

***This is a OGF DFDL Workgroup draft document in preparation.***

***Use with caution. The material herein is subject to change, and URL links to this version will be broken when newer revisions and the final version are released.***

***This draft document is an updated version of GFD-P-R.207 with the errata and clarifications applied to it that have been gathered since its first publication.***

***The updates are indicated with change bars, text coloring, and text strikethrough in the style of this sentence.***

## Data Format Description Language (DFDL) v1.0 Specification

### Status of This Document

Grid Final Draft (GFD)

### Obsoletes

This document obsoletes GFD-P-R.174 dated January 2011 [[OBSOLETE DFDL](#)].

### Copyright Notice

Copyright © Global Grid Forum (2004-2006). Some Rights Reserved. Distribution is unlimited.

Copyright © Open Grid Forum (2006-~~2019~~). Some Rights Reserved. Distribution is unlimited

### Abstract

This document provides a definition of a standard Data Format Description Language (DFDL). This language allows description of text, dense binary, and legacy data formats in a vendor-neutral declarative manner. DFDL is an extension to the XML Schema Description Language (XSDL).

Document Change History

Revision ID	Date	Description
R11	2019-12-09	Edited in minor/typo fixes.
R10	2019-12-01	Erratum 5.38 Updated UML Diagram in Section 4.1 to have missing members on element class. Marked some comments as resolved. Date changed to December 2019
R09	2019-11-25	Errata Merged 5.39 Section 15.1.1 Clarification for choice branches that are zero-occurrences arrays. 5.40 Section 14.1 Empty Sequence and Hidden Group References 5.41 Section 2.6 Dynamic Type error <i>must</i> cause SDE. 5.42 Section 21 utf16Width 'variable' property value is optional. 5.43 Sections 6.3.1, 21, 4.1.2, 6.3.1.1, 13.3, 22.1.1, 2.2.1 Bidirectional text properties other than textBidi="no" removed and reserved for future use. 5.44 Sections, 9.2, 12.3.3, 23.5.3 Data syntax grammar corrections related to ComplexContent. 5.45 Sections 9.2, 13.2.1, 23.5.3. Data syntax grammar corrections related to SimpleContent and SimpleValue. 5.46 New Section 23.6 on Circular Deadlocks. 5.47 Section 11.2.1.1 Detection of Encoding/Decoding Errors. 5.48 Section 23.5.5 fn:error() function. 5.49 Section 5.3 Allowing import of other annotation language schemas. 5.50 Sections 4.1.1, 9.2, 11, 11.1 Remove Byte-order-mark (BOM) processing. 5.51 Section 12.2 Entity %ES; allowed alone in terminator when not delimited. 5.52 Section 15 Clarify choiceDispatchKey evaluating to empty string is SDE. 5.53 Section 13.16 Clarifications on nilValueDelimiterPolicy. 5.54 Section 12.2 Clarifications on emptyValueDelimiterPolicy. 5.55 Section 9.3.2 Clarifications on Establishing Representation. 5.56 Section 9.5, 9.4.2, 9.4.2.2, 9.4.2.3 Clarifications on Element Defaults. 5.57 Section 9.2.2 Clarifications on Empty Representation. 5.58 Section 14.2 Clarifications on Separators and Suppression. 5.59 Section 14.2.1, 14.2.2, 14.2.3 Clarifications on Potentially Trailing Groups and Elements.
R07	2019-10-11	Fixed typos with errata-related changes per SMH email of 2019-10-04.
R06	2019-09-26	Errata Merged

		<p>5.19. Section 16.1.1, 16.1.2 Correct occursCountKind 'fixed' and 'implicit' to state that the unparser stops looking for occurrences in the Infoset once maxOccurs has been reached.</p> <p>5.20. Section 23.1, 13.11, 15, 7.9 outputValueCalc maximum forward reference is implementation-defined. No forward reference from calendarLanguage, choiceDispatchKey, setVariable.</p> <p>5.21. Section 15. choiceBranchKey is case sensitive.</p> <p>5.22. Section 13.2.1. Parsing block escape scheme processing error if the end of data for the element and escapeBlockEnd is not found.</p> <p>5.23 Section 13.16 nilKind, nilValue, and nilValueDelimiterPolicy properties apply to nillable complex type elements as well as simple types.</p> <p>5.24 Sections 13.11.1, 12.3.7.2.6. Fractional-seconds "S" corrections. Packed decimal calendar max digits is implementation defined.</p> <p>5.25 Section 17. Allow inputValueCalc on root of a choice branch.</p> <p>5.26 Section 2. RFC2119 terminology, except for required and optional.</p> <p>5.27 Section 11. Bit order change requires byte alignment both parsing and unparsing.</p> <p>5.28 Section 13.6. textStandardNaNRep and textStandardInfinityRep: Schema Definition Error if empty string.</p> <p>5.29 Sections 23.4, 23.5.3. Add XPath 2.0 'intersect' and "except" operators. Add dfdl:checkRangeInclusivem dfdl:checkRangeExclusive functions.</p> <p>5.30 Section 12.3.2.2. Clarify scanning text can require byte-level scans, character level is insufficient.</p> <p>5.31 Section 9.5. Clarify Evaluation Order for Statement Annotations section on property expression evaluation times, and grammar regions.</p> <p>5.32 Section 9.3.2.1, 9.2.5 ES and WSP* corrections for nil rep</p> <p>5.33 Section 13.9 textBooleanTrueRep cannot be empty string.</p> <p>5.34 Section 13.2 Fix textPadKind omission of dfdl:textBooleanJustification case.</p> <p>5.35 Section 5.2 Remove redundant description from pattern facet in list.</p>
--	--	---

		<p>5.36 Section 15, 9.3.1.1, 9.3.1.3, 9.3.3, 15.1.2 choiceBranchKey is List of DFDL String Literals.</p> <p>5.37 Section 15 Specify SDE for choiceDispatchKey forward reference.</p> <p>5.38 Section 4.1.2 Add element information item [array].</p> <p><a href="#">That is all the errata up to 2019-09-25.</a></p>
R05	2019-09-17	<p>Errata Merged</p> <p>5.4 Section 23.5.2.5 fn:exactly-one definition should match XPath</p> <p>5.5 Section 23 binding fn: prefix to XPath function namespace.</p> <p>5.6 Section 5.1 elementFormDefault and form attributes</p> <p>5.7 Section 13.6 textNumberPadCharacter clarifications</p> <p>5.8 Sections 14.2, 14.2.2, 14.2.3, 16 Clarify the behaviour of separator suppression.</p> <p>5.9 Section 16.1 Clarify the behaviour when maxOccurs or occursCount is zero.</p> <p>5.10 Sections 3, 9.2.2, 9.4.2.2, 9.4.2.3, 12.2, 13.16 Definition for 'fixed-length element' and clarifications related to dfdl:emptyValueDelimiterPolicy (that it is not applicable).</p> <p>5.11 Section 23.4 Clarify unqualified path steps.</p> <p>5.12 Section 12.3.7.2.1 Clarify decimal minimum bits.</p> <p>5.13 Section 13.7.1.1 Clarifications on converting base-2 binary numbers.</p> <p>5.14 Section 13.7 Clarifications for dfdl:binaryDecimalVirtualPoint</p> <p>5.15 Section 12.3 Correction to dfdl:lengthUnits 'bits' allowed for packed calendars.</p> <p>5.16. Section 12.3.3, 12.3.4 Correction that binarySeconds and binaryMilliseconds are not allowable representations for types xs:date and xs:time.</p> <p>5.17. Section 13.11 Correct regular expression for dfdl:calendarTimeZone.</p> <p>5.18. Section 3.11 Correction for dfdl:lengthKind 'expression' – when unparsing the expression is evaluated to obtain the length.</p>
R04	2019-09-16	<p>Upgraded from MS-Word "Compatibility Mode" to "normal" so tools we use don't have to deal with compatibility mode.</p>
R03	2019-09-13	<p>Errata Merged</p> <p>5.1 merged. Section 13.2.1.1 Escape Scheme Examples</p> <p>5.2 merged. Section 23. Clarification of expression evaluation order.</p> <p>5.3 merged. Section 13.2.1. New property escapeCharacterPolicy added to dfdl:escapeScheme. Dfdl:escapeKind language adjusted to take new property into consideration.</p>



R02	2019-09-13	Changes for Minor typographical fixes accepted. Two Comment bubbles remain where formatting changes could not be applied.
R01	2019-09-09	Minor typographical fixes. (Remove this entry once changes accepted.) <a href="https://redmine.ogf.org/issues/233">https://redmine.ogf.org/issues/233</a> through update #28.

Contents

Data Format Description Language (DFDL) v1.0 Specification .....	1
1. Introduction .....	12
1.1 Why is DFDL Needed? .....	13
1.2 What is DFDL? .....	13
1.2.1 Simple Example .....	13
1.3 What DFDL is not .....	16
1.4 Scope of version 1.0 .....	16
1.5 Related standards .....	17
2. Notational and Definitional Conventions .....	18
2.1 Failure Types .....	18
2.2 Schema Definition Error (SDE) .....	18
2.3 Processing Errors .....	18
2.3.1 Ambiguity of Data Formats .....	19
2.4 Validation Errors .....	19
2.5 Recoverable Error .....	20
2.6 Specific Errors Classified .....	20
2.7 Optional Checks and Warnings .....	22
3. Glossary .....	24
4. The DFDL Information Set (Infoset) .....	31
4.1 Information Items .....	31
4.1.1 Document Information Item .....	31
4.1.2 Element Information Items .....	31
4.2 "No Value" .....	33
4.3 DFDL Information Item Order .....	33
4.4 DFDL Infoset Object model .....	33
4.5 DFDL Augmented Infoset .....	34
5. DFDL Schema Component Model .....	36
5.1 DFDL Subset of XML Schema .....	37
5.2 XSD Facets, min/maxOccurs, default, and fixed .....	39
5.2.1 MinOccurs, MaxOccurs .....	39
5.2.2 MinLength, MaxLength .....	40
5.2.3 MaxInclusive, MaxExclusive, MinExclusive, MinInclusive, TotalDigits, FractionDigits 40	
5.2.4 Pattern .....	40
5.2.5 Enumeration .....	40
5.2.6 Default .....	40
5.2.7 Fixed .....	41
5.3 Compatibility with Other Annotation Language Schemas .....	41
6. DFDL Syntax Basics .....	42
6.1 Namespaces .....	42
6.2 The DFDL Annotation Elements .....	42
6.3 DFDL Properties .....	43

6.3.1	DFDL String Literals.....	44
6.3.2	DFDL Expressions.....	48
6.3.3	DFDL Regular Expressions .....	49
6.3.4	Enumerations in DFDL.....	49
7.	Syntax of DFDL Annotation Elements.....	50
7.1	Component Format Annotations .....	50
7.1.1	The dfdl:ref Property .....	50
7.1.2	Property Binding Syntax .....	51
7.1.3	Empty String as a Representation Property Value .....	52
7.2	dfdl:defineFormat - Reusable Data Format Definitions.....	53
7.2.1	Inheritance for dfdl:defineFormat.....	53
7.2.2	Using/Referencing a Named Format Definition .....	53
7.3	The dfdl:assert Statement Annotation Element .....	53
7.3.1	Properties for dfdl:assert.....	54
7.3.2	Controlling the Timing of Statement Evaluation .....	56
7.4	The dfdl:discriminator Statement Annotation Element.....	56
7.4.1	Properties for dfdl:discriminator .....	57
7.5	The dfdl:defineEscapeScheme Defining Annotation Element .....	59
7.5.1	Using/Referencing a Named escapeScheme Definition.....	60
7.6	The dfdl:escapeScheme Annotation Element.....	60
7.7	The dfdl:defineVariable Annotation Element.....	60
7.7.1	Examples .....	61
7.7.2	Predefined Variables.....	61
7.8	The dfdl:newVariableInstance Statement Annotation Element.....	62
7.8.1	Examples .....	62
7.9	The dfdl:setVariable Statement Annotation Element .....	63
7.9.1	Examples .....	63
8.	Property Scoping Rules.....	64
8.1	Providing Defaults for DFDL properties .....	64
8.2	Combining DFDL Representation Properties from a dfdl:defineFormat .....	65
8.3	Combining DFDL Properties from References .....	65
9.	DFDL Processing Introduction .....	69
9.1	Parser Overview.....	69
9.2	DFDL Data Syntax Grammar.....	70
9.2.1	Nil Representation .....	72
9.2.2	Empty Representation .....	72
9.2.3	Normal Representation .....	72
9.2.4	Absent Representation .....	72
9.2.5	Zero-length Representation .....	72
9.2.6	Missing .....	73
9.2.7	Examples of Missing and Empty Representation .....	73
9.2.8	Round Trip Ambiguities.....	74
9.3	Parsing Algorithm .....	74
9.3.1	Known-to-exist and Known-not-to-exist.....	74

9.3.2	Establishing Representation .....	76
9.3.3	Points of Uncertainty .....	77
9.4	Element Defaults .....	78
9.4.1	Definition 'default value' .....	78
9.4.2	Element Defaults When Parsing .....	78
9.4.3	Element Defaults When Unparsing .....	79
9.5	Evaluation Order for Statement Annotations .....	80
9.5.1	Asserts and Discriminators with testKind 'expression' .....	81
9.5.2	Discriminators with testKind 'expression' .....	81
9.5.3	Elements and setVariable .....	81
10.	Core Representation Properties and their Format Semantics .....	82
11.	Properties Common to both Content and Framing .....	83
11.1	Unicode Byte Order Mark (BOM) .....	86
11.2	Character Encoding and Decoding Errors .....	86
11.2.1	Property dfdl:encodingErrorPolicy .....	86
11.2.2	Unicode UTF-16 Decoding/Encoding Non-Errors .....	88
11.2.3	Preserving Data Containing Decoding Errors .....	88
11.3	Byte Order and Bit Order .....	88
11.4	dfdl:bitOrder Example .....	88
11.4.1	Example Using Right-to-Left Display for 'leastSignificantBitFirst' .....	89
11.4.2	dfdl:bitOrder and Grammar Regions .....	89
12.	Framing .....	90
12.1	Aligned Data .....	90
12.1.1	Implicit Alignment .....	91
12.1.2	Mandatory Alignment for Textual Data .....	92
12.1.3	Mandatory Alignment for Packed Decimal Data .....	93
12.1.4	Example: AlignmentFill .....	93
12.2	Properties for Specifying Delimiters .....	93
12.3	Properties for Specifying Lengths .....	98
12.3.1	dfdl:lengthKind 'explicit' .....	99
12.3.2	dfdl:lengthKind 'delimited' .....	100
12.3.3	dfdl:lengthKind 'implicit' .....	101
12.3.4	dfdl:lengthKind 'prefixed' .....	102
12.3.5	dfdl:lengthKind 'pattern' .....	105
12.3.6	dfdl:lengthKind 'endOfParent' .....	106
12.3.7	Elements of Specified Length .....	107
13.	Simple Types .....	112
13.1	Properties Common to All Simple Types .....	112
13.2	Properties Common to All Simple Types with Text representation .....	113
13.2.1	The dfdl:escapeScheme Properties .....	114
13.3	Properties for Bidirectional support for All Simple Types with Text representation ..	119
13.4	Properties Specific to String .....	119
13.5	Properties Specific to Number with Text or Binary Representation .....	121
13.6	Properties Specific to Number with Text Representation .....	121

13.6.1	The dfdl:textNumberPattern Property .....	128
13.6.2	Converting logical numbers to/from text representation .....	134
13.7	Properties Specific to Number with Binary Representation .....	135
13.7.1	Converting Logical Numbers to/from Binary Representation .....	136
13.8	Properties Specific to Float/Double with Binary Representation .....	141
13.9	Properties Specific to Boolean with Text Representation .....	141
13.10	Properties Specific to Boolean with Binary Representation .....	143
13.11	Properties specific to Calendar with Text or Binary Representation .....	143
13.11.1	The dfdl:calendarPattern property .....	145
13.11.2	The dfdl:calendarCheckPolicy Property .....	148
13.12	Properties Specific to Calendar with Text Representation .....	149
13.13	Properties Specific to Calendar with Binary Representation .....	150
13.14	Properties Specific to Opaque Types (xs:hexBinary) .....	151
13.15	Nil Value Processing .....	151
13.16	Properties for Nillable Elements .....	152
14.	Sequence Groups .....	156
14.1	Empty Sequences .....	156
14.2	Sequence Groups with Separators .....	157
14.2.1	Separators and Suppression .....	159
14.2.2	Parsing Sequence Groups with Separators .....	160
14.2.3	Unparsing Sequence Groups with Separators .....	163
14.3	Unordered Sequence Groups .....	165
14.3.1	Restrictions for Unordered Sequences .....	165
14.3.2	Parsing an Unordered Sequence .....	165
14.3.3	Unparsing an Unordered Sequence .....	166
14.4	Floating Elements .....	167
14.5	Hidden Groups .....	168
15.	Choice Groups .....	170
15.1	Resolving Choices .....	171
15.1.1	Resolving Choices via Speculation .....	171
15.1.2	Resolving Choices via Direct Dispatch .....	172
15.1.3	Unparsing Choices .....	172
16.	Properties for Array Elements and Optional Elements .....	173
16.1	The dfdl:occursCountKind property .....	174
16.1.1	dfdl:occursCountKind 'fixed' .....	174
16.1.2	dfdl:occursCountKind 'implicit' .....	174
16.1.3	dfdl:occursCountKind 'parsed' .....	174
16.1.4	dfdl:occursCountKind 'expression' .....	174
16.1.5	dfdl:occursCountKind 'stopValue' .....	175
16.2	Default Values for Arrays .....	175
16.3	Arrays with DFDL Expressions .....	175
16.4	Points of Uncertainty .....	175
16.5	Arrays and Sequences .....	175
16.6	Forward Progress Requirement .....	175

16.7	Parsing Occurrences with Non-Normal Representation .....	176
16.8	Sparse Arrays .....	176
17.	Calculated Value Properties .....	177
17.1	Example: 2d Nested Array .....	178
17.2	Example: Three-Byte Date .....	179
18.	External Control of the DFDL Processor .....	182
19.	Built-in Specifications .....	183
20.	Conformance .....	184
21.	Optional DFDL Features .....	185
22.	Property Precedence .....	187
22.1	Parsing .....	187
22.1.1	dfdl:element (simple) and dfdl:simpleType .....	187
22.1.2	dfdl:element (complex) .....	191
22.1.3	dfdl:sequence and dfdl:group (when reference is to a sequence) .....	192
22.1.4	dfdl:choice and dfdl:group (when reference is to a choice) .....	193
22.2	Unparsing .....	193
22.2.1	dfdl:element (simple) and dfdl:simpleType .....	193
22.2.2	dfdl:element (complex) .....	198
22.2.3	dfdl:sequence and dfdl:group (when reference is a sequence) .....	199
22.2.4	dfdl:choice and dfdl:group (when reference is a choice) .....	200
23.	Expression language .....	201
23.1	Expression Language Data Model .....	202
23.2	Variables .....	202
23.2.1	Rewinding of Variable Memory State .....	203
23.2.2	Variable Memory State Transitions .....	203
23.3	General Syntax .....	204
23.4	DFDL Expression Syntax .....	204
23.5	Constructors, Functions and Operators .....	206
23.5.1	Constructor Functions for XML Schema Built-in Types .....	206
23.5.2	Standard XPath Functions .....	207
23.5.3	DFDL Functions .....	211
23.5.4	DFDL Constructor Functions .....	213
23.5.5	Miscellaneous Functions .....	214
23.6	Unparsing and Circular Expression Deadlock Errors .....	215
24.	DFDL Regular Expressions .....	216
25.	Security Considerations .....	217
26.	Authors and Contributors .....	218
27.	Intellectual Property Statement .....	219
28.	Disclaimer .....	220
29.	Full Copyright Notice .....	221
30.	References .....	222
31.	Appendix A: Escape Scheme Use Cases .....	225
31.1	Escape Character Same as dfdl:escapeEscapeCharacter .....	225
31.2	Escape Character Different from dfdl:escapeEscapeCharacter .....	225

31.3	Escape Block with Different Start and End Characters .....	226
31.4	Escape Block with Same Start and End Characters .....	227
32.	Appendix B: Rationale for Single-Assignment Variables .....	229
33.	Appendix C: Processing of DFDL String literals .....	230
33.1	Interpreting a DFDL String Literal .....	230
33.2	Recognizing a DFDL String Literal .....	230
33.3	Recognizing DFDL String Literal Part .....	230
34.	Appendix D: DFDL Standard Encodings .....	232
34.1	Purpose .....	232
34.2	Conventions .....	232
34.3	Specification Template .....	232
34.4	Encoding X-DFDL-US-ASCII-7-BIT-PACKED .....	232
34.4.1	Name .....	232
34.4.2	Translation table .....	232
34.4.3	Width .....	233
34.4.4	Alignment .....	233
34.4.5	Byte Order .....	233
34.4.6	Example 1 .....	233
34.4.7	Example 2 .....	234
34.5	Encoding X-DFDL-US-ASCII-6-BIT-PACKED .....	235
34.5.1	Name .....	235
34.5.2	Translation Table .....	235
34.5.3	Width .....	236
34.5.4	Alignment .....	236
34.5.5	ByteOrder .....	236
34.5.6	Example 1 .....	236
34.6	References for Appendix D .....	237

## 1. Introduction

Data interchange is critically important for most computing. Grid computing, Cloud computing, and all forms of distributed computing require distributed software and hardware resources to work together. Inevitably, these resources read and write data in a variety of formats. General tools for data interchange are essential to solving such problems. Scalable and High Performance Computing (HPC) applications require high-performance data handling, so data interchange standards must enable efficient representation of data. Data Format Description Language (DFDL) enables powerful data interchange and very high-performance data handling.

