [XML Schema: Datatypes], this escape matches space (U+0020), tab (U+0009), newline (U+000A), or return (U+000D). Since character strings in an RDBMS have not undergone XML line termination normalization, we broaden it to include any character or two-character sequence that is recognized by [Unicode18] as a line terminator. [Subclause 3.5, “Line terminators”](#), discusses this issue further.
- \S Any single character not matched by \s.
- \i Underscore (“_”), colon (“:”) or letter (this is a lot more than just the Latin letters; see [XML 1.0] appendix B, rule [84]).
- \I Any single character not matched by \i.
- \c Any single character matched by NameChar, as defined in [XML 1.0] section 2.3, rule [4].
- \C Any single character not matched by \c.
- \d Any single digit
- \D Any single character not matched by \d.
- \w Any single Unicode character except those classified as “punctuation”, “separator”, or “other”.
- \W The complement of \w.

3.3 Dot

Dot (period, “.”) is a metacharacter that is used to match any single character (the same behavior as “_” in LIKE predicates), or any single character that is not a line terminator. The default is to match anything except a line terminator. The alternative, called *dot-all mode*, is specified using a flag that contains a lowercase “s”.

For example

```
S LIKE_REGEX 'a.b'
```

matches the following:

```
'xa0by'
```

but not the following:

```
'xa
by'
```

because the character between the “a” and the “b” is a line terminator. However, using dot-all mode like this:

3.3 Dot

```
S LIKE_REGEX 'a.b' FLAG 's'
```

would match both examples.

3.4 Anchors

We have seen that regular expressions look for a match anywhere within a string, without needing to match the entire string. But what if you want to require a match of the entire string? For this, you can use *anchors*. The anchors are the metacharacters “^” for the start of a string (or line), and “\$” for the end of a string (or line). For example:

```
S LIKE_REGEX '^xyz$'
```

can only match a string that is precisely 'xyz'.

Anchors may be used separately to require a “begins with” or “end with” match. For example

```
S LIKE_REGEX '^xyz'
```

matches any string that begins with “xyz”, and

```
S LIKE_REGEX 'xyz$'
```

matches any string that ends with “xyz”.

Instead of matching the begin or end of the string, the anchors may be used to anchor a match to the begin or end of a line, by performing the match in *multi-line mode*. Multi-line mode is specified using a flag containing a lowercase “m”. For example:

```
S LIKE_REGEX '^xyz' FLAG 'm'
```

performs an anchored search in multi-line mode, matching any string containing a line that begins with “xyz”. The example above would match the following string:

```
'line one
xyz
line three'
```

3.5 Line terminators

The metacharacters “.”, “^”, and “\$” and the multi-character escape sequences “\s” and “\S” are defined in terms of a “line terminator”. What counts as a line terminator? [XQuery F&O] only recognizes a line feed (U+000A) as a line terminator. This definition works well for XQuery, because XML normalizes the line terminators on various platforms to a line feed.

A closer look shows that XML has two definitions of line handling, in section 2.11, “End-of-line handling”, of [XML 1.0] and [XML 1.1]. So which should we use for SQL?

A first stop in answering this is to look at [SQL/XML WD] Subclause 6.17, “<XML query>”, which requires XML 1.0 as a basic level of support, and permits XML 1.1 support in the form of Feature X211, “XML 1.1

support”. So, we might specify that the character string is normalized according to either XML 1.0 or XML 1.1 as an implementation-defined choice, or perhaps via a conformance feature.

However, some of the line terminators, even in XML 1.0, are two-character sequences. XML normalizes its input, which means that such two-character sequences are converted to a single character. This changes the relative position of every subsequent character, which would cause unexpected results for `POSITION_REGEX`.

Our solution is to look to [Unicode18], a Unicode standard containing guidelines for regular expression processors. This provides a referenceable definition of line terminator that does not require normalizing the subject character string.

3.6 Bracket expressions

So far, we have seen how to match a specific character, or any character from certain predefined sets of characters. Using bracket expressions, you can specify your own group of characters. (XML Schema and XQuery call these *character class expressions*, but the term *bracket expression* is in common use.)

A bracket expression is begun by a left bracket “[” and terminated by a right bracket “]”. Bracket expressions have a different list of special characters, namely

`^ [] \`

For clarity, we will call these *special characters*, in contrast to the metacharacters listed earlier.