We envisage three dominant kinds of data in the future, as follows:

1. Textual data defined by a format specific schema such as XML[XML] or JSON[JSON].
2. Binary data in standard formats.
3. Data with DFDL descriptors.

Textual XML data is the most successful data interchange standard to date. All such data are by definition new, by which we mean created in the XML era. Because of the large overhead that XML tagging imposes, there is often a need to compress and decompress XML data. However, there is a high-cost for compression and decompression that is unacceptable to some applications. Standardized binary data formats are also relatively new, and are suitable for larger data because of the reduced costs of encoding and more compact size. Examples of standard binary formats are data described by modern versions of ASN.1[ASN1], or the use of XDR [XDR]. These techniques lack the self-describing nature of XML-data. Scientific formats, such as NetCDF[NetCDF] and HDF[HDF] are used by some communities to provide self-describing binary data. There are also standardized binary-encoded XML data formats such as EXI [EXI].

It is an important observation that both XML format and standardized binary formats are *prescriptive* in that they specify or prescribe a representation of the data. To use them your applications must be written to conform to their encodings and mechanisms of expression.

DFDL suggests an entirely different scheme. The approach is *descriptive* in that one chooses an appropriate data representation for an application based on its needs and one then describes the format using DFDL so that multiple programs can directly interchange the described data. DFDL descriptions can be provided by the creator of the format, or developed as needed by third parties intending to use the format. That is, DFDL is not a format for data; it is a way of describing any data format. DFDL is intended for data commonly found in scientific and numeric computations, as well as record-oriented representations found in commercial data processing.

DFDL can be used to describe legacy data files, to simplify transfer of data across domains without requiring global standard formats, or to allow third-party tools to easily access multiple formats. DFDL can also be a powerful tool for supporting backward compatibility as formats evolve.

DFDL is designed to provide flexibility and also permit implementations that achieve very high levels of performance. DFDL descriptions are separable and native applications do not need to use DFDL libraries to parse their data formats. DFDL parsers can also be highly efficient. The DFDL language is designed to permit implementations that use lazy evaluation of formats and to support seekable, random access to data. The following goals can be achieved by DFDL implementations:

- Density. Fewest bytes to represent information (without resorting to compression). Fastest possible I/O.
- Optimized I/O. Applications can write data aligned to byte, word, or even page boundaries and to use memory-mapped I/O to insure access to data with the smallest number of machine cycles for common use cases without sacrificing general access.

DFDL can describe the same types of abstract data that other binary or textual data formats can describe and, furthermore, it can describe almost any possible representation scheme for those data. It is the spirit of DFDL to support canonical data descriptions that correspond closely to the original in-memory representation of the data, and also to provide sufficient information to write as well as to read the given format.



## 1.1 Why is DFDL Needed?

The question arises of why DFDL is needed in an era when there are so many standard data formats available. Ultimately, it is because there are a number of social phenomena in the way software is developed that have lead to the situation today where DFDL is needed to standardize descriptions of diverse data formats.

First, programs are very often written speculatively, that is, without any advance understanding of how important they will become. Given this situation, little effort is expended on data formats since it remains easier to program the I/O in the most straightforward way possible with the programming tools in use. Even something as simple as using an XML-based data format is harder than just using the native I/O libraries of a programming language.

In time, however, it is realized that a software program is important because either many people are using it, or it has become important for business or organizational needs to start using it in larger scale deployments. At that point it is often too late to go back and change the data formats. For example, there may be real or perceived business costs to delaying the deployment of a program for a rewrite just to change the data formats, particularly if such rewriting will reduce the performance of the program and increase the costs of deployment. (It takes *longer* to program, but at least it's *slower* when you are done☺)

Additionally, the need for data format standardization for interchange with other software may not be clear at the point where a program first becomes 'important'. Eventually, however, the need for data interchange with the program becomes apparent.

The above phenomena are not something that is going away any time soon. There are, of course, efforts to smoothly integrate standardized data format handling into programming languages. Nevertheless, we see a critical role for DFDL since it allows after-the-fact description of a data format.

## 1.2 What is DFDL?

DFDL is a language for describing data formats. A DFDL description allows data to be read from its native format and to be presented as an instance of an information set or indeed converted to the corresponding XML document. DFDL also allows data to be taken from an instance of an information set and written out to its native format.

DFDL achieves this by leveraging W3C XML Schema Definition Language (XSDL) 1.0. [[XSDL](#)]

An XML schema is written for the logical model of the data. The schema is augmented with special DFDL annotations. These annotations are used to describe the native representation of the data. This is an established approach that is already being used today in commercial systems.

### 1.2.1 Simple Example

Consider the following XML data:

```
<w>5</w>
<x>7839372</x>
<y>8.6E-200</y>
<z>-7.1E8</z>
```

The logical model for this data can be described by the following fragment of an XML schema document that simply provides description of the name and type of each element:

```
<xs:complexType name="example1">
  <xs:sequence>
    <xs:element name="w" type="xs:int"/>
    <xs:element name="x" type="xs:int"/>
    <xs:element name="y" type="xs:double"/>
    <xs:element name="z" type="xs:float"/>
  </xs:sequence>
</xs:complexType>
```

Now, suppose we have the same data but represented in a non-XML format. A binary representation of the data could be visualized like this (shown as hexadecimal):

```
0000 0005 0077 9e8c
169a 54dd 0a1b 4a3f
ce29 46f6
```

To describe this in DFDL, we take our original XML schema document that described the data model and we annotate the type definition as follows:

```
<xs:complexType>
  <xs:sequence>
    <xs:element name="w" type="xs:int">
      <xs:annotation>
        <xs:appinfo source="http://www.ogf.org/dfdl/">
          <dfdl:element representation="binary"
            binaryNumberRep="binary"
            byteOrder="bigEndian"
            lengthKind="implicit"/>
        </xs:appinfo>
      </xs:annotation>
    </xs:element>
    <xs:element name="x" type="xs:int">
      <xs:annotation>
        <xs:appinfo source="http://www.ogf.org/dfdl/">
          <dfdl:element representation="binary"
            binaryNumberRep="binary"
            byteOrder="bigEndian"
            lengthKind="implicit"/>
        </xs:appinfo>
      </xs:annotation>
    </xs:element>
    <xs:element name="y" type="xs:double">
      <xs:annotation>
        <xs:appinfo source="http://www.ogf.org/dfdl/">
          <dfdl:element representation="binary"
            binaryFloatRep="ieee"
            byteOrder="bigEndian"
            lengthKind="implicit"/>
        </xs:appinfo>
      </xs:annotation>
    </xs:element>
    <xs:element name="z" type="xs:float">
      <xs:annotation>
        <xs:appinfo source="http://www.ogf.org/dfdl/">
          <dfdl:element representation="binary"
            byteOrder="bigEndian"
            lengthKind="implicit"
            binaryFloatRep="ieee" />
        </xs:appinfo>
      </xs:annotation>
    </xs:element>
  </xs:sequence>
</xs:complexType>
```

This simple DFDL annotation expresses that the data are represented in a binary format and that the byte order will be big endian. This is all that a DFDL parser needs to read the data.

Consider if the same data are represented in a text format:

```
5,7839372,8.6E-200,-7.1E8
```

Once again, we can annotate the same data model, this time with properties that provide the character encoding, the field separator (comma) and the decimal separator (period):

```
<xs:complexType>
```

```

<xs:sequence>
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:sequence encoding="UTF-8" separator="," />
    </xs:appinfo>
  </xs:annotation>
  <xs:element name="w" type="xs:int">
    <xs:annotation>
      <xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:element representation="text"
          encoding="UTF-8"
          textNumberRep="standard"
          textNumberPattern="####0"
          textStandardDecimalSeparator="."
          lengthKind="delimited"/>
      </xs:appinfo>
    </xs:annotation>
  </xs:element>
  <xs:element name="x" type="xs:int">
    <xs:annotation>
      <xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:element representation="text"
          encoding="UTF-8"
          textNumberRep="standard"
          textNumberPattern="#####0"
          textStandardDecimalSeparator="."
          lengthKind="delimited"/>
      </xs:appinfo>
    </xs:annotation>
  </xs:element>
  <xs:element name="y" type="xs:double">
    <xs:annotation>
      <xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:element representation="text"
          encoding="UTF-8"
          textNumberRep="standard"
          textNumberPattern="0.0E+000"
          textStandardDecimalSeparator="."
          lengthKind="delimited"/>
      </xs:appinfo>
    </xs:annotation>
  </xs:element>
  <xs:element name="z" type="xs:float">
    <xs:annotation>
      <xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:element representation="text"
          encoding="UTF-8"
          textNumberRep="standard"
          textNumberPattern="0.0E0"
          textStandardDecimalSeparator="."
          lengthKind="delimited"/>
      </xs:appinfo>
    </xs:annotation>
  </xs:element>
</xs:sequence>
</xs:complexType>

```

Many properties are repeatedly expressed in the example for the sake of simplicity. Later sections of this specification will define the mechanisms DFDL provides to avoid this repetitiveness.

### 1.3 What DFDL is not

DFDL maps data from a non-XML representation to an instance of an information set. This can be thought of as a data transformation. However, DFDL is not intended to be a general transformation language and, in particular, DFDL does not intend to provide a mechanism to map data to arbitrary XML models. There are two specific limitations on the data models that DFDL can work to:

1. DFDL uses a subset of XML Schema, in particular, you cannot use XML attributes in the data model.
2. The order of the data in the data model must correspond to the order and structure of the data being described.

This latter point deserves some elaboration. The XML schema used must be suitable for describing the physical data format. There must be a correspondence between the XML schema's constructs and the physical data structures. For example, generally the elements in the XML schema must match the order of the physical data. DFDL does allow for certain physically unordered formats as well.

The key concept here is that when using DFDL, you do not get to design an XML schema to your preference and then populate it from data. That would involve describing the data format, and describing a transformation for mapping it to the XML schema you have designed. DFDL is only about the format part of this problem. There are other languages, such as XSLT, which are for transformation. In DFDL, you describe only the format of the data, and this format constrains the nature of the XML schema you must use in its description.

### 1.4 Scope of version 1.0

The goals of version 1.0 are as follows:

1. Leverage XML technology and concepts
2. Support very efficient parsers/formatters
3. Avoid features that require unnecessary data copying
4. Support round-tripping, that is, read and write data in a described format from the same description
5. Keep simple cases simple
6. Simple descriptions should be "human readable" to the same degree that XSDL is.

The general features of version 1.0 are as follows:

- a) Text and binary data parsing and unparsing
- b) Validate the data when parsing and unparsing using XSDL validation.
- c) Defaulted input and output for missing representations
- d) Reference – use of the value of a previously read element in subsequent expressions
- e) Choice – capability to select among format variations
- f) Hidden sequence of elements – A description of an intermediate representation whose corresponding Infoset is not exposed in the final result.
- g) Basic Math – in DFDL expressions
- h) Out-of-type value handling (e.g., The string value 'NIL' to indicate nil for an integer)
- i) Speculative parsing to resolve uncertainty.
- j) Very general parsing capability: Lookahead/Push-back

Version 1.0 of DFDL is a language capable of expressing a wide range of binary and text-based data formats.

DFDL is capable of describing binary data as found in the data structures of COBOL, C, PL1, Fortran, etc. In particular, it is able to describe repeating sub-arrays where the length of an array is stored in another location of the structure.

DFDL is capable of describing a wide variety of textual data formats such as HL7, X12, and SWIFT. Textual data formats often use syntax delimiters, such as initiators, separators and terminators to delimit fields.

DFDL has certain composition properties. I.e., two formats can be nested or concatenated and a working format results.

The following topics have been deferred to future versions of the standard:

- Extensibility: There are real examples of proprietary data format description languages that we use as our base of experience from which to derive standard DFDL. However, there are no examples of extensible format description languages. Therefore, while extensibility is desirable in DFDL, there is not yet a base of experience with extensibility from which to derive a standard.
- Rich Layering: Some formats require data to be described in multiple passes. Combining these into one DFDL schema requires very rich layering functionality. In these layers one element's value becomes the representation of another element. DFDL V1.0 allows description of only a limited kind of layering.

## 1.5 Related standards

1. Prescriptive systems:
  - a. Javascript Object Notation (JSON) [[JSON](#)]
  - b. EXI (binary XML) [[EXI](#)]
  - c. Thrift [[Thrift](#)]
  - d. Avro [[AVRO](#)]
  - e. ASN.1 with any of the prescribed encoding rules: Basic Encoding Rules (BER), Distinguished Encoding Rules (DER), Canonical Encoding Rules (CER) [[ASN1CER](#)] or Packed Encoding Rules (PER) [[ASN1PER](#)]
2. Descriptive systems:
  - a. ASN1 Encoding Control Notation (also known as ITU-T X.692) [[ASN1ECN](#)]
  - b. BFD: Binary Format Description (BFD) Language [[BFD](#)]

## 2. Notational and Definitional Conventions

The key words *must*, *must not*, *required*, *shall*, *shall not*, *should*, *should not*, *recommended*, *may*, *and may not* in this document are to be interpreted as described in [RFC2119]. The terms *required* and *optional* in this document are used as described in Section 3 (Glossary). Note that for reasons of clarity these words are not always capitalized in this document.

Deleted: and optional

Examples are for illustration purposes only and for clarity they will often not include all the necessary DFDL properties.

### 2.1 Failure Types

Where the phrase "must be consistent with" is used, it is assumed that a conforming DFDL implementation must check for the consistency and issue appropriate diagnostic messages when an inconsistency is found.

There are several kinds of failures that can occur when a DFDL processor is handling data and/or a DFDL schema.

### 2.2 Schema Definition Error (SDE)

When the DFDL schema itself contains an error, it implies that the DFDL processor cannot process data because the DFDL schema is not meaningful. It may be ambiguous, or contain conflicting definitions. Equivalently, we can say that there is no possible data that conforms to the schema; hence, the schema cannot be meaningful. All conforming DFDL processors must detect all Schema Definition Errors, and must issue some kind of appropriate diagnostic message. The behavior of a DFDL processor after a Schema Definition Error is detected is out of scope for this specification.

When a Schema definition error can be detected statically, that is given only the schema, it is desirable, though not required by the DFDL standard, that such errors be detected and diagnostic messages issued before any data are processed. Of course not all Schema Definition Errors can be detected without reference to data as some representation properties may obtain their values from the data (see also section 2.3.1 Ambiguity of Data Formats).

The expression language included within DFDL is strongly, statically type checkable. This means that type checking of expressions can be performed without processing data, and implementations are encouraged to perform this checking statically so that Schema Definition Errors having to do with type inconsistencies can be detected before processing data.

Note that Schema Definition Errors cannot be suppressed by points of uncertainty.

#### 2.2.1.1 Schema Component Constraint: Unique Particle Attribution

A DFDL processor MUST implement the Schema Component Constraint: Unique Particle Attribution defined in *XML Schema Part 1: Structures* [XSDLV1] that applies to the DFDL schema subset.

Two elements **overlap** if

- They are both element declaration particles whose declarations have the same name and target namespace.

A schema will violate the unique attribution constraint if it contains two particles which overlap and which either

- Are both in the particles of a *choice* group

Or

- Either may validate adjacent information items and the first has XSDL minOccurs less than XSDL maxOccurs.

### 2.3 Processing Errors

If a DFDL schema contains no Schema Definition Errors, then there is the additional possibility of a *processing error* when processing data using a DFDL schema. A processing error occurs if the

data does not conform to the format described by the schema, that is to say, the data is not well-formed relative to the schema.

Processing errors can be suppressed by a point of uncertainty. See section 9.3.3.

It is expected that DFDL implementations will provide additional implementation-defined mechanisms for dealing with effective processing errors, such as the means of specifying retry points or the means of skipping some data so as to recover from the error in some way.

Exceptions that occur in the evaluation of the DFDL expression language are processing errors.

Non-conformance with the XSDL *minOccurs* or *maxOccurs* constraints is either a processing error or only a validation error depending on the settings of certain DFDL properties (see section 16 below).

### 2.3.1 Ambiguity of Data Formats

A data format using delimiters may be ambiguous if the delimiters are not distinct, and a data format description which has fixed data requirements (that is, where some elements have fixed values) may be ambiguous even with fixed-length elements.<sup>1</sup>

If the delimiter string values are stored within the data, perhaps as elements of a header part of the data, then this ambiguity certainly cannot be examined until the data is available.

Given an ambiguous grammar, a DFDL implementation may successfully parse a particular input data stream. That is, the part of the schema with the ambiguity may not be exercised by a particular data stream, or the data may parse successfully anyway because the ambiguity may not cause any kind of failure or processing error.

Hence, to insure compatible behavior, DFDL v1.0 implementations MUST NOT detect grammar ambiguities as errors. Implementations are of course free to issue warnings to help users identify these situations, but ambiguity is neither a Schema Definition Error nor a Processing Error.

#### 2.3.1.1 Unparsing Must be Unambiguous

Usually, the behavior of the unparser is symmetric to the behavior of the parser; however, there are cases where the DFDL schema will accept several equivalent representations for the same logical data. In this case it would be ambiguous which of these equivalent representations should be produced by the unparser. The DFDL standard contains representation properties which are used to eliminate this ambiguity. It is a Schema Definition Error if a DFDL schema is being used to unparse data and there is any ambiguity about the representation.

## 2.4 Validation Errors

Logical validation checks are constraints expressed in XSDL, and they apply to the logical values of the infoset. Hence, parsing must successfully construct the infoset from the representation in order for validation checks to be meaningful. This implies that validation errors cannot affect the ability of a DFDL processor to successfully parse or unparse data; that is, validation errors are independent of whether the data is well-formed with respect to the DFDL schema.

DFDL processors may provide both validating and non-validating behaviors on either or both of parse and unparse. (A DFDL implementation could support validate on parse, but not support it on unparse and still be considered conforming.)

Validation on unparsing takes place on the augmented infoset that is created by the unparser as a side-effect of creating the output data stream.

When resolving points of uncertainty (during parsing), validation errors are ignored.

---

<sup>1</sup> A very complex analysis is needed to identify this sort of grammar ambiguity in general. While we believe this may be decidable for DFDL v1.0, future versions of DFDL may add features (such as recursive types) which make this analysis undecidable.

The way a validation error is presented to the execution context of a DFDL processor is not specified by the DFDL language. The validity of an element is recorded in the DFDL Infoset, see Section 4 The DFDL Information Set (Infoset).

The following DFDL schema constructs are allowed in DFDL and are checked when validating:

1. XSDL pattern facet - (for XSD string type elements only)
2. XSDL minLength, maxLength
3. XSDL minInclusive, minExclusive, maxInclusive, maxExclusive
4. XSDL enumeration
5. XSDL maxOccurs

Note that validation is distinct from the checking of DFDL assert or discriminator predicates. When a DFDL discriminator or assert is used to discriminate a choice or other point of uncertainty when parsing, then that assert or discriminator is essential to parsing and it is evaluated irrespective of whether validation is enabled or disabled.

There is also a function `dfdl:checkConstraints` available in the DFDL Expression language. This can be used to explicitly include checking of the XSD facet constraints as part of parsing a specific element. Such checking is part of parsing, and does not create validation errors. See Section 23.5.3 DFDL Functions for details.

## 2.5 Recoverable Error

This error type is used with the `dfdl:assert` annotation when parsing to permit the checking of physical format constraints without terminating a parse. For example, some formats will have redundancy by having known lengths, as well as delimiters. A recoverable error can be issued, using an assert to check a physical length constraint when property `lengthKind` is 'delimited'.

Recoverable errors are independent of validation, and when resolving points of uncertainty, recoverable errors are ignored.

## 2.6 Specific Errors Classified

This section clarifies which errors are Schema Definition Errors and which are processing errors.

The following are processing errors:

- Arithmetic Errors
  - Division by zero
  - Integer Arithmetic Underflow
  - Integer Arithmetic Overflow
  - Note: Floating point math can produce NaN (Not a Number) values. This is not an error, nor are properly typed operations on floating point NaN values.
- Expression Errors
  - Dynamic Type Error – unable to convert to target type
    - Example: non-digits found in string argument to `xs:int(...)` constructor.
    - Note: if a DFDL Implementation cannot distinguish Dynamic Type Errors from Static Type Errors, then a Dynamic Type Error **must** cause a Schema Definition Error
  - Index out of bounds error – index not  $\leq$  number of occurrences, or is  $< 1$ .
    - Note: same error for `dfdl:testBit` if `bitPos` is not 1..8, or for character positions in a string-value
  - Indexing of non-array non-optional element
    - Example: `x[1]` when `x` is declared and has both `minOccurs="1"` and `maxOccurs="1"` explicitly, or by not stating either or both of them.
  - Illegal argument value (correct type, illegal value)
- Parse Errors
  - Delimiter not found
  - Data not convertible to type
  - Assertion failed
  - Discriminator failed

Deleted: should



- Required occurrence not found
  - No choice alternative successfully parsed.
  - Character set decoding failure and `dfdl:encodingErrorPolicy` is 'error'
- Unparse Errors
  - Truncation scenarios where truncation is being disallowed
  - Rounding error – rounding needed but not allowed. (Unparsing)
  - No choice alternative successfully unparsed.
  - Character set encoding failure and `dfdl:encodingErrorPolicy` is 'error'
- Implementation-defined Limit Errors - Implementations can have fixed or adjustable limits that some formats and some data may exceed at processing time. This specification does not further specify what these errors are, but some possible examples are:
  - Data longer than allowed for representation of a given data type
    - Example: exceed maximum length of representation of `xs:decimal` when `dfdl:representation` is "text".
  - Expression references too far back into infoset (parsing)
  - Expression references too far forward into infoset (unparsing)
  - Number of array elements exceeds limit.
  - Regular expression exceeds time limit

The following are Schema Definition Errors, regardless of whether they are detected in advance of processing or once processing begins:

- Errors in XML Schema Construction and Structure
  - See XML Schema Specification Part 1, Section 5.1 [\[XSDLV1\]](#)
- Use of XSD constructs outside of DFDL subset
- Implementation-defined Limitations
  - Use of DFDL schema constructs not supported by this implementation.
    - Example: `xs:choice` is an optional part of the DFDL specification (see section 21). If not supported, it must be rejected as a Schema Definition Error.
    - Example: use of packed-decimal when it is not supported by the implementation.
    - Example: use of `dfdl:assert` when it is not supported by the implementation (See Spec section 21 on DFDL Subsets)
    - Note: Unrecognized DFDL properties or property values can produce a warning and an implementation can attempt to process data despite the warning.
  - Exceeding implementation-dependent limits for schema size/complexity
    - Example: schema too large – simply a limit on how large the schema can be, how many files, how many top-level constructs, etc.
- Schema Not Valid
  - See XML Schema Specification Part 1, Section 5.2 [\[XSDLV1\]](#)
- UPA violation (Unique Particle Attribution)
- Reference to DFDL global definition not found
  - Format definition (`dfdl:defineFormat`)
  - Escape schema definition (`dfdl:defineEscapeScheme`)
  - Variable Definition (`dfdl:defineVariable`)
- DFDL Annotations not well-formed or not valid
- DFDL Annotations Incompatible
  - E.g., `dfdl:assert` and `dfdl:discriminator` at same combined annotation point, or more than one format annotation at an annotation point.
- DFDL Properties and their values
  - Property not applicable to DFDL annotation
  - Property value not suitable for property
  - Property conflict
    - Between Element Reference and Element Declaration
    - Between Element Declaration and Simple Type Definition

- Between Simple Type Definition and Base Simple Type Definition
  - Between Group Reference and Sequence/Choice of Group Definition
- Required property not found
- Expressions
  - Expression syntax error
  - Named child element doesn't exist – E.g., /a/b, and there is no child b in existence.
    - Note: no child possible in the schema is a different error, but also a Schema Definition Error, as /a/b would not have a type in that case.
    - Note: This is an SDE, as schema authors are advised to use fn:exists(...) to test for existence of elements when it is possible that they not exist.
  - Variable read but not defined
  - Variable assigned after read
  - Variable assigned more than once
  - Static Type error – type is incorrect for usage
    - Note: if an implementation is unable to distinguish Static Type Errors from Dynamic Type Errors, then both should cause Schema Definition Errors.
  - Path step definition not found – e.g., /a/n:b but no definition for n:b as local or global element.
  - Not enough arguments for function
  - Expression value is not single node
    - Most DFDL expression contexts require an expression to identify a single node, not an array (aka sequence of nodes). There are a few exceptions such as the fn:count(...) function, where the path expression must be to an array or optional element.
  - Expression value is not array element or optional element.
    - Some DFDL expression contexts require an array or an optional element.
    - Example: The fn:count(...) function argument must be to an array or optional element. It is an SDE if the argument expression is otherwise.
- Regular Expressions
  - Syntax error

## 2.7 Optional Checks and Warnings

A DFDL processor:

- That only implements a DFDL parser does not have to check (for Schema Definition Errors) properties that are solely used when unparsing, though it is recommended that it does so for portability reasons.
- That does not implement some optional features does not have to check properties or annotations required by those optional features, but MUST issue a warning that an unrecognized property or annotation has been encountered.
- Need not check global objects as they may legitimately be incomplete due to properties intended to be supplied based on scoping rules and the context at the point of use.  
There are two exceptions to this, which must be checked:
  - Global simple types that are referenced by prefixLengthType property
  - Global elements that are the document root.

Some situations suggest likely errors, but a DFDL processor cannot be certain. In these situations, a DFDL processor may issue warnings to assist a DFDL schema author in identifying likely errors. An important case of this is when the DFDL processor encounters a schema component and annotation where there are explicitly properties that are not relevant to the component as defined. Depending on the specifics of the component and property the DFDL processor can or must take certain actions. If the:

- Property is not applicable to the component's DFDL annotation.

- Schema definition error. Example is lengthKind on xs:sequence.
- Property is not applicable because of simple type.
  - Warning (optional). Example is calendarPatternKind on xs:string.
- Property is not applicable because of another DFDL property setting.
  - Warning (optional). Example is binaryNumberRep when representation is text.
- Invalid value for a property that is unused or ignored.
  - Warning (optional). Example is dfdl:lengthKind is not 'explicit' but dfdl:length is an expression and that expression contains invalid syntax.

### 3. Glossary

**Adjacent** - Two parts of the input/output stream are adjacent if they are at consecutive addresses.

**Addressable Unit** - This is the unit of storage that makes up the input or output stream holding the representation of the data. The units are bits, bytes, or characters.

**Annotation point** - A location within a DFDL schema where DFDL annotation elements are allowed to appear.

**Applicable properties** - All the DFDL properties that apply to that type of schema construct. For example all the DFDL properties that apply to an `xs:simpleType`.

**Array** - The set of adjacent elements whose XSDL element declaration specifies the potential for it to have more than one occurrence (XSD property `maxOccurs > '1'` or 'unbounded'). Of course any given array can have any number of element occurrences, including zero elements or exactly 1 element as long as the occurrence constraints are met. If XSD property `maxOccurs` is 'unbounded' then there is no constraint to the maximum number of occurrences, though implementations may have implementation-defined maximum capabilities. An optional element (where XSD property `maxOccurs` is '1', `minOccurs` is '0') is not considered to be an array as described in this document. Note that a sequence is not to be confused with an array. A sequence is a complex tuple type for an element; the children of a sequence can be of different types. All elements of an array have the same type and have the same information item members except for the value member.

**Array Element** - An element declaration or reference with XSD property `maxOccurs > '1'` or 'unbounded'.

**Augmented Infoset** - When unparsing one begins with the DFDL schema and conceptually with the logical infoset. As the values of items are filled in by defaulting, and by use of the DFDL `outputValueCalc` property (including on hidden items), these new item values augment the infoset. The resulting infoset is called the augmented infoset.

**Binary** - There are two meanings for this word depending on context.

- Data is divided into two broad categories of representations, which are *text* and *binary*. Hence, binary representation includes any kind of non-text representation.
- Within binary (not text) data, we distinguish base-10 representations which are called *packed decimal*, from base-2 representations which are called *binary*. The common *twos-complement* representation used for signed integers is a base-2 binary representation.

**Binary Representation** - Of type `xs:hexBinary`, or of other type with property `fdfl:representation` 'binary'. Note that type `xs:string` can never have binary representation.

**Bit Order** - Within a binary integer, if the most-significant bit is assigned bit position 1, then the bit order is *most-significant-bit first*. If the least-significant bit is assigned bit position 1, then the bit order is *least-significant-bit first*. When the bit order is most-significant-bit first, then the least-significant bit of byte N is considered to be adjacent to the most-significant bit of byte N+1. When the bit order is least-significant-bit first, then the most-significant bit of byte N is considered to be adjacent to the least-significant bit of byte N+1.

**Bit Position** - The data stream is assumed to be a collection of consecutively numbered unsigned bytes. Each byte is a numeric value from 0 to 255. The bits of a byte are referred to by their numerical significance as the  $2^n$  bit, for  $n$  from 0 to 7. Hence, the byte value  $255 = 2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0$ . The  $2^7$ -bit is the most-significant bit, and the  $2^0$ -bit is the least significant bit. The bits within each byte are assigned numbered bit positions 1 to 8 according to the *bit order*. Given a bit-order, every bit in the data stream has a unique bit position.

**Bit String** - The ordered set of bits from a first bit with *bit position* N, to bit position N+M is a bit string of length M bits.

**Byte** - The term "byte" refers to an 8-bit octet. Can also refer to an integer with value from 0 to 255 inclusive.

**CCSID** - see Coded Character Set Identifier.

**Character** - An ISO10646 [\[ISO10646\]](#) character having a unique character code as its identifier. This concept is independent of font, typeface, size, and style, so 'F', 'F', 'F', are all the same character 'F'.

**Character Code** - The canonical integer used to identify a character in the ISO10646 [\[ISO10646\]](#) standards. This number identifies the character, but can be independent of any specific character set encoding of the character. Example: The 'f' character known in Unicode [\[Unicode\]](#) as LEFT CURLY BRACKET has character code U+007B. However, depending on the character set encoding, the value 0x7B may or may not appear in the representation of that character.

**Character Set** - An abstract set of characters that are assigned (or mapped to) a representation by a particular character set encoding. For most character set encodings their character set is a subset of the Unicode character set.

**Character Set Encoding** - Often abbreviated to just 'encoding'. A specific representation of a character set as bytes or bits of data. A character set encoding is usually identified by a standard character set encoding name or a recognized alias name, or by a coded character set identifier or CCSID [\[CCSID\]](#). These identifiers are standardized. The names and aliases are standardized by the IANA [\[IANA\]](#) (where unfortunately, they are called character set names). CCSIDs are an industry standard. Examples of character set encoding names are UTF-8, USASCII, GB2312, ebcdic-cp-it, ISO-8859-5, UTF-16BE, Shift\_JIS. There are also additional DFDL standard character set encodings, see *DFDL Standard Encoding*. The DFDL standard also allows for implementation-defined character set encodings to be supported..

**Character Width** - The number of code units or alternatively the number of bytes or bits used to represent a character in a specific character set encoding is called the character width. Encodings are either fixed width (all characters encoded using the same width), or variable-width (different characters are encoded using different widths). For example the UTF-32 character set encoding has 4-byte character width, whereas USASCII has a 1-byte character width. UTF-8 is variable width, and any specific character has width 1, 2, 3, or 4 bytes. See also *Fixed-Width Character Encoding* and *Variable-Width Character Encoding*

**Code Point** - The integer that identifies a character within a character set encoding. A code point is represented by one or more code units. When a character set is fixed width, then there is no distinction between a code unit and a code point. For unicode character set encodings, there is no distinction between a character code and a code point. Examples:

- € - character code U+20AC
  - IBM01148 encoding - the code point is 0x9F, and this encoding is fixed width so there is no distinction between the code point 0x9F and the code unit 0x9F that represents the encoded character.
  - UTF-8 encoding - there is no distinction between the character code 0x20AC and the code point 0x20AC. However, it is represented by 3 code units 0xE2 0x82 0xAC

**Code Unit** - When a character set encoding uses differing variable width representations for characters, the units making up these variable width representations are called code units. For example the UTF-8 encoding uses between 1 and 4 code units to represent characters, and for UTF-8, the individual code units are single bytes. DFDL's interpretation of the UTF-16 encoding is either fixed or variable width. When format property `dfdl:utf16Width 'variable'` then UTF-16 is variable width and this encoding uses either one or two code units per character, but in this case each individual code unit is a 16-bit value. When a character set is fixed width, then there is no distinction between a code unit and a code point.

**Coded Character Set Identifier (CCSID)** - An alternate identifier of a character set encoding. Originally created by IBM, CCSIDs are a broadly used industry standard. See [\[CCSID\]](#).

**Component** - A construct within a DFDL schema that may contain a DFDL annotation.

**Content** - The content is the bits of data that are interpreted to compute a logical value.

**Content Model** - Used in describing the syntactic structure of XSD and DFDL annotations of it. An element of a schema can have empty, simple, or element-only content. An element

declaration for an element of complex type containing a `xs:sequence` element is said to have a sequence in its content model.

**Contiguous** - An element has a contiguous representation if all parts of its representation are adjacent in the input/output stream. Most simple types have contiguous representations naturally. Groups containing elements that are themselves contiguous are also considered to have contiguous representations irrespective of alignment fill or padding of any kind that exists within the group. Similarly, arrays of elements that are themselves contiguous are also contiguous. An example of a non-contiguous representation would be a nillable element, where a flag is used to determine whether or not the element is nil, and the location of that flag is not adjacent to the value representation.

**Count** - The number of occurrences of an element.

**Data Stream** - Data where the format is being described by a DFDL schema. This use of 'stream' implies only that there is a numbering scheme that specifies a unique *bit position* for every bit within the data. This use of 'stream' does not imply anything about whether the data is persistently stored or not, nor does it imply anything about whether there are sequential or random access capabilities for access to the data.

**DBCS** - See Double-Byte Character Set

**Decimal** - This term is used several different ways distinguished by context:

1. Base 10. When data has text representation, a decimal number has base-10 digits.
2. Type `xs:decimal` - a logical type of number that has an integer component and an optional base-10 fractional component. This type subsumes all integer types, as they are of type `xs:decimal` but with the further restriction that the fractional part doesn't exist. Note that a base-10 fraction has different rounding properties than a base-2 or floating point numeric fraction; hence, `xs:decimal` is the type commonly used to represent currency/money in data.
3. Packed Decimal - A binary data representation. See separate glossary entry below.

**Defining Annotations** - The annotation elements `dfdl:defineFormat`, `dfdl:defineVariable`, and `dfdl:defineEscapeScheme`

**Delimiter** - A character or string used to separate, or mark the start and end of, items of data. In DFDL, `dfdl:lengthKind` 'delimited' scans the data for initiators, separators, and terminators.

**Delimiter scanning** - When parsing, the process of scanning for a specific item in the input data which either marks the end of an item or the beginning of a subsequent item. Delimiter scanning also takes into account escape schemes so as to allow the delimiters to appear within data if properly escaped.

**DFDL** - Data Format Description Language

**DFDL Processor** - A program that uses DFDL schemas in order to process data described by them.

**DFDL Schema** - An XML schema containing DFDL annotations to describe data format.

**DFDL Standard Encoding** - A character set for which there is no IANA name or CCSID but the name and definition of which DFDL implementations must agree on. See Section 34 [Appendix D: DFDL Standard Encodings](#).

**Double-Byte Character Set (DBCS)** - A character set encoding where each character code consists of one code unit which uses exactly 2 bytes.

**Dynamic extent** - This is a characteristic of the data stream. When parsing data corresponding to a schema component, the collection of bits within the data stream that contain any aspect of the representation of that schema component make up the component's dynamic extent.

**Dynamic scope** - This is a characteristic of parts of the DFDL schema. When a definition or declaration contains or references another declaration or definition, then the contained definition or declaration is said to be in the dynamic scope of the enclosing one. The important characteristic of dynamic scoping is that it traverses references. When parsing, the dynamic

**Deleted:** Appendix D: DFDL Standard Character Set Encodings...

scope of an element declaration includes all definitions and declarations used as part of parsing that element.

**Element** - A part of the data described by an element declaration in the schema and presented as an element information item in the infoset.

**Encoding** - See Character Set Encoding.

**Explicit properties** - The explicit properties are the combination of any defined locally on the annotation and any defined by a `dfdl:defineFormat` annotation referenced by a local `dfdl:ref` property.

**Fixed-Length Element** - An element of specified length where `dfdl:lengthKind` is 'explicit' but `dfdl:length` is not an expression, or `dfdl:lengthKind` is 'implicit' (of simple type only). Note that choice branches where `dfdl:choiceLengthKind` is 'explicit' are also referred to as 'fixed-length' but are not necessarily elements.

**Fixed-Width Character Encoding** - A character set encoding where all characters are encoded using a single code unit for their representation. Note that a code unit is not necessarily a single byte. It may be more than one byte, or some number of bits less than a byte. Examples of different fixed widths are:

- 1-byte wide: ASCII, ebcdic-cp-us, iso-8859-1. See also SBCS (Single-Byte Character Set)
- 2-bytes wide: UTF-16 when `dfdl:utf16Width` is 'fixed'. See also DBCS (Double-Byte Character Set)
- 4-bytes wide: UTF-32.
- 7-bits wide: X-DFDL-US-ASCII-7-BIT-PACKED<sup>2</sup>.

**Fixed Array Element** - An array element where XSDL `minOccurs` is equal to XSDL `maxOccurs`.

**Format annotations** - The annotation elements `dfdl:format`, `dfdl:element`, `dfdl:simpleType`, `dfdl:group`, `dfdl:sequence`, `dfdl:choice`, and `dfdl:escapeScheme`.

**Format property** - A DFDL property carried on a DFDL format annotation.

**Framing** - The term used to describe the delimiters, length fields, and other parts of the data stream which are present, and may be necessary to determine the length or position of the content of an element.

**Implementation-defined feature** - A feature where the implementation has discretion in how it is performed, and the implementation must document how it is performed.

**Implementation-dependent feature** - A feature where the implementation has discretion in how it is performed, but the implementation is not required to document how the feature is performed.

**Index** - The position of an occurrence in a count, starting at 1.

**Item** - A DFDL information set consists of a number of **information items**; or just *items* for short.

**Least-Significant Bit** - Often abbreviated to **LSB**. In a binary integer the least significant bit is the bit having the least place value. Within an 8-bit unsigned byte, the bit with place value  $2^0$  (or 1) is the least significant bit.