3.6.1 Listing characters

If a bracket expression does not contain any of the special characters, then the bracket expression matches any single character that is listed between the brackets. For example,

```
S LIKE_REGEX '[abc]'
```

matches any of the following:

```
'say'
'boy'
'lack'
```

All backslash escape sequences are available for use within a bracket expression. For example, to match either a caret or a backslash, you can use

```
S LIKE_REGEX '[\\^\\]'
```

To match all letters or digits, one might use

```
S LIKE_REGEX '[\p{L}\p{Nd}]'
```

where “`\p{L}`” is the escape matching any letter and “`\p{Nd}`” is the escape matching any digit.

3.6.2 Matching a range

A minus sign “-” is used to specify a character range. For example:

```
S LIKE_REGEX '[sa-my]'
```

matches the lowercase letters “s”, all the letters between “a” and “m” inclusive, and “y”. Ranges are defined in terms of the UCS code point ordering. When there are multiple ranges, the bracket expression matches the union of the ranges. For example:

```
S LIKE_REGEX '[a-me-z]'
```

matches all lowercase letters.

Using a special character in a range is sometimes permitted, but tricky. Rather than present the rules here, our advice is to use a backslash escape if the start or end point of a range must be a special character.

3.6.3 Negation

A caret “^” is a special character when it is the first character of a bracket expression, where it indicates that the set of characters is anything not listed by the following bracket expression. For example:

```
S LIKE_REGEX '[^aj-m]'
```

is *True* if S contains any character that is not “a”, “j”, “k”, “l”, or “m”.

3.6.4 Character class subtraction

A bracket expression may conclude with a minus sign “-” followed by a nested bracket expression. This is called a *character class subtraction*, and indicates that any character matched by the nested bracket expression is to be removed from the set of characters that might be a match. For example:

```
S LIKE_REGEX '[a-z-[m-p]]'
```

matches anything between “a” and “z”, except for the letters between “m” and “p”, inclusive. This example is equivalent to:

```
S LIKE_REGEX '[a-lq-z]'
```

Seemingly you can nest character class subtractions indefinitely. This concludes the presentation of bracket expressions.

3.7 Alternation

You can specify a choice of regular expressions using a vertical bar “|”. For example:

```
S LIKE_REGEX 'a|b'
```

is *True* if S contains either an “a” or a “b”.

Alternation has lower precedence than concatenation. Thus

```
S LIKE_REGEX 'ab|xyz'
```

is *True* if S contains either “ab” or “xyz”. To override this precedence, you can use parentheses, such as this example:

```
S LIKE_REGEX 'a(b|xy)z'
```

The preceding example is *True* if S contains either “abz” or “axyz”.

3.8 Quantifiers

Quantifiers are metacharacters that specify a match for some number of repetitions of a regular expression. There are two sets of quantifiers, the greedy and the reluctant. The *greedy quantifiers* are:

- { n } Exactly n repetitions.
- { n, } n or more repetitions.
- { n,m } Between n and m repetitions, inclusive.
- ? 0 (zero) or 1 (one) repetition; equivalent to {0,1}.
- * 0 (zero) or more repetitions; equivalent to {0,}.
- + 1 (one) or more repetitions; equivalent to {1,}.

The *reluctant quantifiers* are formed by suffixing a question mark to a greedy quantifier. Thus, “*?” is the reluctant form of “*”, and “??” is the reluctant form of “?”. The greedy quantifiers try to match as much as possible, whereas the reluctant quantifiers try to match as little as possible (while still allowing the overall regular expression to match). There is no difference in behavior between the greedy and reluctant quantifiers for LIKE_REGEX. We will look at this difference for the other operators shortly.

Examples:

```
S LIKE_REGEX 'a{3}'
```

is equivalent to

```
S LIKE_REGEX 'aaa'
```

and matches any string containing at least three consecutive instances of “a”. Note that if S contains more than three consecutive instances of “a”, it still matches; to test whether S contains a substring of three consecutive instances of “a” and no more is a lot harder, since you have to also require something other than an “a” at both ends of the substring.

```
S LIKE_REGEX 'ab+c'
```

is equivalent to

3.8 Quantifiers

```
S LIKE_REGEX 'ab{1,}c'
```

and matches any string that contains a substring consisting of an “a”, one or more “b”s, and then a “c”.