**Length** - When discussing data items and their representations, the term 'length' is used to refer to the measure of the size of the representation of an item in units of bits, bytes, or characters. The length of an array is the number of bits, bytes, or characters making up its representation, and has nothing to do with the number of occurrences of the array. Any element occurrence has length. Only array elements and optional elements have numbers of occurrences other than 1.

**Lexical scope** - In a DFDL Schema document, the lexical scope of any element is the collection of schema declarations, definitions, and annotations contained within the element textually.

<sup>2</sup> X-DFDL-US-ASCII-7-BIT-PACKED is a DFDL standard encoding, which uses the US-ASCII characters, but each code unit is stored occupying only 7 bits, not a whole 8-bit-byte. DFDL standard encodings are defined in a separate specification. See Section 34 Appendix D: DFDL Standard Encodings.

**Local properties** – Local properties are the properties defined on an annotation in either short, attribute or element form

**Logical layer** - A DFDL Schema with all the DFDL annotations ignored is an ordinary XSDL schema. The logical structure described by this XSDL is called the DFDL *logical layer*.

**Most-Significant Bit** - Often abbreviated to **MSB**. In a binary integer the most significant bit is the bit having the greatest place value. Within an 8-bit unsigned byte, the bit with place value  $2^7$  is the most significant bit.

**Nibble** - 4 bits. A single hexadecimal digit (0 to 9, A to F) is often referred to as a nibble as it can be represented in exactly 4 bits.

**Node** - The term Node is a shorter equivalent to Element Information Item of the DFDL Infoset described in Section 4.1.2 Element Information Items.

**Non-representation property** – A format property that is not a representation property, specifically `dfdl:ref`, `dfdl:hiddenGroupRef`, `dfdl:choiceBranchKey`, `dfdl:choiceDispatchKey`, `dfdl:inputValueCalc`, `dfdl:outputValueCalc`. See also representation property.

**Occurrence** - An instance of an element in the data, or an item in the DFDL Infoset.

**Optional Element** - An element declaration or reference where XSDL `minOccurs` is equal to zero.

**Optional Occurrence** - An occurrence with an index greater than XSDL `minOccurs`.

**Packed decimal** – A physical representation of a decimal and integer numbers where each digit is packed into one nibble (4 bits) of a byte. There are several variants, some also include a sign nibble and some include a padding nibble. The term covers all the following enums of the `dfdl:binaryNumberRep` and `dfdl:binaryCalendarRep` properties – 'packed' (IBM 390 packed), 'bcd' (standard binary coded decimals or BCDs) and 'ibm4690Packed' (IBM 4690 packed).

**Potentially represented** - An element declaration in the schema describes a *potentially represented* item if that element declaration does not have a `dfdl:inputValueCalc` property. Whether the element declaration describes an occurrence that is actually represented or not depends on whether the element declaration is for an optional element, and whether the element has a corresponding value in the augmented infoset.

**Physical Layer** - A DFDL Schema adds DFDL annotations onto an XSDL language schema. The annotations describe the physical representation or *physical layer* of the data.

**Point of Uncertainty** - A point of uncertainty occurs in the data stream when there is more than one schema component that might occur at that point.

**Representation property** - A format property that is used to describe a physical characteristic of a component. Such a property will apply to one or more grammar regions of the component. See also non-representation property.

**Required Element** - An element declaration or reference where XSDL `minOccurs` is greater than zero.

**Required Occurrence** - An occurrence with an index less than or equal to XSDL `minOccurs`.

**Required Property** – A DFDL property that must have a value. The required properties for each `xs:schema` component are listed in the Property Precedence tables in section 23.

**Resolved set of annotations** - When DFDL annotations appear on a group reference and the sequence or choice of the referenced global group, or appear among an element reference, an element declaration, and its type definition, then they are combined together and the resulting set of annotations is referred to as the *resolved set of annotations* for the schema component.

**SBCS** - See Single Byte Character Set

**Scan** – Examine the input data looking for delimiters such as separators and terminators, or matches to regular expressions.

**Single-Byte Character Set (SBCS)** - A character set encoding where each character code consists of one code unit which is exactly a single byte (8 bits).



**Schema** - The set of all declarations and definitions in the schema, including all included and imported schemas taken together. This includes both the XSDL declarations and definitions, and the DFDL definitions provided in the top-level DFDL annotations.

**Schema Component Designator** (SCD) - A notation for referring to one of the components of a DFDL Schema. This is being standardized by W3C. See <http://www.w3.org/TR/xmlschema-ref>.

**Schema Definition Order** - The order that the schema components are defined in a schema document.

**Specified length** - An item has specified length when `dfdl:lengthKind` is "implicit" (simple type only), "explicit", or "prefixed".

**Speculative Parsing** - When the parser reaches a point of uncertainty it attempts to parse each option in turn until one is known-to-exist or known-not-to-exist.

**Statement annotations** - The annotation elements `dfdl:assert`, `dfdl:discriminator`, `dfdl:setVariable`, and `dfdl:newVariableInstance`. Also called DFDL Statements.

**Statically** - A DFDL Implementation can analyze a DFDL schema and determine the presence of many kinds of errors. This is called static analysis, compilation of the schema, or determining the presence of the error statically.

**Surrogate Pair** - A Unicode character whose character code value is greater than 0xFFFF can be encoded into variable-width UTF-16BE or UTF-16LE (which are variable-width encodings when the DFDL property `utf16Width` is 'variable'). In this case the representation uses two adjacent code units each of which is called a surrogate, and the pair of which is called a surrogate pair.

**Target length** - When unparsing, the length (in `dfdl:lengthUnits`) of an item's representation is the target length. The length of the content corresponding to a logical data value in the infoset may be shorter or longer than the target length, in which case padding or truncation may be necessary to make the logical data content conform to the target length. Rules for when padding and truncation occur, and how they are applied are specific to simple data types, and are controlled by a number of DFDL format properties.

**Text** - Consisting of characters in some character set encoding. Normally we think of text data as being human readable, but many character set encodings contain special control characters that are not human readable but we call data containing these *text* anyway. The `dfdl:encoding` property is required in order to decode/encode the text.

**Text Representation** - Of type `xs:string`, or of other types (except `xs:hexBinary`) with property `dfdl:representation` 'text'. Note that type `xs:hexBinary` never has text representation. This term specifically refers to the representation of the SimpleContent region being textual.

**Textual** - See Text.

**Twos-Complement** - A very common scheme for representing binary integers within data. A positive integer consisting of N bits is represented as its base-2 absolute value. A negative integer is represented as the complement (all bits inverted) of its absolute value plus 1.

**Unicode** - A character set defined by the Unicode Consortium, and standardized at the International Standards Organization (ISO) as ISO10646.

**Unit** - See Addressable Unit.

**Unpadded length** - This is the length of the content of an item of the infoset, prior to any filling or padding which might be introduced due to `dfdl:lengthKind` "prefixed" or `dfdl:lengthKind` "explicit". It is equal to or smaller than the target length.

**Validity** - A DFDL Infoset is said to be valid with respect to a DFDL schema if each Infoset item is valid with respect to its corresponding DFDL schema component. Validity is about the Infoset and the values it holds. It is independent of the data representation when parsing or unparsing. See Section 2.4 Validation Errors, for a list of the specific value checks that are performed when validating a DFDL Infoset against a DFDL schema.

**Variable-Width Character Encoding** - A character set encoding where characters are encoded using one or more code units for their representation depending on which specific character is being encoded. Examples with their ranges of varying width:

- 1 to 4 bytes: UTF-8
- 1 or 2 16-bit code units: UTF-16 when property `dfdl:utf16Width` is 'variable'
- 1 or 2 bytes: Shift-JIS

**Well-formed** - Data is said to be well-formed with respect to a DFDL schema if a DFDL processor can parse the data into a DFDL Infoset, or a DFDL processor can unparse to that data from a DFDL Infoset. The validity of values in the infoset is not necessary for data to be well-formed.

**Width** - See Character Width.

#### 4. The DFDL Information Set (InfoSet)

This section defines an abstract data set called the **DFDL Information Set (InfoSet)**. Its purpose is to define the abstract data structure that must be provided:

- To an invoking application by a DFDL parser when parsing DFDL-described data using a DFDL Schema;
- To a DFDL unparser by an invoking application when generating DFDL-described data using a DFDL Schema

The DFDL InfoSet contains enough information so that a DFDL schema can be defined that will unparse the infoSet and reparse the resultant datastream to produce the same infoSet.

There is no requirement for DFDL-described data to be valid in order to have a DFDL information set.

DFDL information sets may be created by methods (not described in this specification) other than parsing DFDL-described data.

A DFDL information set consists of a number of **information items**; or just **items** for short. The information set for any well-formed DFDL-described data will contain at least a document information item and one element information item. An information item is an abstract description of a part of some DFDL-described data: each information item has a set of associated named **members**. In this specification, the member names are shown in square brackets, **[thus]**. The types of information item are listed in Section 4.1 [Information Items](#).

The DFDL Information Set does not require or favor a specific interface or class of interfaces. This specification presents the information set as a modified tree for the sake of clarity and simplicity, but there is no requirement that the DFDL Information Set be made available through a tree structure; other types of interfaces, including (but not limited to) event-based and query-based interfaces, are also capable of providing information conforming to the DFDL Information Set.

The terms "information set" and "information item" are similar in meaning to the generic terms "tree" and "node", as they are used in computing. However, the former terms are used in this specification to reduce possible confusion with other specific data models.

The DFDL Information Set is similar in purpose to the XML Information Set [\[XML InfoSet\]](#), however, it is not identical, nor a perfect subset, as there are **important differences**.

##### 4.1 Information Items

An information set contains two different types of information items, as explained in the following sections. Every information item has members. For ease of reference, each member is given a name, indicated **[thus]**.

###### 4.1.1 Document Information Item

There is exactly one **document information item** in the information set, and all other information items are accessible through the **[root]** member of the document information item.

There is no specific DFDL schema component that corresponds to this item. It is a concrete artifact describing the information set.

The document information item has the following members:

**[root]** The element information item corresponding to the root element declaration of the DFDL Schema.

**[dfdlVersion]** String. The version of the DFDL specification to which this information set conforms. For DFDL V1.0 this is 'dfdl-1.0'

**[schema]** String. This member is reserved for future use.

###### 4.1.2 Element Information Items

There is an **element information item** for each value parsed from the non-hidden DFDL-described data. This corresponds to an occurrence of a non-hidden element declaration of simple type in the DFDL Schema and is known as a **simple element information item**.

*what are they? Do you address, specifically?*

**Deleted: [unicodeByteOrderMark]** Enum. When the encoding of the root element of the document is exactly UTF-8, UTF-16, or UTF-32 (or CCSID equivalent), the member value indicates whether the document starts with a Byte-order-mark (BOM), and what the value of the mark was. If there is a BOM at the start of the data stream, then for UTF-8 encoding the value is 'UTF-8'; for UTF-16 encoding the value is 'UTF-16LE' or 'UTF-16BE'; for UTF-32 the value is 'UTF-32LE' or 'UTF-32BE'. If there is no BOM then the member value is empty. When the encoding of the root element of the document is any other encoding, the member value is empty. When unparsing, if this member is not empty and the encoding is UTF-8, UTF-16, or UTF-32, then this member's value is used to determine the specific byte-order mark written, and for UTF-16 and UTF-32, the byte order used when characters are encoded to the output data stream.¶

There is an **element information item** for each explicitly declared structure in the DFDL-described data. This corresponds to an occurrence of an element declaration of complex type in the DFDL Schema and is known as a **complex element information item**.

In this information set, as in an XML document, an array is just a set of adjacent elements with the same name and namespace. (To represent the array explicitly, introduce a new complex type element to contain the array elements only.)

One of the element information items is the [root] member of the document information item, corresponding to the root element declaration of a DFDL Schema and all other element information items are accessible by recursively following its [children] member.

An element information item has the following members:

**[namespace]** String. The namespace, if any, of the element. If the element does not belong to a namespace, the value is the empty string.

**[name]** String. The local part of the element name.

**[document]** The document information item representing the DFDL information set that contains this element. This element is empty except in the root element of an information set.

**[dataType]** String. The name of the XML Schema 1.0 built-in simple type to which the value corresponds. DFDL supports a subset of these types listed in section 5.1 DFDL Subset of XML Schema. In a complex element information item this member has no value.

**[dataValue]** The value in the value space (as defined by XML Schema Part 2: Datatypes [XSDLV1]) of the [dataType] member. In a complex element information item this member has no value. If the [nilled] member is true, then this member has no value.

For information items of datatype xs:string, the value is an ordered collection of unsigned 16-bit integer codepoints each having any value from 0x0000 to 0xFFFF. Where defined, these are interpreted as the ISO646 character codes. Codepoints disallowed by ISO 10646, such as 0xD800 to 0xDFFF are explicitly allowed by the DFDL infoset. DFDL's infoset represents Unicode characters with character codes beyond 0xFFFF by way of surrogate pairs (2 adjacent codepoints) in a manner consistent with the UTF-16 encoding of ISO 10646. The value can have length 0, in which case the value may be referred to as an 'empty string'.

For information items of datatype xs:hexBinary, the value is an ordered collection of unsigned 8-bit bytes each having value from 0 to 255. The length of this collection can be 0 in which case the value may be referred to as an 'empty hexBinary'.

**[nilled]** Boolean. True if the nillable item is nil. False if the nillable item is not nil. If the element is not nillable this member has no value. If this member is true then for a simple element the [dataValue] member has no value, and for a complex element the [children] member has no value. If this member is true then the Infoset item is said to be nil or nilled.

**[array]** Boolean. True if the item is an array, meaning that it corresponds to an element having maxOccurs value greater than 1, or 'unbounded'.

**[children]** An ordered set of zero or more element information items. The order they appear in the set is the order implied by the DFDL Schema. 'Ordered set' is not formally defined here, but two operations are assumed: 'count' gives the number of information items, and 'at (index)' gives the element at ordinal position 'index' starting from 1. In a simple element information item this member has no value. In a document information item this member contains exactly one element information item. If the [nilled] member is true then this member has no value.

**[parent]** The complex element information item which contains this information item in its [children] member. In the root element of an information set this member is empty.

**[schema]** String. A reference to a schema component associated with this information item, if any. If not empty, the value must be an absolute or relative Schema Component Designator [SCD].

COMMENTS START HERE

Random thought: It doesn't have to be. It could be convenient to have a persistent reference back to the root of the tree in every node

This is unexpected to me. It appears that you are defining a format for the DFDL infoset. I was expecting that it would be a purely logical object. Meaning that if an element was a string...

**Deleted:** The codepoints of the string are stored in 'implicit' (also known as logical), left-to-right bidirectional ordering and orientation.

... for example, that in Java I would get a String object. Now that is actually represented in memory is presumably a choice for the JVM. Why does the infoset need a physical representation?

**[valid]** Boolean<sup>3</sup>. True if the element is valid as determined by a DFDL implementation that performs validation checking. A complex element information item is not valid if any of its **[children]** are not valid. Empty if validation is not enabled.

**[unionMemberSchema]**<sup>4</sup> String. For simple element information items, this member contains an SCD reference to the member of the union that matched the value of the element. Empty if validation is not enabled. Empty if the element's type is not a union.

On unparsing, any non-empty values for the **[valid]** or **[unionMemberSchema]** members are ignored. However, in the augmented infoset which is built during the unparse operation **[valid]** will have a value, and **[unionMemberSchema]** may have a value.

#### 4.2 "No Value" - worth disambiguating from "nil"/"nillable"?

Some members may sometimes have the value **no value**, and it is said that such a member has no value. This value is distinct from all other values. In particular it is distinct from the empty string, the empty set, and the empty list, each of which simply has no members.

#### 4.3 DFDL Information Item Order

On parsing and unparsing information items will be presented in the order they are defined in the DFDL Schema.

#### 4.4 DFDL Infoset Object model

By way of illustration, the DFDL information set is presented below as an object model using a Unified Modeling Language (UML) class diagram, augmented using the Object Constraint Language (OCL) [\[UML\]](#).

The structure of the information set follows the Composite design pattern. In case of inconsistency or ambiguity, the preceding discussion takes precedence.

DFDL is able to describe the format of the physical representation for data whose structure conforms to this model. Note that this model allows hierarchically nested data, but does not allow representation of arbitrary connected graphs of data objects.

*I wonder if this concept should be introduced earlier. I have just read about a lot of members that can have empty values. I went back and re-read because I was unsure if I had understood. Also this paragraph is a bit mysterious it says nothing about why this concept is needed*

<sup>3</sup> The purpose of this member is to support construction of a W3C standard Post Schema Validation Infoset (PSVI) from a DFDL Infoset.

<sup>4</sup> Also to support PSVI construction.

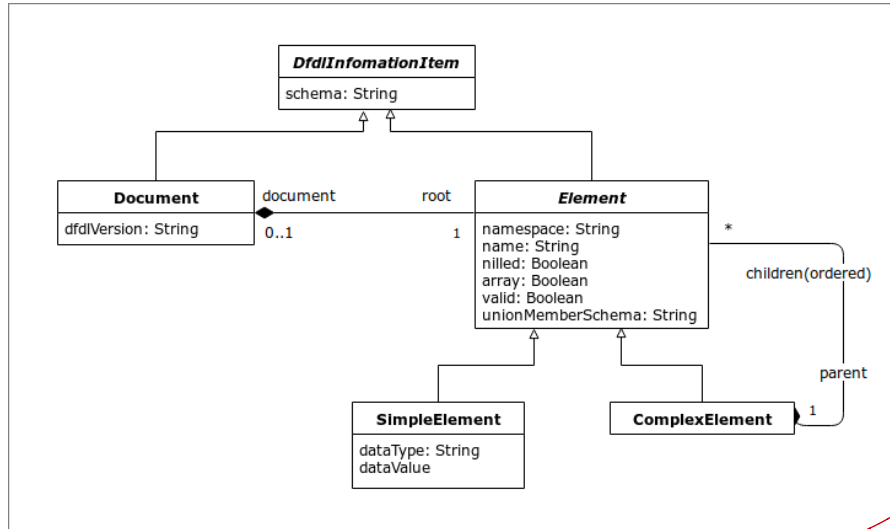


Figure 1 DFDL Infoset Object Model

#### 4.5 DFDL Augmented Infoset

When unparsing, one begins with the DFDL schema and conceptually with the logical infoset. As the values of items are filled in by defaulting, and by use of the `dfdl:outputValueCalc` property (including on hidden items) (see section 17 Calculated Value Properties), these new item values augment the infoset. The resulting infoset is called the augmented infoset.

An element declaration in the schema describes a *potentially represented* item if that element declaration does not have a `dfdl:inputValueCalc` property (see section 17 Calculated Value Properties). Whether the element declaration describes an item that is actually represented or not depends on whether the element declaration is for an optional element and whether the element has a corresponding value in the augmented infoset.

In expressions, the function `dfdl:contentLength()` and `dfdl:valueLength()` can be called to determine the length of an item. If an element declaration is not potentially represented, then these functions are defined to return 0.

When unparsing, an element declaration and the infoset are considered as follows. An implementation may use any technique consistent with this algorithm:

a) If the element declaration has a `dfdl:outputValueCalc` property then the expression which is the `dfdl:outputValueCalc` property value is evaluated and the resulting value becomes the value of the element item in the augmented infoset. Any pre-existing value for the infoset item is superseded by this new value.

References to other augmented infoset items from within the `dfdl:outputValueCalc` expression must obtain their values from the augmented infoset directly (when the value is already present) or by recursively using these methods (a) and (b) as needed.

b) If the element declaration has no corresponding value in the augmented infoset, and the element declaration is for a *required* occurrence, and it *has a default value specified*, then an element item having the default value is created in the augmented infoset.

c) If any Infoset item's value is requested recursively as a part of (a) above and (a) does not apply, and the corresponding value is not present, and (b) does not apply then it is a processing error.

I like this. Consider leading with it?

I'm not keen on "Augmented" hard to change so not suggesting that, but important to get the concept right. Is it reasonable to say that the Infoset is like a Class and an Augmented Infoset is a corresponding Object?

This is super-confusing. I do not understand the concept of "potentially represented" or how it relates to actually represented.

I think there's a jump here. The previous section on the "infoset" was talking about logical members represented: [member]. In this section on augmented infoset we seem to suddenly be talking in terms of properties like "dfdl:foo". I think the transition needs to be smoothed a little.

Makes sense but I'm confused why are you talking re now?

Given this augmented infoset, then if the potentially represented element declaration has a corresponding infoset item then that item is converted to its representation according to its DFDL properties. If the element declaration is for a required occurrence, and there is no value in the augmented infoset then it is a processing error.

Because rule (a) above is used even if the augmented infoset item already exists and has a value, it is possible for a `dfdl:outputValueCalc` expression to be evaluated multiple times. DFDL implementations are free to cache values and avoid this repeated evaluation for efficiency, as the semantics of DFDL require that the `dfdl:outputValueCalc` expression return the same value every time it is evaluated.

*I think a big part of why this is confusing is that I don't think I have yet been introduced to `dfdl:outputValueCalc`*

## 5. DFDL Schema Component Model

When using DFDL, the format of data is described by means of a *DFDL Schema*.

The DFDL Schema Component Model is shown in conceptual UML in Figure 2. First we show the model for elements, groups and the top of the type hierarchy.

The shaded boxes have direct corresponding element syntax and therefore appear in DFDL schema

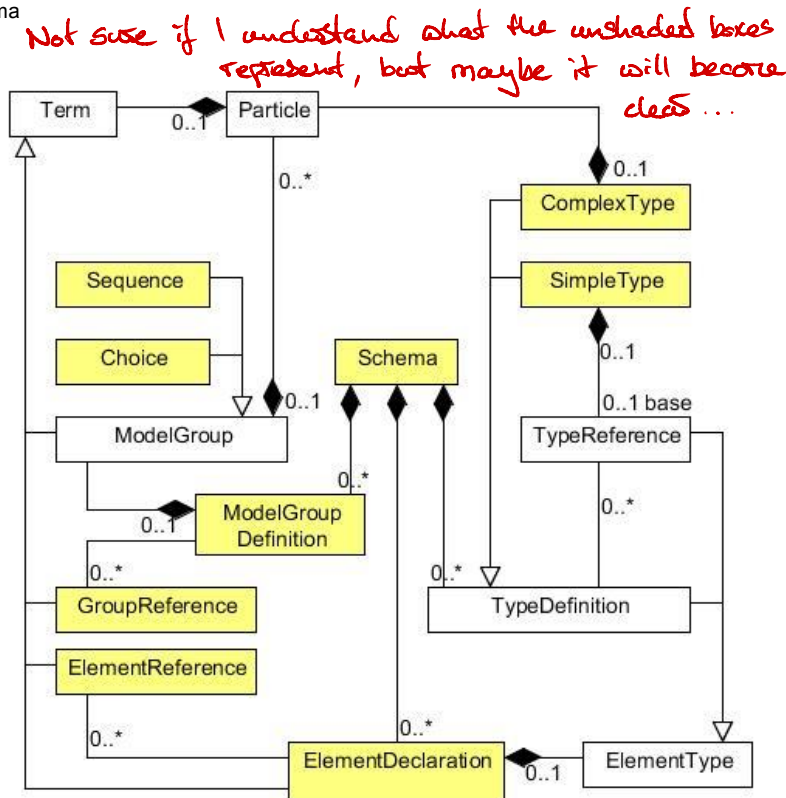


Figure 2 DFDL Schema UML diagram

The simple types are shown in Figure 3. The graph shows all the types defined by XML Schema version 1.0, and the subset of these types supported by DFDL are shown as shaded.

*All we seem to be back to describing the logical model maybe a lot of that discussion under the Augmented Infoset should be pushed to after this?*



## DFDL built-in types

*Significance of shading note?*

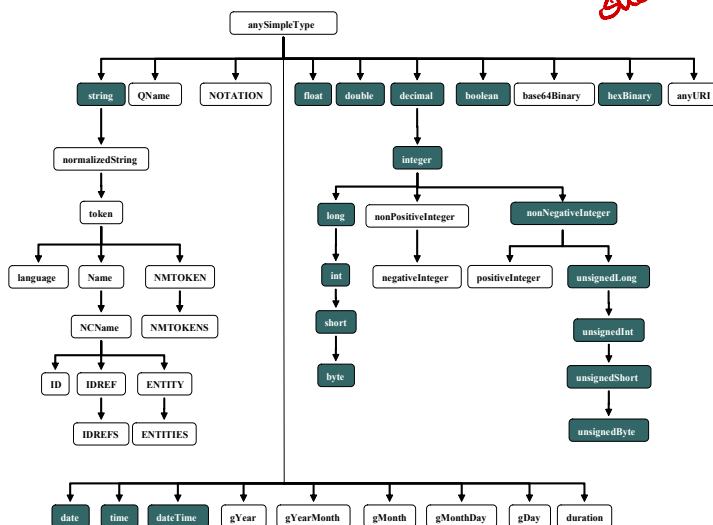


Figure 3 DFDL simple types

These types are defined as they are in XML Schema, with exceptions for:

- String – In DFDL a string can contain any character codes. None are reserved. (Including the character with character code U+0000, which is not permitted in XML documents.)

Each object defined by a class in the above UML is called a *DFDL Schema component*.

We express the DFDL Schema Model using a subset of the XML Schema Description Language (XSDL). XSDL provides a standardized schema language suitable for expressing the DFDL Schema Model.

A DFDL Schema is an XML schema containing only a restricted subset of the constructs available in full W3C XML Schema Description Language. Within this XML schema, special DFDL annotations are distributed that carry the information about the data's format or representation.

A DFDL Schema is a valid XML schema. However, the converse is not true since the DFDL Schema Model does not include many concepts that appear in XML schema.

### 5.1 DFDL Subset of XML Schema

The DFDL subset of XSDL is a general model for hierarchically-nested data. It avoids the XSDL features used to describe the peculiarities of XML as a syntactic textual representation of data, and features that are simply not needed by DFDL.

The following lists detail the similarities and differences between general XSDL and this subset.

DFDL Schemas consist of:

- Standard XSDL namespace management

- Standard XSDL import and management for multiple file schemas
- Local element declarations with dimensionality via maxOccurs and minOccurs.
- Global element declarations
- ComplexType definitions with empty or element-only content models.
- DFDL appinfo annotations describing the data format
- These simple types: string, float, double, decimal, integer, long, int, short, byte, nonNegativeInteger, unsignedLong, unsignedInt, unsignedShort, unsignedByte, boolean, date, time, dateTime, hexBinary
- These facets: minLength, maxLength, minInclusive, maxInclusive, minExclusive, maxExclusive, totalDigits, fractionDigits, enumeration, pattern (for xs:string type only)
- Fixed values
- Default values
- 'sequence' model groups (without minOccurs and maxOccurs or with both minOccurs="1" and maxOccurs="1")
- 'choice' model groups (without minOccurs and maxOccurs or with both minOccurs="1" and maxOccurs="1")
- Simple type derivations derived by restriction from the allowed built-in types
- Reusable Groups: named model group definitions can only contain one model group
- Element references with dimensionality via maxOccurs and minOccurs.
- Group references without dimensionality
- Nillable attribute is "true" (that is, nillable="true" in the element declaration.)
- Appinfo annotations for sources other than DFDL are permitted and ignored
- Unions; the memberTypes must be derived from the same simple type. DFDL annotations are not permitted on union members.<sup>5</sup>
- XML Entities
- The xs:schema "elementFormDefault" attribute
- The xs:element "form" attribute

Note: xs:nonNegativeInteger is treated as an unsigned xs:integer.

The following constructs from XML Schema are not used as part of the DFDL Schema Model of DFDL v1.0 schemas; however, they are all reserved<sup>6</sup> for future use since the data model may be extended to use them in future versions of DFDL:

- Attribute declarations (local or global)
- Attribute references
- Attribute group definitions
- Complex type derivations where the base type is not xs:anyType.
- Complex types having mixed content models or simple content models
- List simple types
- Union simple types where the member types are not derived from the same simple type.
- These atomic simple types: normalizedString, token, Name, NCName, QName, language, positiveInteger, nonPositiveInteger, negativeInteger, gYear, gYearMonth,

<sup>5</sup> The purpose of unions is to allow multiple constraints via facets such as multiple independent range restrictions on numbers. This enhances the ability to do rich validation of data.

<sup>6</sup> By reserved we mean that conforming DFDL v1.0 implementations MAY NOT assign semantics to them.

gMonth, gMonthDay, gDay, ID, IDREF, IDREFS, ENTITIES, ENTITY, NMTOKEN, NMTOKENS, NOTATION, anyURI, base64Binary

- maxOccurs and minOccurs on model groups (except if both are '1')
- minOccurs = '0' on branches of xs:choice model groups
- Identity Constraints
- Substitution Groups
- xs:all groups
- xs:any element wildcards
- Redefine - This version of DFDL does not support xs:redefine. DFDL schemas must not contain xs:redefine directly or indirectly in schemas they import or include.
- whitespace facet
- Recursively-defined types and elements (defined by way of type, group, or element references)

## 5.2 XSD Facets, min/maxOccurs, default, and fixed

XSD element declarations and references can carry several properties that express constraints on the described data. These constraints are mainly used for validation. These properties include:

- the facets
- minOccurs, maxOccurs
- default
- fixed

The facets and the types they are applicable to are:

- minLength maxLength (for types xs:string, and xs:hexBinary)
- pattern
- enumeration (all types except xs:boolean)
- maxInclusive, maxExclusive, minExclusive, minInclusive (for types xs:float, xs:double, xs:date, xs:time, xs:dateTime, xs:decimal and all integer types descending from xs:decimal in Figure 3)
- totalDigits (for type xs:decimal and all integer types descending from xs:decimal in Figure 3)
- fractionDigits (for type xs:decimal)

The facets (but not maxOccurs nor minOccurs) are also checked by the dfdl:checkConstraints DFDL expression language function.

The following sections describe these in more detail.

### 5.2.1 MinOccurs, MaxOccurs

The XSDL minOccurs property is used:

- To determine if an element declaration or reference is an array, an optional element, or neither.
- For some values of the property dfdl:occursCountKind, to determine the required minimum number of occurrences of an array both when parsing and unparsing.
- If validating, to determine the minimum valid number of occurrences of an array both when parsing and unparsing.

The XSDL maxOccurs property is used:

- To determine if an element declaration or reference is an array, an optional element, or neither.

**Deleted:** (for type xs:string and all types descending from xs:string in Figure 3 DFDL simple types)

Dropping specific DFDL behavior details into the discussion of XML schema tags is hard to follow. At this point I don't really know enough about these tags to follow.

maybe the fix is to mention the DFDL tags and their relationship to the current context but leave the specification of behaviour to the section on those tags

- When `dfdl:occursCountKind` is "fixed", then the XSDL `maxOccurs` value is the fixed number of occurrences of the array element, which is then called a Fixed Array Element. It is a Schema Definition Error if XSDL `minOccurs` is not equal to XSDL `maxOccurs`.
- When `dfdl:occursCountKind` is "implicit" then XSDL `maxOccurs` value (if not unbounded) is used to determine the maximum number of element occurrences both when parsing or unparsing.
- If validating, to determine the maximum valid number of occurrences of an array both when parsing and unparsing.

For some values of `dfdl:occursCountKind` such as 'implicit', it is a processing error when an array is found to have a number of occurrences not conforming to XSDL `minOccurs` in the absence of a default value specification. For other values of `dfdl:occursCountKind` such as 'parsed', it is only a validation error if an array is found to have fewer than XSDL `minOccurs` occurrences. See Section 16, Properties for Array Elements and Optional Elements, for more details.

#### 5.2.2 MinLength, MaxLength

These facets are used:

- When `dfdl:lengthKind` is "implicit" and type is `xs:string` or `xs:hexBinary`. In that case the length is given by the value of the XSD `maxLength` facet. In this case the XSD `minLength` facet is required to be equal to the XSD `maxLength` facet (Schema Definition Error otherwise).
- For validation of variable length string elements.

#### 5.2.3 MaxInclusive, MaxExclusive, MinExclusive, MinInclusive, TotalDigits, FractionDigits

- Used for validation only

The format of numbers is not derived from these facets. Rather `dfdl` properties are used to specify the format.

#### 5.2.4 Pattern

- Allowed only on elements of type `xs:string` or types derived from it in Figure 3 DFDL simple types.
- Used for validation only

It is important to avoid confusion of the pattern facet with other uses of regular expressions that are needed in DFDL (for example, to determine the length of an element by regular expression matching).

Note: in XSD, pattern is about the lexical representation of the data, and since all is text there, everything has a lexical representation. In DFDL only strings are guaranteed to have a lexical and logical value that is identical.

#### 5.2.5 Enumeration

Enumerations are used to provide a list of valid values in XSD.

- Used for validation only

Note: in DFDL we do not use XSD enumeration as a means to define symbolic constants. These are captured using `dfdl:defineVariable` constructs so they can be referenced from expressions.

#### 5.2.6 Default

The 'default' property is used:

- To provide the logical value of a required element while parsing when the element is missing. See 9.4 Element Defaults.
- To provide the logical value of a required element when unparsing when element is missing. See 9.4 Element Defaults.

Note that the 'fixed' and 'default' properties are mutually exclusive on an element declaration.

### 5.2.7 Fixed

The 'fixed' property is used in the same ways as the 'default' property but in addition:

- To constrain the logical value of an element when validating.

Note that the 'fixed' and 'default' properties are mutually exclusive on an element declaration.

### 5.3 Compatibility with Other Annotation Language Schemas

A DFDL Schema only applies annotations on a subset of the XML Schema constructs. However, a DFDL schema may be annotated not only by DFDL annotations, but by other XML annotation languages.

The XML schemas of those other annotation languages MAY use any constructs of XML Schema, including those prohibited by DFDL (such as attribute declarations.)

A DFDL implementation MUST ignore any schema file included or imported by a DFDL schema if the top level xs:schema element of that schema does not have an XML namespace binding for the DFDL namespace.

*I feel like i'm missing something. Does this make it clear what my DFDL implementation behavior should be if I feed it one of these hybrid XML / other annotation schemas?*

## 6. DFDL Syntax Basics

Using DFDL, a data format is described by placing special annotations at various positions within an XML schema. This XML schema conveys the basic structure of the data format, while the annotations fill in the detail. Annotations are used to describe aspects such as the file encoding and byte ordering, as well as declaring variables for reference elsewhere, and specifying properties that govern the capabilities of the DFDL processor. A DFDL processor requires these annotations, along with the structural information of the enclosing XML schema, to make sense of the physical data model.

### 6.1 Namespaces

The `xs:appinfo` source URI `http://www.ogf.org/dfdl/` is used to distinguish DFDL annotations from other annotations.

The element and attribute names in the DFDL syntax are in a namespace defined by the URI `http://www.ogf.org/dfdl/dfdl-1.0/`. All symbols in this namespace are reserved. DFDL implementations may not provide extensions to the DFDL standard using names in this namespace. Within this specification, the namespace prefix for DFDL is "dfdl" referring to the namespace `http://www.ogf.org/dfdl/dfdl-1.0/`.

Attributes on DFDL annotations that are not in the DFDL namespace or in no namespace are ignored.

A DFDL Schema document contains XML schema annotation elements that define and assign names to parts of the format specification. These names are defined using the target namespace of the schema document where they reside, and are referenced using QNames in the usual manner. A DFDL schema document can include or import another schema document, and namespaces work in the usual manner for XML schema documents. The schema is the schema including all additional schemas referenced through import and include. Generally, in this specification, when we refer to the DFDL Schema we mean the schema. When we refer to a specific document we will use the term DFDL Schema document.

### 6.2 The DFDL Annotation Elements

DFDL annotations must be positioned specifically where DFDL annotations are allowed within an XML schema document. These positions are known as *annotation points*. When an annotation is positioned at an annotation point, it binds some additional information to the schema component that encloses it. The description of a data format is achieved by correctly placing annotations on the structural components of the schema.

DFDL specifies a collection of annotations for different purposes. They are organized into three different annotation types: Format Annotations, Statement Annotations, and Defining Annotations.

At any single annotation point of the schema there can be only one format annotation, but there can be several statement annotations although there are rules about which of those are allowed to co-exist as well which will be described in sections about those specific annotation types.

Annotation Type	Annotation Element	Description
Format Annotation	<code>xs:choice</code>	Defines the physical data format properties of an <code>xs:choice</code> group. See section 7.1.
	<code>xs:element</code>	Defines the physical data format properties of an <code>xs:element</code> and <code>xs:element</code> reference. See section 7.1.
	<code>xs:format</code>	Defines the physical data format properties for multiple DFDL schema constructs. Used on an <code>xs:schema</code> and as a child of a <code>dfdl:defineFormat</code> annotation. This includes aspects such as the encodings, separators, and many more. See section 7.1.
	<code>xs:group</code>	Defines the physical data format properties of an <code>xs:group</code> reference. See section 7.1.

	<i>xs:</i> property	Used in the syntax of format annotations. See section 7.1.2.2.
	<i>xs:</i> sequence	Defines the physical data format properties of an xs:sequence group. See section 7.1.
	<i>xs:</i> simpleType	Defines the physical data format properties of an xs:simpleType. See section 7.1.
	<i>dfdl:</i> escapeScheme	Defines the scheme by which quotation marks and escape characters can be specified. This is for use with delimited text formats. See section 7.6.
Statement Annotation	<i>xs:</i> assert	Defines a test to be used to ensure the data are well formed. Assert is used only when parsing data. See section 7.3
	<i>dfdl:</i> discriminator	Defines a test to be used when resolving a point of uncertainty such as choice branches or optional elements. A dfdl:discriminator is used only when parsing data to resolve the point of uncertainty to one of the alternatives. See section 7.4
	newVariableInstance	Creates a new instance of a variable. See section 7.8
	setVariable	Sets the value of a variable whose declaration is in scope See section 7.9
Defining Annotation	defineEscapeScheme	Defines a named, reusable escapeScheme See section 7.5
	defineFormat	Defines a reusable data format by collecting together other annotations and associating them with a name that can be referenced from elsewhere. See section 7.2
	defineVariable	Defines a variable that can be referenced elsewhere. This can be used to communicate a parameter from one part of processing to another part. See section 7.7

Table 1 - DFDL Annotation Elements

### 6.3 DFDL Properties

Properties on DFDL annotations may be one or more of the following types

- DFDL string literal:  
The property represents a sequence of literal bytes or characters which appear in the **data stream**. The value type is a restriction of the XSDL xs:token that further disallows the space character. DFDL entities must be used to express whitespace in a DFDL String Literal.
- DFDL expression  
The property is an **xs:string** the value of which is a DFDL subset XPath 2.0 expression that returns a value derived from other property values and/or from the DFDL infoset. Leading and trailing whitespace is trimmed for DFDL expressions.
- DFDL regular expression  
The property is an **xs:string** the value of which is a regular expression that can be used as a pattern to calculate the length of an element by applying that pattern to the sequence of literal bytes or characters which appear in the data stream. Note that leading and trailing whitespace is not trimmed and is part of the regular expression value.

*you get the idea I need to know these these live.*

*I have no idea what a variable is (actually its ok in this context)*

*Maybe my confusion here is I don't know what a Property is? Is it the value of a DFDL annotation?*

*examples would help?*

*Whoa! I thought these were properties in the DFDL Schema text? or I misunderstood?*

*are you using XSDL to describe text in the XMC Schema? This is feeling very tricky! Needs clarification?*

- **Enumeration**  
The property value is an XSDL xs:token the value of which is one of the allowed values listed in the property description.
- **Logical Value.**  
The property value is a string that describes a logical value. The type of the logical value is one of the XML schema simple types. The string must conform to the XML schema lexical representation for the type.
- **QName**  
The property value is an XML Qualified Name as specified in "Namespaces in XML" [\[XMLNamespaces\]](#)

Some properties accept a list or union of types

- **List of DFDL String Literals or Logical Values**  
The property value is a whitespace separated list of the specified type. When parsing, if more than one string literal in the list matches the portion of the data stream being evaluated then the longest matching value in the list must be used. When unparsing, the first value in the list must be used. String literals containing whitespace or string literals representing the empty string must use character class entities in their syntax.
- **Union of types and expressions.**  
The property value is a union of DFDL expression and exactly one of the other types. The expression must resolve to a value of the other type.
- **Union of types.**  
The property value is a union of two or more types. The type is dependent on the value of another property. For example dfdl:nilValue can be a List of DFDL String Literals or a List of Logical Values depending on dfdl:nilKind

### 6.3.1 DFDL String Literals

DFDL String Literals represent a sequence of literal bytes or characters which appear in the data stream. This presents the following challenges

- the literal characters in the data stream might not be in the same encoding as the DFDL schema
- it may be necessary to specify a literal character which is not valid in an XML document
- it may be necessary to specify one or more raw byte values
- A DFDL string literal can describe any of the following types of literal data in any combination:
  - a single literal character in any encoding
  - a string of literal characters in any encoding
  - a bi-directional character string
  - one or more characters from a set of related characters ( e.g. end-of-line characters)
  - a literal byte value

A DFDL string literal is therefore able to describe any arbitrary sequence of bytes and characters.