3.9 Locating a match

LIKE_REGEX only cares whether a match exists; the other operators care about where a match is located in the string. Consider the regular expression “a+” and the string “a1aa2aaa3”. There are ten possible matches for “a+”, indicated by the underlining on the following lines:

```
'a1aa2aaa3' -- position 1, length 1
' a1aa2aaa3' -- position 3, length 1
' a1aa2aaa3' -- position 3, length 2
' a1aa2aaa3' -- position 4, length 1

' a1aa2aaa3' -- position 6, length 1
' a1aa2aaa3' -- position 6, length 2
' a1aa2aaa3' -- position 6, length 3
' a1aa2aaa3' -- position 7, length 1
' a1aa2aaa3' -- position 8, length 2
' a1aa2aaa3' -- position 9, length 1
```

Notice that some of the matches are substrings of other matches. The rules of XQuery regular expressions are designed to ignore certain matches, so that the recognized matches are mutually disjoint. Obviously there are many ways to do this, so the rules provide priorities in determining the recognized matches. There are three priorities:

- 1) The top priority is to find a match as early in the string as possible. This is commonly called the *leftmost rule*.
- 2) The second priority is to find the first alternative of an alternation, if possible. We are unaware of a common name for this rule.
- 3) The last priority is to find the longest possible match for greedy quantifiers, and the shortest match for reluctant quantifiers. In the case of greedy quantifiers, this is commonly called the *longest rule*; we are unaware of a common name for the rule regarding reluctant quantifiers.

[Historical note: POSIX only has a leftmost longest rule. There were no reluctant quantifiers, and the priority for matching alternations was the longest match rather than the first alternative.]

These rules will be illustrated by examples:

Subject string	regular expression	match(es) underlined	priority
baaaaaa	ba a*	<u>ba</u> aaaaa ba <u>aaaaa</u>	leftmost (even though baaaaaa would be longer); second match must start after the first match
ab	a ab	<u>a</u> b	first alternative (rather than matching ab)

Subject string	regular expression	match(es) underlined	priority
abcabbabc	ab*	<u>abc</u> abbabc abc <u>ab</u> abc abcabb <u>abc</u>	leftmost longest (greedy quantifier consumes two "b"s) longest
abcabbabc	ab*?	<u>abc</u> abbabc abc <u>ab</u> abc abcabb <u>abc</u>	shortest (no need to match "b") shortest shortest

3.10 Capture and back-reference

A *parenthesized sub-expression* is a portion of a regular expression that is enclosed in parentheses. Parenthesized sub-expressions are numbered in order of their left parenthesis. For example, in the regular expression

```
((a)|(b))
```

there are three sub-expressions:

- 1) ((a)|(b))
- 2) (a)
- 3) (b)

A sub-expression can be referenced later in a regular expression using a back-reference, taking the form of a backslash followed by one or more digits. Thus the three sub-expressions in the example can be referenced as “\1”, “\2”, and “\3”. For example, consider the regular expression:

```
\p{Z}(\p{L}*)\p{Z}*\1\p{Z}
```

The first and only parenthesized sub-expression (“\p{L}”*) matches any sequence of letters that is bounded by some kind of space character (“\p{Z}”) before and after the sequence of letters. The back-reference (“\1”) matches whatever sequence of letters was captured by the first sub-expression. This regular expression might be used to search for occurrences of a repeated word (perhaps caused by a cut-and-paste error). Here is an example of a subject string, with underlining to indicate the match for the entire regular expression:

```
Hello Dolly you're looking looking swell
```

When a back-reference references a parenthesized group with a quantifier, then the back-reference matches the last iteration of the quantified sub-expression. For example, consider the regular expression:

```
'(ab*)*c\1'
```

and the subject string:

```
'abbbabbabcabbbbb'
```

The matches to “(ab*)” are shown by underlining below:

```
'abbbabbabcabbbbb'
```

3.10 Capture and back-reference

```
'abbbabbabcbabbbbbb'  
'abbbabbabcbabbbbbb'
```

These three iterations of “(ab*)” are matched by “(ab*)*” and then the “c” is matched. Next, we need to match “\1”. The last match for the first parenthesized sub-expression is “ab”, so the overall match is indicated by underlining below:

```
'abbbabbabcbabbbbbb'
```

In the event that a sub-expression is unmatched, a back-reference to it matches the zero-length string. For example, consider the regular expression:

```
'((a*)|(b*))c??\3'
```

and the subject string:

```
'xyzaaccb'
```

In this example, the alternation “((a*)|(b*))” matches the “aa”, which is a match for the first alternative. Thus there is no match for the second alternative, “(b*)”. The “c?” prefers to match a zero-length string (though it could match the “c”), and the “\3” must match a zero-length string. Thus, the complete substring that is matched is underlined below:

```
'xyzaaaccb'
```

3.11 Precedence

The precedence of operators outside bracket expressions is as follows (from highest to lowest):

— Highest precedence: atoms, defined as:

- Parentheses.
- Individual characters.
- Escape sequences.
- Dot (“.”)
- Anchors (“^”, “\$”)
- Bracket expressions.

— Quantifiers.

— Concatenation.

— Alternation (lowest).

Examples:

1) Quantifiers have higher precedence than concatenation:

ab* is equivalent to a(b*)

2) Concatenation outranks alternation:

$ab|cd$ is equivalent to $(ab)|(cd)$

3.12 Modes

The preceding discussion has mentioned two of the flags, “s” to specify dot-all mode, and “m” to specify multi-line mode. There are two additional flags, “i” for case-insensitive mode, and “x” to disregard white space in regular expressions for readability. The complete set of modes is:

- "s" Specifies dot-all mode, in which a period matches any character. If “s” is not specified, then a period matches any single character except a line terminator.
- "m" Specifies multi-line mode, in which the anchors match the beginning or end of a line. If “m” is not specified, then the anchors match the beginning or end of the subject string.
- "i" Specifies case-insensitive mode.
- "x" Specifies that white space characters in a regular expression are ignored. This allows you to set off portions of a regular expression for greater readability.

(Blank page)

4 Operators using regular expressions

We propose five operators using the XQuery regular expression syntax:

- 1) LIKE_REGEX — predicate that returns *True* if a substring of a string matches a regular expression.
- 2) OCCURRENCES_REGEX — numeric function returning the number of matches for a regular expression in a string.
- 3) POSITION_REGEX — numeric function returning the position of the start of a match for a regular expression in a string, or the position of the next character after a match.
- 4) SUBSTRING_REGEX — character string function returning a substring that matches a regular expression in a string.
- 5) TRANSLATE_REGEX — character function that performs a replacement operation on one or all matches to a regular expression in a string.

4.1 LIKE_REGEX

LIKE_REGEX is a predicate that returns *True* if a substring of a string matches a regular expression.

The syntax is:

```
<regex like predicate> ::=
  <row value predicand>
  [ NOT ] LIKE_REGEX <XQuery pattern>
  [ FLAG <XQuery option flag> ]
```

where

- <row value predicand> is the subject string to be searched for matches to the <XQuery pattern>.
- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the \$flags argument of the [XQuery F&O] function `fn:match`.

The result is *Unknown* if any of the operands is the null value, *True* if there is a substring that matches the <XQuery pattern> in the <row value predicand>, and *False* if there is no match.

Note that unlike LIKE, LIKE_REGEX can return *True* without matching the entire string. The usual convention for regular expression matching is to search for a match somewhere within the searched string, without necessarily matching the entire string. The user may use anchors to require a match to the entire string.

Exceptional cases:

- If any of the parameters is the null value, the result is *Unknown*.
- If the pattern or flag is not valid, then an exception condition is raised.

4.1 LIKE_REGEX

Examples:

'abcde' LIKE_REGEX 'c' evaluates to *True*.

'abcde' LIKE_REGEX 'x' evaluates to *False*.

'abcde' LIKE_REGEX CAST (NULL AS CHAR(10)) evaluates to *Unknown*.

'abcde' LIKE_REGEX '\ ' raises an exception condition. In this example, “\” is not a well-formed regular expression.

'abcde' LIKE_REGEX 'x' FLAG '?' raises an exception condition. In this example, the flag “?” is invalid.

4.2 OCCURRENCES_REGEX

OCCURRENCES_REGEX is a numeric function returning the number of matches for a regular expression in a string. The syntax is:

```
<regex occurrences function> ::=
    OCCURRENCES_REGEX <left paren>
        <XQuery pattern> [ FLAG <XQuery option flag> ]
        IN <regex subject string>
        [ FROM <start position> ]
        [ USING <char length units> ] <right paren>
```

where:

- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the \$flags argument of the [XQuery F&O] function fn:match.
- <regex subject string> is the character string to be searched for matches to the <XQuery pattern>.
- <start position> is an optional exact numeric value with scale 0 (zero) specifying the position at which to start the search (the default is position 1 (one)).
- <char length units> is CHARACTERS or OCTETS, indicating the unit in which <start position> is measured (the default is to measure in CHARACTERS).