**Empty String:** The special DFDL entity %ES; is provided for describing an empty string or an empty byte sequence. The %ES; entity is the only way to do this. A DFDL string literal with value "" (the empty string) is usually invalid. There are a few properties that explicitly allow an empty DFDL String Literal, and these properties assign a property-specific meaning to the empty string value.

**Whitespace:** When whitespace must be used as part of a property value, the DFDL string literal must use entities (such as %WSP;) to represent the whitespace. (This allows a property to represent lists of DFDL string literals by using literal spaces to separate list elements.)

extra  
included?



### 6.3.1.1 Character strings in DFDL String Literals

A literal string in a DFDL Schema is written in the character set encoding specified by the XML directive that begins all XML documents:

```
<?xml version="1.0" encoding="UTF-8" ?>
```

In this example, the DFDL schema is written in UTF-8, so any literal strings contained in it, and particularly string literals found in its representation property bindings in the format annotations, are expressed in UTF-8.

However, these strings are being used to describe features of text data that are commonly in other character set encodings. For example, we may have EBCDIC data that is comma separated. A comma in EBCDIC has a single-byte code unit of 0x6B in the data, the numeric value of which does not correspond to the Unicode character code for comma which is U+002C. However, when we indicate that an item is "," (comma) separated and we specify this using a string literal along with specifying the 'encoding' property to be 'ebcdic-cp-us' then this means that the data are separated by EBCDIC commas regardless of what character set encoding is used to write the DFDL Schema.

```
<?xml version="1.0" encoding="UTF-8">
<xs:schema ... >
  ...
  <dfdl:format encoding="ebcdic-cp-us" separator=","/>
  ...
</xs:schema>
```

When a DFDL processor uses the separator expressed in this manner, the string literal "," is *translated* into the character set encoding of the data it is separating as specified by the encoding representation property. Hence, in this case we would be searching the data for a character with codepoint 0x6B (the EBCDIC comma), not a UTF-8 or Unicode (0x2C) comma which is what exists in the DFDL schema document file.

### 6.3.1.2 DFDL Character Entities, Character Class Entities, and Byte Values in String Literals

DFDL character entities specify a single Unicode character and provide a convenient way to specify code points that appear in the data stream but would be difficult to specify in XML strings. For example, common non-printable characters or code points, such as 0x00, that are not valid in XML documents. DFDL entities are based on XML entities, which can also be used in a DFDL schema.

The following grammar gives the syntax of DFDL String Literals generally, including the various kinds of entities.

DfdlStringLiteral	::= (DfdlStringLiteralPart)+   DfdlESEntity
DfdlStringLiteralPart	::= LiteralString   DfdlCharEntity   DfdlCharClass   ByteValue
LiteralString	::= A string of literal characters
DfdlCharEntity	::= DfdlEntity   DecimalCodePoint   HexadecimalCodePoint
DfdlCharClass	::= '%' DfdlCharClassName ';' ;
ByteValue	::= '%#r' [0-9a-fA-F]{2} ';' ;
DfdlEntity	::= '%' DfdlEntityName ';' ;
DecimalCodePoint	::= '%#' [0-9]+ ';' ;
HexadecimalCodePoint	::= '%#x' [0-9a-fA-F]+ ';' ;
DfdlEntityName	::= 'NUL'   'SOH'   'STX'   'ETX'

So, based on this example a "property" is a value applied to a DFDL attribute which can correspond to things in the data.

**Deleted:** Character strings can include bidirectional data. ¶

OK I understand the principle but this is an example to make it concrete

		'EOT' 'ENQ' 'ACK' 'BEL'
		'BS' 'HT' 'LF' 'VT' 'FF'
		'CR' 'SO' 'SI' 'DLE'
		'DC1' 'DC2' 'DC3' 'DC4'
		'NAK' 'SYN' 'ETB' 'CAN'
		'EM' 'SUB' 'ESC' 'FS'
		'GS' 'RS' 'US' 'SP'
		'DEL' 'NBSP' 'NEL' 'LS'
DfdlCharClassName	::=	DfdlNLEntity   DfdlWSPEntity   DfdlWSPStarEntity   DfdlWSPPlusEntity
DfdlNLEntity	::=	'NL'
DfdlWSPEntity	::=	'WSP'
DfdlWSPStarEntity	::=	'WSP*'
DfdlWSPPlusEntity	::=	'WSP+'
DfdlESEntity	::=	'ES'

**Table 2 DFDL Character Entity, Character Class Entity, and Byte Value Entity Syntax**

Using %% inserts a single literal "%" into the string literal. This "%" is subject to character set encoding translation as is any other character.

A HexadecimalCodePoint provides a hexadecimal representation of the character's code point in ISO/IEC 10646.

A DecimalCodePoint provides a decimal representation of the character's code point in ISO/IEC 10646.

A DfdlEntityName is one of the mnemonics given in the following tables.

Mnemonic	Meaning	Unicode Character Code
NUL	null	U+0000
SOH	start of heading	U+0001
STX	start of text	U+0002
ETX	end of text	U+0003
EOT	end of transmission	U+0004
ENQ	enquiry	U+0005
ACK	acknowledge	U+0006
BEL	bell	U+0007
BS	backspace	U+0008
HT	horizontal tab	U+0009
LF	line feed	U+000A
VT	vertical tab	U+000B
FF	form feed	U+000C
CR	carriage return	U+000D
SO	shift out	U+000E

SI	shift in	U+000F
DLE	data link escape	U+0010
DC1	device control 1	U+0011
DC2	device control 2	U+0012
DC3	device control 3	U+0013
DC4	device control 4	U+0014
NAK	negative acknowledge	U+0015
SYN	synchronous idle	U+0016
ETB	end of transmission block	U+0017
CAN	cancel	U+0018
EM	end of medium	U+0019
SUB	substitute	U+001A
ESC	escape	U+001B
FS	file separator	U+001C
GS	group separator	U+001D
RS	record separator	U+001E
US	unit separator	U+001F
SP	space	U+0020
DEL	delete	U+007F
NBSP	no break space	U+00A0
NEL	Next line	U+0085
LS	Line separator	U+2028

**Table 3 DFDL Entities****6.3.1.3 DFDL Character Class Entities in DFDL String Literals**

The following DFDL character classes are provided to specify one or more characters from a set of related characters.

Mnemonic	Meaning	Unicode Character Code(s)
NL	Newline On parse any one of the single characters CR, LF, NEL or LS or the character combination CRLF. On unparse the value of the dfdl:outputNewLine property is output, which must specify one of the single characters %CR;, %LF;, %NEL;, or %LS; or the character combination %CR;%LF;.	U+000A LF U+000D CR U+000D U+000A CRLF U+0085 NEL U+2028 LS
WSP	Single whitespace On parse any whitespace character On unparse a space (U+0020) is output	U+0009-U+000D (Control characters) U+0020 SPACE U+0085 NEL

		U+00A0 NBSP U+1680 OGHAM SPACE MARK U+180E MONGOLIAN VOWEL SEPARATOR U+2000-U+200A (different sorts of spaces) U+2028 LSP U+2029 PSP U+202F NARROW NBSP U+205F MEDIUM MATHEMATICAL SPACE U+3000 IDEOGRAPHIC SPACE
WSP*	Optional Whitespaces On parse whitespace characters are ignored. On unparse nothing is output	Same as WSP
WSP+	Whitespaces On parse one or more whitespace characters are ignored. It is an processing error if no whitespace character is found. On unparse a space (U+0020) is output.	Same as WSP
ES	Empty String Used in whitespace separated lists when empty string is one of the values.	

Table 4 DFDL Character Class Entities

#### 6.3.1.4 DFDL Byte Value Entities in DFDL String Literals

DFDL byte value entities provide a way to specify a single byte as it appears in the data stream without any character set encoding translation. To specify a string of byte values, a sequence of two or more byte value entities must be used. The syntax is in Table 2 DFDL Character Entity, Character Class Entity, and Byte Value Entity Syntax above. Example:

```
%#rFF;
```

#### 6.3.2 DFDL Expressions

Some DFDL properties allow DFDL expressions (see Section 23 Expression language) to be used so that the property can be set dynamically at processing-time.

The general syntax of expressions is "{ expression }"

The rules for recognizing DFDL expressions are

- Discard any leading and trailing whitespace.
- Must start with a '{' in the first position and end with '}' in the last position.
- '{' in any position other than the first is treated as a literal.
- '}' in any position other than the last position is treated as a literal.

- '{' as the first characters are treated as the literal '{' and not as the start of a DFDL expression.

DFDL expressions reference other items in the infoset or augmented infoset using absolute or relative paths. Relative paths are evaluated when the component containing the expression is referenced not when it is declared. For example a global element may have a DFDL property which is an expression that contains a relative path to another element. The relative path is evaluated when the global element is referenced from an element reference.

DFDL expressions that are used to provide the value of DFDL properties in the `dfdl:format` annotation on the top level `xs:schema` declaration MAY NOT contain relative paths.

#### 6.3.3 DFDL Regular Expressions

The DFDL `lengthPattern` property expects a regular expression to be specified. The DFDL Regular Expression language is defined in the section 24 DFDL Regular Expressions.

#### 6.3.4 Enumerations in DFDL

Some DFDL properties accept an enumerated list of valid values. It is a Schema Definition Error if a value other than one of the enumerated values is specified. The case of the specified value must match the enumeration. An enumeration is of type string unless otherwise stated.

## 7. Syntax of DFDL Annotation Elements

This section describes the syntax of each of the DFDL annotation elements along with discussion of their basic meanings.

The DFDL annotation elements are listed in **Table 1 - DFDL Annotation Elements**

### 7.1 Component Format Annotations

A data format can be 'used' or put into effect for a part of the schema by use of the component format annotation elements.

There are specific annotations for each type of schema component that supports only the representation properties applicable to that component. The table below gives the specific annotation for each schema component.

Schema component	DFDL annotation
xs:choice	dfdl:choice
xs:element	dfdl:element
xs:element reference	dfdl:element
xs:group reference	dfdl:group
xs:schema	dfdl:format
xs:sequence	dfdl:sequence
xs:simpleType	dfdl:simpleType

**Table 5 DFDL Component Format Annotations**

In addition the dfdl:format annotation is used inside a dfdl:defineFormat annotation to define a named reusable set of representation properties that can be referenced from any component specific format annotation or from other named format definitions.

A dfdl:format annotation at the top level of a schema, that is as an annotation child element on the xs:schema, provides a set of default properties for the lexically enclosed schema document. See 8.1 Providing Defaults for DFDL properties.

Example of DFDL component format annotation:

```
<xs:schema ...>
...
<xs:element name="root">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:element ref="aBaseConfig"
        representation="text"
        encoding="UTF-8"/>
    </xs:appinfo>
  </xs:annotation>
</xs:element>
...
</xs:schema>
```

*These seems to be a lot going on here would it be worth walking through?*

*I don't remember what xs:appinfo does*

#### 7.1.1 The dfdl:ref Property

A named, reusable, dfdl:defineFormat definition is used by referring to its name from a format annotation using the 'ref' property. For example:

```
<dfdl:element ref="reusableDef" encoding="ebcdic-cp-us" />
```

The behavior of this dfdl:defineFormat definition is as if all representation properties defined by the named dfdl:defineFormat definition were instead written directly on this format annotation;

*Ah, of course! But I needed to know this earlier. It feels like the doc was originally written with logical ordering of concepts, but over time details have crept into sections earlier on before they can be understood by a first-time reader.*

*I'm lost. I don't understand how this section is different, or relates to the one above.*

*No idea what this means there is no dfdl:defineFormat in the example?*

however, these are superseded by any representation properties that are defined here such as the encoding property in the example above.

### 7.1.2 Property Binding Syntax

The format properties may be specified in one of three forms:

1. Attribute form
2. Element form
3. Short form

A DFDL property may be specified using any form with the following exceptions

- The dfdl:ref property may be specified in attribute or short form
- The dfdl:escapeSchemeRef property may be specified in attribute or short form
- The dfdl:hiddenGroupRef property may be specified in attribute or short form
- The dfdl:prefixLengthType property may be specified in attribute or short form
- Short form is not allowed on the xs:schema element.

It is a Schema Definition Error if the same property is specified in more than one form in the resolved set of annotations for an annotation point.

#### 7.1.2.1 Property Binding Syntax: Attribute Form

Within the format annotation elements are bindings for properties of the form:

```
Property="Value"
```

For example:

```
<xs:annotation>
  <xs:appinfo source="http://www.ogf.org/dfdl/">
    <dfdl:format encoding="utf-8" separator="%NL;"/>
  </xs:appinfo>
</xs:annotation>
```

The **Property** is the name of the property. The **Value** is an XML string literal corresponding to a value of the appropriate type.

#### 7.1.2.2 Property Binding Syntax: Element Form

The representation properties can sometimes have complex syntax, so an element form for representation property bindings is provided as element content within the format element content model. This is provided to ease syntactic expression difficulties. The element is called dfdl:property and it has one attribute 'name' which provides the name of the property.

For example:

```
<xs:annotation>
  <xs:appinfo source="http://www.ogf.org/dfdl/">
    <dfdl:format>
      <dfdl:property name='encoding'>utf-8</dfdl:property>
      <dfdl:property name='separator'>%NL;</dfdl:property>
    </dfdl:format>
  </xs:appinfo>
</xs:annotation>
```

Element form is mostly used for properties that themselves contain the quotation mark characters and escape characters so that they can be expressed without concerns about confusion with the XSDL syntax use of these same characters. CDATA encapsulation can be used so as to allow malformed XML and mismatched quotes to be easily used as representation property values:

```
<dfdl:property name='initiator'><CDATA[<!-- ]]></dfdl:property>
```

I'm coming back after a few days but I don't think you told me what a property is.

??? Two examples above neither makes sense to me

what is an annotation point

what is meant by annotation

I see finally!  
A Property is a specific DFDL attribute that tells the DFDL implementation something about the data format locally or globally

circulas?

### 7.1.2.3 Property Binding Syntax: Short Form

To save textual clutter, short-form syntax for format annotations is also allowed on `xs:element`, `xs:sequence`, `xs:choice`, `xs:group` (for group references only), and `xs:simpleType` schema elements. (The `xs:schema` element cannot carry short-form annotations). Attributes which are in the namespace '<http://www.ogf.org/dfdl/dfdl-1.0/>' and whose local name matches one of the DFDL representation properties are assumed to be equivalent to specific DFDL attribute form annotations.

For example the two forms below are equivalent in that they describe the same data format. The first is the short form of the second:

```
<xs:element name="elem1">
  <xs:complexType>
    <xs:sequence dfdl:separator="%HT;" >
      ...
    </xs:sequence>
  </xs:complexType>
</xs:element>

<xs:element name="elem2">
  <xs:complexType>
    <xs:sequence>
      <xs:annotation><xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:sequence separator="%HT;" />
      </xs:appinfo></xs:annotation>
      ...
    </xs:sequence>
  </xs:complexType>
</xs:element>
```

*Very helpful example*

Another example:

```
<xs:sequence dfdl:separator=",">
  <xs:element name="elem1" type="xs:int" maxOccurs="unbounded"
    dfdl:representation="text"
    dfdl:textNumberRep="standard"
    dfdl:initiator="["
    dfdl:terminator="]"/>

  <xs:element name="elem2" type="xs:int" maxOccurs="unbounded">
    <xs:annotation><xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:element representation="text"
        textNumberRep="standard"
        initiator="["
        terminator="]"/>
    </xs:appinfo></xs:annotation>
  </xs:element>
</xs:sequence>
```

Because short form syntax is not allowed on the `xs:schema` element, an attribute form `dfdl:format` annotation must be used instead. *Does this refer to the example? Orphaned text?*

### 7.1.3 Empty String as a Representation Property Value

DFDL provides no mechanism to un-set a property. Setting a representation property's value to the empty string doesn't remove the value for that property, but sets it to the empty string value. This may not be appropriate as a value for certain properties.

For example, in delimited text data formats, it is sensible for the separator to be defined to be the empty string. This turns off use of separator delimiters. For many other string-valued properties, it is a Schema Definition Error to assign them the empty string value. For example, the character set encoding property (`dfdl:encoding`) cannot be set to the empty string.



## 7.2 dfdl:defineFormat - Reusable Data Format Definitions

One or more dfdl:defineFormat annotation elements can appear within the annotation children of the xs:schema element. DFDL defining annotation elements may only appear as annotation children of the xs:schema element.

The order of their appearance does not matter, nor does their position relative to other non-annotation children of the xs:schema.

Each dfdl:defineFormat has a required name attribute.

The construct creates a named data format definition. The value of the name attribute is of XML type NCName. The format name will become a member of the schema's target namespace. These names must be unique within the namespace.

If multiple format definitions have the same 'name' attribute, in the same namespace, then it is a Schema Definition Error.

Here is an example of a format definition:

```
<xs:schema ...>
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:defineFormat name="myConfig" >
        <dfdl:format representation="text"
          ref="textSpecialFormat1" />
      </dfdl:defineFormat>
    </xs:appinfo>
  </xs:annotation>
  ...
</xs:schema>
```

A dfdl:defineFormat serves only to supply a named definition for a format for reuse from other places. It does not cause any use of the representation properties it contains to describe any actual data.

### 7.2.1 Inheritance for dfdl:defineFormat

A dfdl:defineFormat declaration can inherit from another named format definition by use of the dfdl:ref property of the dfdl:format annotation. This allows a single-inheritance hierarchy that reuses definitions. When one definition extends another in this way, any property definitions contained in its direct elements override those in any inherited definition.

Conceptually, the dfdl:ref inheritance chains can be *flattened* and removed by copying all inherited property bindings and then superseding those for which there is a local binding. Throughout this document we will assume inheritance is fully flattened. That is, all dfdl:ref inheritance is first removed by flattening before any other examination of properties occurs.

It is a Schema Definition Error if use of the dfdl:ref property results in a circular path.

### 7.2.2 Using/Referencing a Named Format Definition

See Section 7.1.1 The dfdl:ref Property.

## 7.3 The dfdl:assert Statement Annotation Element

The dfdl:assert statement annotation element is used to assert truths about a DFDL model that are used when parsing to ensure that the data are well-formed. They are not used when unparsing. These checks are separate from validation checking and are performed even when validation is off. This distinction is needed to ensure that switching validation off does not affect parsing.

Examples of dfdl:assert elements are below:

```
<dfdl:assert message="Value is not zero." test="{ ../x eq 0}" />

<dfdl:assert message="Precondition violation." >
  { ../x le 0 and ../y ne "-->" and ../y ne "<!--" }
</dfdl:assert>
```

```
<dfdl:assert message="Postcondition violation." testKind='expression'>
  {../x ne ""}
</dfdl:assert>
```

### 7.3.1 Properties for dfdl:assert

A dfdl:assert annotation contains a test expression or a test pattern. The dfdl:assert is said to be successful if the test expression evaluates to true or the test pattern returns a non-zero length match, and unsuccessful if the test expression evaluates to false or the test pattern returns a zero length match. An unsuccessful dfdl:assert causes either a processing error or a recoverable error to be issued, as specified by the failureType property of the dfdl:assert.