The <regex subject string> is searched for matches to the <XQuery pattern>, starting from position <start position>, which is measured in the units specified by <char length units>, either CHARACTERS or OCTETS. The result is the number of matches.

Exceptional cases:

- If any of the parameters is the null value, then the result is the null value.
- If the pattern or flag is not valid, then an exception condition is raised.
- If a starting position is given in octets, but it is not the first octet of a character, then the result is implementation-dependent. The use of OCTETS is discussed under POSITION_REGEX.
- If any of the numeric parameters is too large or too small, then the result is -1. This includes the following cases:

- The starting position is less 1 (one).
- The starting position is greater than the length of the string (measured in CHARACTERS or OCTETS as specified by <char length units>).

Examples:

OCCURRENCES_REGEX ('a' IN 'what is that?') evaluates to 2.

OCCURRENCES_REGEX ('a' IN 'what is that?' FROM 5) evaluates to 1 (one).

OCCURRENCES_REGEX ('A' FLAG 'i' IN 'what is that') evaluates to 2.

OCCURRENCES_REGEX ('A' IN 'what is that') evaluates to 0 (zero).

4.3 POSITION_REGEX

POSITION_REGEX is a numeric function returning the position of the start of a match, or one plus the end of a match, for a regular expression in a string. The syntax is:

```
<regex position expression> ::=
    POSITION_REGEX <left paren> [ <regex position start or after> ]
    <XQuery pattern> [ FLAG <XQuery option flag> ]
    IN <regex subject string>
    [ FROM <start position> ]
    [ USING <char length units> ]
    [ OCCURRENCE <regex occurrence> ]
    [ GROUP <regex capture group> ] <right paren>

<regex position start or after> ::=
    START
    | AFTER
```

where:

- START indicates that the starting position of the match to the regular expression is desired; AFTER indicates that the character position immediately following the match is desired (START is the default). If the match consumes the last character of the subject string, then AFTER returns the length of the string plus 1 (one).
- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the \$flags argument of the [XQuery F&O] function fn:match.
- <regex subject string> is the character string to be searched for matches to the <XQuery pattern>.
- <start position> is an optional exact numeric value with scale 0 (zero), identifying the character position at which to start the search (the default is 1 (one)).
- <char length units> is CHARACTERS or OCTETS, indicating the unit in which <start position> is measured, and the unit in which the returned position is measured (the default is to measure in CHARACTERS).
- <regex occurrence> is an optional exact numeric value with scale 0 (zero) indicating which occurrence of a match is desired (the default is 1 (one)).

PDTR 19075-1:2011(E)
4.3 POSITION_REGEX

— <regex capture group> is an optional exact numeric value with scale 0 (zero) indicating which capture group of a match is desired (the default is 0 (zero), indicating the entire occurrence).

The <regex subject string> is searched for matches to the <XQuery pattern>. If there are at least *RO* matches, where *RO* is the value of <regex occurrence>, then either the starting position of the *RO*-th match, or the position immediately following the *RO*-th match, is returned (for the *START* or *AFTER* options, respectively). Positions are measured in the units specified by <char length units>, either *CHARACTERS* or *OCTETS*. If a <regex capture group> *CAP* is specified, then the position at the start or immediately following the substring that matches the *CAP*-th parenthesized subexpression is used.

With *AFTER*, note that the position returned is the one after the match. If the match consumes the last character of the string, then the position returned is actually one plus the length of the string (in characters or octets, as requested by <char length units>). The rationale for providing the position that is 1 (one) after the end of the match is that this is the correct place to begin a search for the next match. If the user wishes to process the subject string in a loop, the loop can continue until the *AFTER* position is greater than the length of the subject string. However, when doing this, the user must beware of a pitfall: if the regular expression matches a zero-length string, then the *AFTER* position and the *START* position are the same, and resuming the search at the *AFTER* position will simply find the same zero-length match again.