The testKind property specifies whether an expression or pattern is used by the dfdl:assert. The expression or pattern can be expressed as an attribute or as a value.

```
<dfdl:assert test="{test expression}" />

<dfdl:assert>
  {test expression}
</dfdl:assert>
```

It is a Schema Definition Error if a property is specified in more than one form.

It is a Schema Definition Error if both a test expression and a test pattern are specified.

A dfdl:assert can appear as an annotation on:

- an xs:element declaration (local or global)
- an xs:element reference
- an xs:group reference
- an xs:sequence
- an xs:choice
- an xs:simpleType definition (local or global)

If the resolved set of statement annotations for a schema component contains multiple dfdl:assert statements, then those with testKind 'pattern' are executed before those with testKind 'expression' (the default). However, within each group the order of execution among them is not specified.

If one of the resolved set of asserts for a schema component is unsuccessful, and the failureType of the assert is 'processingError', then no further asserts in the set are executed.

Property Name	Description
testKind	Enum (optional) Valid values are 'expression', 'pattern' Default value is 'expression' Specifies whether a DFDL expression or DFDL regular expression is used in the dfdl:assert. Annotation: dfdl:assert
test	DFDL Expression Applies when testKind is 'expression' A DFDL expression that evaluates to true or false. If the expression evaluates to true then parsing continues. If the expression evaluates to false then a processing error is raised. Any element referred to by the expression must have already been processed or must be a descendent of this element.

(pattern)

	<p>If a processing error occurs during the evaluation of the test expression then the <code>dfdl:assert</code> also fails.</p> <p>It is a Schema Definition Error if <code>testKind</code> is 'expression' or not specified, and an expression is not supplied by either the value of the <code>dfdl:assert</code> element or the value of the test attribute.</p> <p>Annotation: <code>dfdl:assert</code></p>
testPattern	<p>DFDL Regular Expression</p> <p>Applies when <code>testKind</code> is 'pattern'</p> <p>A DFDL regular expression that is applied against the data stream starting at the data position corresponding to the beginning of the representation. Consequently the framing (including any initiator) is visible to the pattern at the start of the component on which the <code>dfdl:assert</code> is positioned.</p> <p>If the pattern matching of the regular expression reads data that cannot be decoded into characters of the current encoding, then the behavior is controlled by the <code>dfdl:encodingErrorPolicy</code> property. See Section 11.2.1 Property <code>dfdl:encodingErrorPolicy</code> for details.</p> <p>If the length of the match is zero then the <code>dfdl:assert</code> evaluates to false and a processing error is raised.</p> <p>If the length of the match is non-zero then the <code>dfdl:assert</code> evaluates to true.</p> <p>If a processing error occurs during the evaluation of the test regular expression then the <code>dfdl:assert</code> also fails.</p> <p>It is a Schema Definition Error if <code>testKind</code> is 'pattern', and a pattern is not supplied by either the value of the <code>dfdl:assert</code> element or the value of the <code>testPattern</code> property.</p> <p>It is a Schema Definition Error if there is no value for the <code>dfdl:encoding</code> property in scope.</p> <p>It is a Schema Definition Error if <code>dfdl:leadingSkip</code> is other than 0.</p> <p>It is a Schema Definition Error if the <code>dfdl:alignment</code> is not 1 or 'implicit'</p> <p>Annotation: <code>dfdl:assert</code></p>
message	<p>String or DFDL Expression</p> <p>Defines text to be used as a diagnostic code or for use in an error message, when the assert is unsuccessful.</p> <p>The DFDL Expression must return type <code>xs:string</code>. Any element referred to by the message expression must have already been processed or must be a descendent of this element. There is special treatment for errors that occur while evaluating the message expression. See below for details.</p> <p>Annotation: <code>dfdl:assert</code></p>
failureType	<p>Enum (optional)</p> <p>Valid values are 'processingError', 'recoverableError'.</p> <p>Default value is 'processingError'.</p> <p>Specifies the type of failure that occurs when the <code>dfdl:assert</code> is unsuccessful.</p> <p>When 'processingError', a processing error is raised.</p> <p>When 'recoverableError', a recoverable error is raised.</p> <p>If an error occurs while evaluating the test expression, a processing error occurs, not a recoverable error.</p> <p>Recoverable errors do not cause backtracking like processing errors.</p> <p>Annotation: <code>dfdl:assert</code></p>

**Table 6 dfdl:assert properties**

Example of a dfdl:assert with a message expression:

```
<dfdl:assert message="{ fn:concat('unknown case ', ../data1) }">
{ if (...pred1...) then ...expr1...
  else if (...pred2...) then ...expr2...
  else fn:false()
}
</dfdl:assert>
```

The message specified by the message property is issued only if the dfdl:assert is unsuccessful, that is, the test expression evaluates to false or the test pattern returns a zero-length match. If so, and the message property is an expression, the message expression is evaluated at that time.

If a processing error or Schema Definition Error occurs while evaluating the message expression, a recoverable error is issued to record this error (containing implementation-dependent content), then processing of the assert continues as if there was no problem and in a manner consistent with the failureType property, but using an implementation-dependent substitute message.

*assertion? - what's a statement?*

### 7.3.2 Controlling the Timing of Statement Evaluation

Schema authors can insert xs:sequence constructs to control the timing of evaluation of statements more precisely. For example:

```
<xs:sequence dfdl:separator=",">
...
<xs:element ref="a" .../>
<xs:sequence>
  <xs:sequence>
    <xs:annotation><xs:appinfo source="http://www.ogf.org/dfdl/" >
      <dfdl:assert test="{test expression}" />
    </xs:appinfo></xs:annotation>
  </xs:sequence>
  <xs:element ref="b" .../>
</xs:sequence>
...
</xs:sequence>
```

In the above, the assert test expression is evaluated after parsing element 'a', and before parsing element 'b'. The use of two nested interior sequences surrounding element 'b' in this manner insures that the outermost sequence's separator usage is not disrupted.

### 7.4 The dfdl:discriminator Statement Annotation Element

DFDL discriminators are used during parsing to resolve points of uncertainty that cannot be resolved by speculative parsing. Discriminators are not used during unparsing. They can also be used to force a resolution earlier during the parsing of a group so that subsequent parsing errors are treated as processing errors of a known component rather than a failure to find a component.

A discriminator determines the existence or non-existence of a component. If the discriminator is successful then the component is known to exist and any subsequent errors will not cause backtracking at points of uncertainty. If a discriminator is unsuccessful then the component is known not to exist and backtracking occurs immediately.

If the complex type of an element contains a sequence group as its content model then if the sequence group is known not to exist, then the element is known not to exist.

Examples of dfdl:discriminator annotation are below :

```
<dfdl:discriminator>
{ ../recType eq 0 }
</dfdl:discriminator>
```

*have you defined 'component'? This is a value from the dataset?*

*I don't think I understand the parsing algo enough to read this section*

```
<dfdl:discriminator test="{ ../recType eq 0}" />
```

When the discriminator's expression evaluates to "false", then it causes a processing error, and the discriminator is said to fail.

#### 7.4.1 Properties for dfdl:discriminator

A DFDL discriminator contains a test expression that is an expression that evaluates to true or false. The discriminator is said to be successful if the test evaluates to true and unsuccessful (or fails) if the test evaluates to false.

The testKind property specifies whether an expression or pattern is used by the dfdl:discriminator. The expression or pattern can be expressed as an attribute or as a value.

```
<dfdl:discriminator test="{test expression}" />

<dfdl:discriminator>
  { test expression }
</dfdl:discriminator>
```

It is a Schema Definition Error if a property is specified in more than one form.

It is a Schema Definition Error if both a test expression and a test pattern are specified.

A dfdl:discriminator can be an annotation on

- an xs:element declaration (local or global)
- an xs:element reference
- an xs:group reference
- an xs:sequence
- an xs:choice
- an xs:simpleType definition (local or global)

The resolved set of statement annotations for a schema component can contain only a single dfdl:discriminator or one or more dfdl:assert annotations, but not both. To clarify: dfdl:assert annotations and dfdl:discriminator annotations are exclusive of each other. It is a Schema Definition Error otherwise.

Property Name	Description
testKind	Enum Valid values are 'expression', 'pattern' Default value is 'expression' Specifies whether a DFDL expression or DFDL regular expression is used in the dfdl:discriminator. Annotation: dfdl:discriminator
test	DFDL Expression Applies when testKind is 'expression' A DFDL expression that evaluates to true or false. If the expression evaluates to true then the discriminator succeeds and parsing continues. If the expression evaluates to false then the discriminator fails and a processing error is raised. If a processing error occurs during the evaluation of the test expression then the discriminator also fails. Any element referred to by the expression must have already been processed or is a descendent of this element. The expression must have been evaluated by the time this element and its descendents have been processed or when a processing error occurs when processing this element or its descendents.

	<p>It is a Schema Definition Error if testKind is 'expression' or not specified, and an expression is not supplied by either the value of the dfdl:discriminator element or the value of the test attribute.</p> <p>Annotation: dfdl:discriminator</p>
testPattern	<p>DFDL Regular Expression</p> <p>Applies when testKind is 'pattern'</p> <p>A DFDL regular expression that is applied against the data stream starting at the data position corresponding to the beginning of the representation. Consequently the framing (including any initiator) is visible to the pattern at the start of the component on which the dfdl:discriminator is positioned.</p> <p>If the pattern matching of the regular expression reads data that cannot be decoded into characters of the current encoding, then the behavior is controlled by the dfdl:encodingErrorPolicy property. See Section 11.2.1 Property dfdl:encodingErrorPolicy for details.</p> <p>If the length of the match is zero then the dfdl:discriminator evaluates to false and a processing error is raised.</p> <p>If the length of the match is non-zero then the dfdl:discriminator evaluates to true.</p> <p>It is a Schema Definition Error if testKind is 'pattern', and a pattern is not supplied by either the value of the dfdl:discriminator element or the value of the testPattern property.</p> <p>It is a Schema Definition Error if there is no value for the dfdl:encoding property in scope.</p> <p>It is a Schema Definition Error if dfdl:leadingSkip is other than 0.</p> <p>It is a Schema Definition Error if the dfdl:alignment is not 1 or 'implicit'</p> <p>Annotation: dfdl:discriminator</p>
message	<p>String or DFDL Expression</p> <p>Defines text to be used as a diagnostic code or for use in an error message, when the discriminator is unsuccessful.</p> <p>The DFDL Expression must return type xs:string. Any element referred to by the message expression must have already been processed or must be a descendent of this element. There is special treatment for errors that occur while evaluating the message expression. See below for details.</p> <p>Annotation: dfdl:discriminator</p>

**Table 7 dfdl:discriminator properties**

The message specified by the message property is issued only if the discriminator is unsuccessful, that is, the test expression evaluates to false or the test pattern returns a zero-length match. If so, and the message property is an expression, the message expression is evaluated at that time.

If a processing error or Schema Definition Error occurs while evaluating the message expression, a recoverable error is issued to record this error (containing implementation-dependent content), then processing of the discriminator continues as if there was no problem, but in the case of failure using an implementation-dependent substitute message.

Examples of dfdl:discriminator annotations:

```
<xs:sequence>
  <xs:choice>
    <xs:element name='branchSimple' >
      <xs:annotation>
        <xs:appinfo source="http://www.ogf.org/dfdl/">
```

```

        <dfdl:discriminator test='{. eq "a"}' />
      </xs:appinfo>
    </xs:annotation>
  </xs:element>

  <xs:element name='branchComplex' >
    <xs:annotation>
      <xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:discriminator test='{./identifier eq "b"}' />
      </xs:appinfo>
    </xs:annotation>
    <xs:complexType >
      <xs:sequence>
        <xs:element name='identifier' />
        ...
      </xs:sequence>
    </xs:complexType>
  </xs:element>

  <xs:element name='branchNestedComplex' >
    <xs:annotation>
      <xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:discriminator test='{./Header/identifier eq "c"}' />
      </xs:appinfo>
    </xs:annotation>
    <xs:complexType >
      <xs:sequence>
        <xs:element name='Header' />
        <xs:complexType >
          <xs:sequence>
            <xs:element name='identifier' />
            ...
          </xs:sequence>
        </xs:complexType>
      </xs:element>
    </xs:sequence>
  </xs:complexType>
</xs:choice>
</xs:sequence>

```

### 7.5 The dfdl:defineEscapeScheme Defining Annotation Element

One or more dfdl:defineEscapeScheme annotation elements can appear within the annotation children of the xs:schema. The dfdl:defineEscapeScheme elements may only appear as annotation children of the xs:schema.

The order of their appearance does not matter, nor does their position relative to other annotation or non-annotation children of the xs:schema.

Each dfdl:defineEscapeScheme has a required name attribute and a required dfdl:escapeScheme child element.

The construct creates a named escape scheme definition. The value of the name attribute is of XML type NCName. The name will become a member of the schema's target namespace. These names must be unique within the namespace among escape schemes.

If multiple dfdl:defineEscapeScheme definitions have the same 'name' attribute, in the same namespace, then it is a Schema Definition Error.

Each dfdl:defineEscapeScheme annotation element contains a dfdl:escapeScheme annotation element as detailed below.

Here is an example of an escapeScheme definition:

```
<xs:schema ...>
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:defineEscapeScheme name="myEscapeScheme">
        ...
        <dfdl:escapeScheme escapeKind="escapeCharacter"
          escapeCharacter="/" />
        ...
      </dfdl:defineEscapeScheme>
    </xs:appinfo>
  </xs:annotation>
  ...
</xs:schema>
```

A `dfdl:defineEscapeScheme` serves only to supply a named definition for a `dfdl:escapeScheme` for reuse from other places. It does not cause any use of the representation properties it contains to describe any actual data.

### 7.5.1 Using/Referencing a Named `escapeScheme` Definition

A named, reusable, escape scheme is used by referring to its name from a `dfdl:escapeSchemeRef` property on an element. For example:

```
<xs:element name="foo" type="xs:string" >
  <xs:annotation><xs:appinfo source="http://www.ogf.org/dfdl/">
    <dfdl:element representation="text"
      escapeSchemeRef="myEscapeScheme"/>
  </xs:appinfo></xs:annotation>
</xs:element>
```

### 7.6 The `dfdl:escapeScheme` Annotation Element

The `dfdl:escapeScheme` annotation is used within a `dfdl:defineEscapeScheme` annotation to group the properties of an escape scheme and allows a common set of properties to be defined that can be reused.

An escape scheme defines the properties that describe the text escaping rules in force when data such as text delimiters are present in the data. There are two variants on such schemes,

- The use of a single escape character to cause the next character to be interpreted literally. The escape character itself is escaped by the escape escape character.
- The use of a pair of escape strings to cause the enclosed group of characters to be interpreted literally. The ending escape string is escaped by the escape escape character.

On parsing, the escape scheme is applied after pad characters are trimmed and on unparsing before pad characters are added.

DFDL does not perform any substitutions for ampersand notations like `&lt;`.

The syntax of `dfdl:escapeScheme` is defined in Section 13.2.1. [Table 26 Properties Common to All Simple Types with Text Representation](#)

[The `dfdl:escapeScheme` Properties.](#)

### 7.7 The `dfdl:defineVariable` Annotation Element

Variables provide a means for communication within a set of DFDL schema. They are defined as top-level elements in a schema and therefore have global scope.

A new variable is introduced using `dfdl:defineVariable`:

```
<dfdl:defineVariable
  name = NCName
  type? = QName
  defaultValue? = logical value or dfdl expression
  external? = 'false' | 'true' >
  <!-- Contains: logical value or dfdl expression (default value) -->
```

Deleted: Table 26 Properties Common to All Simple Types with Text Representation¶  
The `dfdl:escapeScheme` Properties

really? how do schema interpreters

Explain why variables are needed.  
Before getting into details of syntax - explain the model?



```
</dfdl:defineVariable>
```

The name of a newly defined variable is placed into the target namespace of the schema containing the annotation. Variable names are distinct from format and escape scheme names and so cannot conflict with them. A variable can have any type from the DFDL subset of XML schema simple types. If no type is specified, the type is xs:string.

The defaultValue is optional. This is a literal value or an expression which evaluates to a constant, and it can be specified as an attribute or as the element value. If specified the default value must match the type of the variable (otherwise it is a Schema Definition Error).

Note that the syntax supports both a defaultValue attribute and the default value being specified by the element value. Only one or the other may be present (otherwise it is a Schema Definition Error). To set the default value to "" (empty string), the defaultValue attribute syntax must be used, or the expression { "" } must be used as the element value.

Note the value of the name attribute is an NCName. The name of a variable is defined in the target namespace of the schema containing the definition. If multiple dfdl:defineVariable definitions have the same 'name' attribute in the same namespace then it is a Schema Definition Error.

A default *instance* of the variable is created (with global scope). Further instances of the variable may subsequently be created on schema elements. If the variable has a default value, this will be used as the default value for any *instances* of the variable (unless overridden when the instance is created).

The external property is optional. If not specified it takes the default value 'false'. If true the value may be provided by the DFDL processor and this external value will be used as the global default value (overriding any defaultValue specified on the dfdl:defineVariable). The mechanism by which the processor provides this value is implementation-defined.

There is no required order between dfdl:defineVariable and other schema level defining annotations or a dfdl:format annotation that may refer to the variable.

A defaultValue expression is evaluated before processing the data stream begins.

A defaultValue expression can refer to other variables but not to the infoset (so no path locations). The referenced variable must either have a defaultValue or be external. It is a Schema Definition Error otherwise.

If a defaultValue expression references another variable then that prevents the referenced variable's value from ever changing, that is, it is considered to be a read of the variable's value.

If a defaultValue expression references another variable and this causes a circular reference, it is a Schema Definition Error.

It is a Schema Definition Error if the type of the variable is a user-defined simple type restriction.

#### 7.7.1 Examples

```
<dfdl:defineVariable name="EDIFACT_DS" type="xs:string"
  defaultValue="," />

<dfdl:defineVariable name="codepage" type="xs:string"
  external="true">utf-8</dfdl:defineVariable>
```

#### 7.7.2 Predefined Variables

The following variables are predefined

Name	Namespace URI	Type	Default value	External
encoding	http://www.ogf.org/dfdl/dfdl-1.0/	xs:string	'UTF-8'	true
byteOrder	http://www.ogf.org/dfdl/dfdl-1.0/	xs:string	'bigEndian'	true
binaryFloatRep	http://www.ogf.org/dfdl/dfdl-1.0/	xs:string	'ieee'	true
outputNewLine	http://www.ogf.org/dfdl/dfdl-1.0/	xs:string	'%LF;'	true

**Table 8 Pre-defined variables**

These variables are expected to be commonly set externally so are predefined for convenience.

```
<xs:element name="title" type="xs:string">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:element encoding="{ $dfdl:encoding}" />
    </xs:appinfo>
  </xs:annotation>
</xs:element>
```

*- Is this a property of a variable or what's the difference?*

**7.8 The dfdl:newVariableInstance Statement Annotation Element**

Scoped instances of defined variables are created using dfdl:newVariableInstance:

```
<dfdl:newVariableInstance
  ref = QName
  defaultValue? = logical value or dfdl expression >
  <!-- Contains: logical value or dfdl expression (value) -->
</dfdl:newVariableInstance>
```

Since an initial instance is created when the variable is defined, the use of dfdl:newVariableInstance is optional. It would be used if an instance with restricted scope is needed.

The dfdl:newVariableInstance annotation can be used on a group reference, sequence or choice only. It is a Schema Definition Error otherwise.

The scope of the instance of a variable is the *dynamic scope* of the schema component and its content model and so is inherited by any contained constructs or construct references.

The ref property is a QName. That is, it may be qualified with a namespace prefix.

An optional default value *for the instance* may be specified. It can be specified as an attribute or as the element value. The expression must not contain forward references to elements which have not yet been processed nor to the current component. If specified the default value must match the type of the variable as specified by dfdl:defineVariable. If the instance is not assigned a new default value then it will inherit the default value specified by dfdl:defineVariable or externally provided by the DFDL processor. If a default value is not specified (and has not been specified by dfdl:defineVariable) then the value of this instance is undefined until explicitly set (using dfdl:setVariable).

If a default value is specified this initial value of the instance will be set when the instance is created. The value will override any (global) default value which was specified by dfdl:defineVariable or which was provided externally to the DFDL processor. A variable instance with a valid value (specified or default) can be referenced anywhere within the scope of the element on which the instance was created.