OCTETS is provided for efficient processing for those UCS encodings that do not have a fixed character width. It is expected that the user will use the output of *POSITION_REGEX* (... *USING OCTETS* ...) to learn the position of some occurrence within a string, measured in octets. That value is then known to be the first octet of a character, and may be used as a starting position in other function invocations. If the user picks an arbitrary octet number, it may be other than the first octet of a character. Naturally, beginning a regular expression match at such an octet can produce unpredictable results. Therefore we say that the result is implementation-dependent if a starting octet is not the first octet of a character.

Exceptional cases:

- If any of the parameters is the null value, the result is the null value.
- If the pattern or flag is not valid, then an exception condition is raised.
- If a starting position is given in octets, but it is not the first octet of a character, then the result is implementation-dependent.
- If any of the numeric parameters is too large or too small, then the result is 0 (zero). This includes the following cases:
 - The starting position is less 1 (one).
 - The starting position is greater than the length of the string (measured in *CHARACTERS* or *OCTETS* as specified by <char length units>).
 - There are not at least *RO* matches.
 - There are not *CAP* parenthesized subexpressions.

Examples:

```
POSITION_REGEX ( 'a' IN 'what is that?' ) evaluates to 3.
```

```
POSITION_REGEX ( START 'a' IN 'what is that?' ) evaluates to 3.
```

```
POSITION_REGEX ( AFTER 'a' IN 'what is that?' ) evaluates to 4.
```

```
POSITION_REGEX ( AFTER 'a' IN 'a' ) evaluates to 2.
```

POSITION_REGEX ('a' IN 'what is that?' FROM 5) evaluates to 11.
 POSITION_REGEX ('a' IN 'what is that?' OCCURRENCE 2) evaluates to 11.
 POSITION_REGEX ('(a)(t)' IN 'what is that?' GROUP 2) evaluates to 4.
 POSITION_REGEX ('A' FLAG 'i' IN 'what is that') evaluates to 3.
 POSITION_REGEX ('A' IN 'what is that') evaluates to 0.

4.4 SUBSTRING_REGEX

SUBSTRING_REGEX is a character string function returning a substring that matches a regular expression in a string. The syntax is:

```
<regex substring function> ::=
  SUBSTRING_REGEX <left paren>
    <XQuery pattern> [ FLAG <XQuery option flag> ]
    IN <regex subject string>
    [ FROM <start position> ]
    [ USING <char length units> ]
    [ OCCURRENCE <regex occurrence> ]
    [ GROUP <regex capture group> ] <right paren>
```

where:

- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the \$flags argument of the [XQuery F&O] function fn:match.
- <regex subject string> is the character string to be searched for matches to the <XQuery pattern>.
- <start position> is an optional exact numeric value with scale 0 (zero), indicating the character position at which to start the search (the default is position 1 (one)).
- <char length units> is CHARACTERS or OCTETS, indicating the unit in which <start position> is measured (the default is to measure in CHARACTERS).
- <regex occurrence> is an optional exact numeric value with scale 0 (zero) indicating which occurrence of a match is desired (the default is 1 (one)).
- <regex capture group> is an optional exact numeric value with scale 0 (zero) indicating which capture group of a match is desired (the default is 0 (zero), indicating the entire occurrence).

The <regex subject string> is searched for matches to the <XQuery pattern>. If there are at least *RO* matches, where *RO* is the value of <regex occurrence>, then the result is the substring that is the *RO*-th match. If <regex capture group> *CAP* is specified, then the result is the substring that matches the *CAP*-th parenthesized substring within the substring that is the *RO*-th match. If there are not at least *RO* matches, or at least *CAP* parenthesized subexpressions, the result is the null value.

The exceptional cases are:

- If any of the parameters is the null value, the result is the null value.

PDTR 19075-1:2011(E)
4.4 SUBSTRING_REGEX

- If the pattern or flag is not valid, then an exception condition is raised.
- If a starting position is given in octets, but it is not the first octet of a character, then the result is implementation-dependent.
- If any of the numeric parameters is too large or too small, then the result is the null value. This includes the following cases:
 - The starting position is less than 1 (one).
 - The starting position is greater than the length of the string (measured in CHARACTERS or OCTETS as specified by <char length units>).
 - There are not at least *RO* matches.
 - There are not *CAP* parenthesized subexpressions.

Examples:

SUBSTRING_REGEX ('\p{L}*' IN 'what is that?') evaluates to the string “what”.

SUBSTRING_REGEX ('\p{L}*' IN 'what is that?' FROM 2) evaluates to the string “hat”.

SUBSTRING_REGEX ('\p{L}*' IN 'what is that?' OCCURRENCE 2) evaluates to the string “is”.

SUBSTRING_REGEX ('(is) (\p{L}*)' IN 'what is that?' GROUP 2) evaluates to the string “that”.

4.5 TRANSLATE_REGEX

TRANSLATE_REGEX is a character string function that performs a replacement operation on one or all matches to a regular expression in a string. The syntax is:

```
<regex transliteration> ::=
  TRANSLATE_REGEX <left paren>
    <XQuery pattern> [ FLAG <XQuery option flag> ]
    IN <regex subject string>
    [ WITH <regex replacement string> ]
    [ FROM <start position> ]
    [ USING <char length units> ]
    [ OCCURRENCE <regex transliteration occurrence> ] <right paren>
```

```
<regex transliteration occurrence> ::=
  <regex occurrence>
  | ALL
```

where:

- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the `$flags` argument of the [XQuery F&O] function `fn:match`.
- <regex subject string> is the character string to be searched for matches to the <XQuery pattern>.

- <regex replacement string> is a character string whose value is suitable for use as the \$replacement argument of the [XQuery F&O] function `fn:replace`. The special syntax for replacement strings is discussed below. The default is the zero-length string.
- <start position> is an optional exact numeric value with scale 0 (zero), indicating the character position at which to start the search (the default position is 1 (one)).
- <char length units> is CHARACTERS or OCTETS, indicating the unit in which <start position> is measured (the default is to measure in CHARACTERS).
- <regex transliteration occurrence> is either the keyword ALL, or an exact numeric value with scale 0 (zero), indicating which occurrence of a match is desired (the default is ALL).

The <regex subject string> is searched for matches to the <XQuery pattern>. If ALL is specified or implied, then every match is replaced by the value of <regex replacement string>. If a numeric <regex transliteration occurrence> is specified, then only that match is replaced.

Exceptional cases:

- If any of the parameters is the null value, then the result is the null value.
- If the pattern, flag or replacement string is not valid, then an exception condition is raised.
- If the pattern matches a zero-length string, then an exception condition is raised (this is the behavior of XQuery's `fn:replace` in this case)
- If a starting position is given in octets, but it is not the first octet of a character, then the result is implementation-dependent.
- If any of the numeric parameters is too large or too small, then the result is the null value. This includes the following cases:
 - The starting position is less than 1 (one).
 - The starting position is greater than the length of the string (measured in CHARACTERS or OCTETS as specified by <char length units>).
 - A numeric <regex transliteration occurrence> is specified, and there are not at least that many matches.

First, here are some examples with no replacement string. In these examples, any matched substring is replaced by a zero-length string, effectively removing the matched substring.

`TRANSLATE_REGEX ('a' IN 'what was that?')` evaluates to the string “wht ws tht?”

`TRANSLATE_REGEX ('a' IN 'what was that?' OCCURRENCE ALL)` evaluates to the string “wht ws tht?”

`TRANSLATE_REGEX ('a' IN 'what was that?' FROM 5)` evaluates to the string “what ws tht?”

`TRANSLATE_REGEX ('a' IN 'what was that?' OCCURRENCE 2)` evaluates to the string “what ws that?”

`TRANSLATE_REGEX ('A' IN 'what was that?')` evaluates to the string “what was that?”

`TRANSLATE_REGEX ('A' FLAG 'i' IN 'what was that?')` evaluates to the string “wht ws tht?”

Next, here are some examples in which the matched substrings are replaced with constant text:

4.5 TRANSLATE_REGEX

`TRANSLATE_REGEX ('a' IN 'what was that?' WITH 'U')` evaluates to the string “whUt wUs thUt?”

`TRANSLATE_REGEX ('a' IN 'what was that?' WITH 'U' OCCURRENCE ALL)` evaluates to the string “whUt wUs thUt?”

`TRANSLATE_REGEX ('a' IN 'what was that?' WITH 'U' OCCURRENCE 2)` evaluates to the string “what wUs that?”

`TRANSLATE_REGEX ('a' IN 'what was that?' WITH 'U' FROM 5)` evaluates to the string “what wUs thUt?”

`TRANSLATE_REGEX ('A' FLAG 'i' IN 'what was that?' WITH 'U')` evaluates to the string “whUt wUs thUt?”

Finally, this is an example of a special capability of the replacement string, to replace a matched substring with the value of a captured sub-expression. `$0` in a replacement string is used to represent the entire matched substring, as in this example:

`TRANSLATE_REGEX ('\p{L}*' IN 'what was that?' WITH '<$0>')` evaluates to the string “<what> <was> <that>?”

The preceding example uses “`\p{L}`” to find maximal substrings of letters. Each such substring is then enclosed in angle brackets using the replacement string. The following example makes this replacement only on the second occurrence of a match:

`TRANSLATE_REGEX ('\p{L}*' IN 'what was that?' WITH '<$0>' OCCURRENCE 2)` evaluates to the string “what <was> that?”

`$n` in a replacement string represents the value of the n -th captured sub-expression. For example:

`TRANSLATE_REGEX ('([\p{L}-[aeiou]]*)([aeiou]*)([\p{L}-[aeiou]])' IN 'what was that?' WITH '$3-$2-$1')` evaluates to the string “t-a-wh s-a-w t-a-th?”

In this example, “[aeiou]” is the set of English vowels and “[\p{L}-[aeiou]]” is the set of other letters (consonants for the purpose of this example). The regular expression looks for a string of consonants followed by a string of vowels followed by a string of consonants. These three strings constitute the three parenthesized subexpressions. The replacement string rearranges the three strings and places minus signs between them.

To represent a literal dollar sign (“\$”) in a replacement string, you must write “`\$`”. To represent a literal backslash (“\”) in a replacement string, you must write “`\\`”. There are no other escape sequences in a replacement string. For example:

`TRANSLATE_REGEX ('\p{L}' IN 'what was that?' WITH '\p{L}')`

is invalid because “`\p`” is not a recognized escape sequence. This example will raise an exception condition.

Index

Index entries appearing in **boldface** indicate the page where the word, phrase, or BNF nonterminal was defined; index entries appearing in *italics* indicate a page where the BNF nonterminal was used in a Format; and index entries appearing in roman type indicate a page where the word, phrase, or BNF nonterminal was used in a heading, Function, Syntax Rule, Access Rule, General Rule, Leveling Rule, Table, or other descriptive text.

— A —

AFTER • 6, 19, 20
 ALL • 22, 23
 anchors • 8
 AS • 18

— B —

block escape • 7
 bracket expression • 9

— C —

category escape • 6
 CHAR • 18
 character class expressions • 9
 character class subtraction • 10
 CHARACTERS • 18, 19, 20, 21, 22, 23
 complementary block escape • 7
 complementary category escape • 6

— D —

dot-all mode • 7

— E —

escape sequence • 6

— F —

FLAG • 8, 17, 18, 19, 21, 22, 23, 24
 FROM • 18, 19, 21, 22, 23, 24

— G —

greedy quantifier • 11
 GROUP • 19, 21, 22

— I —

IN • 5, 6, 18, 19, 20, 21, 22, 23, 24

— L —

leftmost rule • 12
 LIKE • 5, 7, 17
 LIKE_REGEX • 1, 5, 6, 7, 8, 9, 10, 11, 12, 17, 18
 longest rule • 12

— M —

metacharacters • 6
 multi-character escape sequence • 7
 multi-line mode • 8

— N —

NOT • 17
 NULL • 18

— O —

OCCURRENCE • 5, 6, 19, 21, 22, 23, 24
 OCCURRENCES_REGEX • 1, 5, 17, 18, 19
 OCTETS • 18, 19, 20, 21, 22, 23

— P —

parenthesized sub-expression • 13
 POSITION_REGEX • 1, 5, 6, 9, 17, 18, 19, 20, 21

— R —

<regex like predicate> • 17
 <regex occurrences function> • 18
 <regex position expression> • 19
 <regex position start or after> • 19
 <regex substring function> • 21
 <regex transliteration> • 22
 <regex transliteration occurrence> • 22, 23
 reluctant quantifier • 11

— S —

single character escape sequence • 6

special characters • 9

START • 19, 20

SUBSTRING_REGEX • 1, 17, 21, 22

— T —

TRANSLATE_REGEX • 1, 17, 22, 23, 24

— U —

USING • 18, 19, 21, 22

— W —

WITH • 22, 24

— X —

Feature X211, “XML 1.1 support” • 8