Note that the syntax supports both a default value attribute and the default value being specified by the element value. Only one or the other may be present. (Schema definition error otherwise.)

To set the default value to "" (empty string), the default value attribute syntax must be used, or the expression { "" } must be used as the element value.

The resolved set of annotations for a component may contain multiple dfdl:newVariableInstance statements. They must all be for unique variables, it is a Schema Definition Error otherwise. The order of execution is specified in Section 9.5 Evaluation Order for Statement Annotations.

There is no short form syntax for creating variable instances.

**7.8.1 Examples**

```
<dfdl:newVariableInstance ref="EDIFACT_DS" defaultValue="","/>

<dfdl:newVariableInstance ref="lengthUnitBits">
  { if (../hdr/fmtCode eq "bits") then 1 else 8 }
</dfdl:newVariableInstance>
```

### 7.9 The dfdl:setVariable Statement Annotation Element

Variable instances get their values either by default, by external definition, or by subsequent assignment using the dfdl:setVariable statement annotation. *(I needed to know this at the top of the variable section!)*

```
<dfdl:setVariable
  ref = QName
  value? = logical value or dfdl expression >
  <!-- Contains: logical value or dfdl expression (value) -->
</dfdl:setVariable>
```

The dfdl:setVariable annotation can be used on a simpleType, group reference, sequence or choice. It may be used on an element or element reference only if the element is of simple type. It is a Schema Definition Error if dfdl:setVariable appears on an element of complex type, or an element reference to an element of complex type. *Why?*

The ref property is a QName. That is, it may be qualified with a namespace prefix.

The syntax supports both a value attribute and the 'value' being specified by the element value. Only one or the other may be present (otherwise it is a Schema Definition Error). To set the value to "" (empty string), the value attribute syntax must be used, or the expression { "" } must be used as the element value.

The value must match the type of the variable as specified by dfdl:defineVariable.

A dfdl:setVariable value expression may refer to the value of this element using a relative path value ".". Use of relative path expressions is recommended wherever possible as this will allow the behavior of the parser to be more effectively scoped. However this practice is not enforced and there may be situations in which use of an absolute path is in fact necessary.

*Recall*  
The expression must not contain forward references to elements which have not yet been processed.

The declaration of a variable must be in scope at the point of the assignment, and at the point of reference.

In normal processing, the value of an instance can only be set once using dfdl:setVariable. Attempting to set the value of the variable instance for a second time is a Schema Definition Error. In addition, if a reference to the variable's value has already occurred and returned a default or an externally supplied value, then no assignment (even a first one) can occur. An exception to this behavior occurs whenever the DFDL processor backtracks because it is processing multiple branches of a choice or as a result of speculative parsing. In this case the variable state is also rewind.

A dfdl:setVariable will override any default value specified on either dfdl:defineVariable or dfdl:newVariableInstance, or externally. *was this defined?*

The resolved set of annotations for an **annotation point** may contain multiple dfdl:setVariable statements. They must all be for unique variables and it is a Schema Definition Error otherwise. The order of execution is specified in Section 9.5 Evaluation Order for Statement Annotations.

There is no short form syntax for variable assignment.

#### 7.9.1 Examples

```
<xs:element name="ds" type="xs:string">
  <xs:annotation>< xs:appinfo source="http://www.ogf.org/dfdl/">
    <dfdl:setVariable ref="EDI:EDIFACT_DS" value="{.}" />
    <dfdl:setVariable ref="delta" {..} </dfdl:setVariable>
  </xs:appinfo></xs:annotation>
</xs:element>
```

In the above example, the element named "ds" contains the string to be used as the EDI:EDIFACT\_DS delimiter at other places in the data, so the above defines the value of the EDI:EDIFACT\_DS variable to take on the value of this element.

*For clarity, with mentioning that you are also setting variable "delta" to the same value*

## 8. Property Scoping Rules

This section describes the rules that govern the scope over which DFDL representation properties apply

The scope of the representational properties on each of the component format annotations is given in **Table 9 DFDL annotation scoping**

Annotation Point	Property Scope
Schema declaration	dfdl:format representation properties apply <i>lexically</i> as default properties over all components in the schema
Element declaration	dfdl:element properties apply locally
Element reference	dfdl:element properties apply locally
Simple type definition	dfdl:simpleType properties apply locally
Sequence	dfdl:sequence properties apply locally
Choice	dfdl:choice properties apply locally
Group reference	dfdl:group properties apply locally

**Table 9 DFDL annotation scoping**

Note: This table lists DFDL annotations on schema components. DFDL annotations can also be placed on other DFDL annotations, such as a dfdl:format within a dfdl:defineFormat, to provide a named reusable format definition. In this case the annotation applies only where the named format is referenced.

DFDL representation properties explicitly defined on annotations, other than a dfdl:format on an xs:schema declaration, apply locally to that component only. The explicitly defined properties are the combination of any defined locally on the annotation and any defined on the dfdl:defineFormat annotation referenced by a local dfdl:ref property. When a property is defined both locally and on the dfdl:defineFormat, the locally defined property takes precedence.

The dfdl:format annotation on the top level xs:schema declaration provides defaults for the DFDL representation properties at every DFDL-annotatable component contained in the schema document. They do not apply to any components in any included or imported schema document (these may have their own defaults).

### 8.1 Providing Defaults for DFDL properties

A dfdl:format annotation on the top level xs:schema declaration may provide defaults for some or all the DFDL representation properties at every annotation point within the schema document. The default properties may be specified in attribute or element form. (Short form is not allowed on the xs:schema element.)

The dfdl:ref property is not a representation property so no default can be set.

The dfdl:escapeSchemeRef property provides a default reference to a dfdl:defineEscapeScheme, the properties of dfdl:escapeScheme are not defaulted individually.

DFDL representation properties defined explicitly on a component apply only to that component and override the default value of that property provided by a default format specified by an xs:schema dfdl:format annotation.

The example below demonstrates the overriding of the encoding property. The value 'ASCII' is the default value for the `title` element, but then it is overridden by the locally defined `utf-8` value for the encoding property, which takes precedence.

```
<xs:schema>
```

```

<xs:annotation>
  <xs:appinfo source="http://www.ogf.org/dfdl/">
    <dfdl:format encoding="ASCII" />
  </xs:appinfo>
</xs:annotation>

<xs:element name="book">
  <xs:complexType>
    <xs:sequence>
      <xs:element name="title" type="xs:string">
        <xs:annotation>
          <xs:appinfo source="http://www.ogf.org/dfdl/">
            <dfdl:element encoding="utf-8" />
          </xs:appinfo>
        </xs:annotation>
      </xs:element>
      <xs:element name="pages" type="xs:int"/>
    </xs:sequence>
  </xs:complexType>
</xs:element>
</xs:schema>

```

## 8.2 Combining DFDL Representation Properties from a dfdl:defineFormat

The DFDL representation properties contained in a referenced dfdl:defineFormat are combined with any DFDL representation properties defined locally on a construct as if they had been defined locally. If the same property is defined locally in and in the referenced dfdl:defineFormat then the local property takes precedence. The combined set of explicit DFDL properties has precedence over any defaults set by a dfdl:format on the xs:schema.

```

<xs:schema>
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:defineFormat name='myFormat'>
        <dfdl:format encoding="ASCII" />
      </dfdl:defineFormat>
    </xs:appinfo>
  </xs:annotation>

  <xs:element name="book">
    <xs:complexType>
      <xs:sequence>
        <xs:element name="title" type="xs:string">
          <xs:annotation>
            <xs:appinfo source="http://www.ogf.org/dfdl/">
              <dfdl:element ref='myFormat' encoding="UTF-8" />
            </xs:appinfo>
          </xs:annotation>
        </xs:element>
        <xs:element name="pages" type="xs:int"/>
      </xs:sequence>
    </xs:complexType>
  </xs:element>
</xs:schema>

```

The example above demonstrates the overriding of an encoding property. The 'ASCII' format encoding from the 'myFormat' is overridden by the UTF-8 format encoding, which as a locally defined property takes precedence.

## 8.3 Combining DFDL Properties from References

The DFDL properties from the following types of reference are combined using the rules below:

- An xs:element and its referenced xs:simpleType restriction,
- An xs:element reference and its referenced global xs:element

- An xs:group reference and an xs:sequence or xs:choice in its referenced global xs:group
- An xs:simpleType restriction and its base xs:simpleType restriction

## Rules

1. Create <sup>(1)</sup> an empty working set of "explicit" properties. <sup>and (2)</sup> Create an empty working set of "default" properties.
2. Move to the innermost schema component in the chain of references.
3. Assemble its applicable "explicit" properties from its local dfdl:ref (if present) and its local properties (if present), the latter overriding the former (that is, local wins over referenced).
4. Combine these with the current working set of "explicit" properties. It is a Schema Definition Error if the same property appears twice. The result is a new working set of "explicit" properties.
5. Obtain applicable "default" properties from a dfdl:format annotation on the xs:schema that contains the component (if such annotation is present). Combine these with the current working set of "default" properties, the latter overriding the former (that is, inner wins). Result is a new working set of "default" properties.
6. Move to the schema component that references the current component, and repeat starting at step 3. If there is no referencing component, carry out step 5 and then go to step 7.
7. Combine the resultant sets of properties. The "explicit" properties take priority, "defaults" only used when no "explicit" is present. It is a Schema Definition Error if a required property is in neither the "explicit" nor the "default" working sets.

"Applicable" properties are all the DFDL properties that apply to that schema component. For example all the DFDL properties that apply to a particular xs:simpleType (as defined by section 13).

```
<xs:simpleType name="newType">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:simpleType alignment="16"/>
    </xs:appinfo>
  </xs:annotation>
  <xs:restriction base="xs:integer">
    <xs:maxInclusive value="10"/>
  </xs:restriction>
</xs:simpleType>

<xs:element name="testElement1" type="newType">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:element representation="binary"/>
    </xs:appinfo>
  </xs:annotation>
</xs:element>
```

*bold face*

The locally defined dfdl:alignment property with value '16' from the xs:simpleType 'newType' is combined with the locally defined dfdl:representation property with value 'binary' and applied to element 'testElement1'.

```
<xs:simpleType name="otherNewType">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:simpleType alignment="64"/>
    </xs:appinfo>
  </xs:annotation>
  <xs:restriction base="newType">
    <xs:maxInclusive value="5"/>
  </xs:restriction>
</xs:simpleType>
```

```

<xs:simpleType name="newType">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:simpleType representation='binary' />
    </xs:appinfo>
  </xs:annotation>
  <xs:restriction base="xs:int">
    <xs:maxInclusive value="10"/>
  </xs:restriction>
</xs:simpleType>

```

*Bdf face*

The locally defined dfdl:representation property with value 'binary' is combined with the locally defined dfdl:alignment property with value '64' from the xs:simpleType restriction 'otherNewType'.

```

<xs:sequence>
  <xs:element ref="testElement1">
    <xs:annotation>
      <xs:appinfo source="http://www.ogf.org/dfdl/">
        <dfdl:element binaryNumberRep="binary" />
      </xs:appinfo>
    </xs:annotation>
  </xs:element>
</xs:sequence>

<xs:element name="testElement1" type="newType">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:element representation="binary" />
    </xs:appinfo>
  </xs:annotation>
</xs:element>

<xs:simpleType name="newType">
  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:simpleType alignment="16" />
    </xs:appinfo>
  </xs:annotation>
  <xs:restriction base="xs:int">
    <xs:maxInclusive value="10"/>
  </xs:restriction>
</xs:simpleType>

```

*64*

The locally defined dfdl:alignment property with value '16' from the xs:simpleType 'newType' is combined with the locally defined dfdl:representation property with value 'binary' and locally defined dfdl:binaryNumberRep with a value of 'binary'

```

<!-- SCHEMA1 -->
<xs:schema targetNamespace="" xmlns:tns1="http://tns1">

  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:format encoding="ASCII" byteOrder="littleEndian"
        initiator="" terminator=""
        sequenceKind="ordered" />
    </xs:appinfo>
  </xs:annotation>

  <xsd:import namespace="http://tns2" schemaLocation="SCHEMA2.xsd"/>

  <xs:element name="book">
    <xs:complexType>
      <xs:group ref="tns2:ggrp1" dfdl:separator=","/>
    </xs:complexType>
  </xs:element>

```

```
</xs:complexType>
</xs:element>

</xs:schema>
```

```
<!-- SCHEMA2 -->
<xs:schema targetNamespace="" xmlns:tns2="http://tns2">

  <xs:annotation>
    <xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:format encoding="UTF-8" byteOrder="littleEndian"
        initiator=""
        sequenceKind="ordered" />
    </xs:appinfo>
  </xs:annotation>

  <xs:group name="ggrp1" >
    <xs:sequence dfdl:separatorPosition="infix" >
      <xs:element name="customer" type="xs:string"
        dfdl:length="8" dfdl:lengthKind="explicit" />
    </xs:sequence>
  </xs:group>
</xs:schema>
```

The DFDL properties applied to the `xs:sequence` in `xs:group "ggrp1"` in SCHEMA2 when referenced from the group reference in SCHEMA1 are

1. `dfdl:separator ","` from the group reference in SCHEMA1
2. `dfdl:separatorPosition "infix"` from the group declaration in SCHEMA2
3. `dfdl:encoding "UTF-8"`, `dfdl:initiator ""` from the default `dfdl:format` annotation in SCHEMA2
4. `dfdl:terminator ""` from the default `dfdl:format` annotation in SCHEMA1



## 9. DFDL Processing Introduction

*Should major headings be bigger?  
easy to lose the structure*

A *DFDL Parser* is an application or code library that takes as input:

- A DFDL annotated schema
- A data stream.

It is able to use the DFDL schema description to interpret the data stream and realize the DFDL Information Model. This information set could then be written out (for example it could be realized as an XML text string) or it could be accessed by an application through an API (for example, a DOM-like tree could be created in memory for access by applications).

Symmetrically, there is a notion of a *DFDL Unparser*. The unparser works from an instance of the DFDL Information Model, a DFDL annotated schema and writes out to a target data stream in the appropriate representation formats.

Often both parser and unparser would be implemented in the same body of software and so we do not always distinguish them. Collectively they are called a *DFDL Processor*. The parser and unparser may, of course, be different bodies of software. Conforming DFDL processors may implement only a parser, because the unparser is an optional feature of DFDL.

### 9.1 Parser Overview

The DFDL logical parser is a recursive-descent parser<sup>7</sup> having guided, but potentially unbounded look ahead that is used to resolve *points of uncertainty*. A DFDL parser reads a specification (the DFDL schema) and it recursively walks down and up the schema as it processes the data. This is done in a manner consistent with the scoping of properties and variables described in Section 8 Property Scoping Rules.

*Can't remember if this was defined yet.*

The unbounded look ahead means that there are situations where the parser must speculatively attempt to parse data where the occurrence of a processing error causes the parser to suppress the error, back out and make another attempt.

Implementations of DFDL may provide control mechanisms for limiting the speculative search behavior of DFDL parsers. The nature of these mechanisms is beyond the scope of the DFDL specification which defines the behavior of conforming parsers *only on correct data*. That is, data that can be parsed without any effective processing errors. Any such control mechanisms must be documented by the implementation and are thus implementation-defined.

*new concept?*

The logical parser recursively descends the DFDL schema beginning with the *distinguished global element declaration* (specified for the processor in an implementation-defined manner, see Section 18). Depending on the kind of schema construct that is encountered and the DFDL annotations on it, and the pre-existing context, the parser performs specific parsing operations on the data stream. These parsing operations typically recognize and consume data from the stream and construct values in the logical model. For values of complex types and for arrays, these logical model values may incorporate values created by recursive parsing.

DFDL Implementations are free to use whatever techniques for parsing they wish so long as the semantics are equivalent to that of the speculative recursive-descent logical parser described in this specification. It is required that implementations distinguish the various kinds of errors (Schema Definition Error, processing error, etc.) no matter what time they are detected. Some implementations may not detect certain Schema Definition Errors until data are being parsed; however, they must still distinguish Schema Definition Errors (which indicate that the schema itself is not meaningful), from parsing errors (which indicate that the input data doesn't satisfy the requirements of the schema), or unparsing errors (which indicate that the infoset does not satisfy the requirements of the schema).

*Confused. "Correct" data could contain ambiguities that require backtracking to parse. An implementation could limit the search space even on correct data.*

<sup>7</sup> A "top-down" parser built from a set of mutually-recursive procedures or a non-recursive equivalent where each such procedure usually implements one of the productions of the grammar. Thus the structure of the resulting program closely mirrors that of the grammar it recognises. See [RDP].

9.2 DFDL Data Syntax Grammar

Data in a format describable via a DFDL schema obeys the grammar given here. A given DFDL schema is read by the DFDL processor to provide specific meaning to the terminals and decisions in this grammar.

The bits of the data are divided into two broad categories:

- 1. Content
- 2. Framing

The content is the bits of data that are interpreted to compute a logical value.

*Framing* is the term used to describe the delimiters, length fields, and other parts of the data stream which are present, and may be necessary to determine the length or position of the content of DFDL Infoset items.

Note that sometimes the framing is not strictly necessary for parsing, but adds useful redundancy to the data format, allowing corrupt data to be more robustly detected, and sometimes the framing adds human readability to the data format.

In the grammar tables below, the terminal symbols are shown in bold italic font.

Productions
Document = <b>SimpleElement</b>   <b>ComplexElement</b>
SimpleElement = SimpleLiteralNilElementRep   SimpleEmptyElementRep   SimpleNormalRep
SimpleEnclosedElement = SimpleElement   AbsentElementRep
ComplexElement = ComplexLiteralNilElementRep   ComplexNormalRep   ComplexEmptyElementRep
ComplexEnclosedElement = ComplexElement   AbsentElementRep
EnclosedElement = SimpleEnclosedElement   ComplexEnclosedElement
AbsentElementRep = <b>Absent</b>
SimpleEmptyElementRep = EmptyElementLeftFraming EmptyElementRightFraming
ComplexEmptyElementRep = EmptyElementLeftFraming EmptyElementRightFraming
EmptyElementLeftFraming = LeadingAlignment <b>EmptyElementInitiator</b> PrefixLength
EmptyElementRightFraming = <b>EmptyElementTerminator</b> TrailingAlignment
SimpleLiteralNilElementRep = NilElementLeftFraming [ <b>NilLiteralCharacters</b>   NilElementLiteralContent] NilElementRightFraming
ComplexLiteralNilElementRep = NilElementLeftFraming <b>NilLiteralValue</b> NilElementRightFraming
NilElementLeftFraming = LeadingAlignment <b>NilElementInitiator</b> PrefixLength

Deleted: = UnicodeByteOrderMark DocumentElement¶  
DocumentElement

Deleted:

Deleted: ¶

NilElementRightFraming = <b>NilElementTerminator</b> TrailingAlignment
NilElementLiteralContent = <b>LeftPadding NilLiteralValue</b> RightPadOrFill
SimpleNormalRep = LeftFraming PrefixLength SimpleContent RightFraming ComplexNormalRep = LeftFraming PrefixLength ComplexContent <b>RightFraming</b>
LeftFraming = LeadingAlignment <b>Initiator</b> RightFraming = <b>Terminator</b> TrailingAlignment
PrefixLength = SimpleContent   PrefixPrefixLength SimpleContent PrefixPrefixLength = SimpleContent
SimpleContent = <b>LeftPadding</b> [ <b>SimpleLogicalValue</b> ] RightPadOrFill <b>SimpleLogicalValue</b> = <b>SimpleNormalValue</b>   <b>NilLogicalValue</b>
<b>ComplexContent</b> = <b>ComplexValue</b> <b>ElementUnused</b> <b>ComplexValue</b> = Sequence   Choice
Sequence = LeftFraming SequenceContent RightFraming SequenceContent = [ <b>PrefixSeparator</b> EnclosedContent [ <b>Separator</b> EnclosedContent ]* <b>PostfixSeparator</b> ]
Choice = LeftFraming ChoiceContent RightFraming ChoiceContent = [ EnclosedContent ] <b>ChoiceUnused</b>
EnclosedContent = [ EnclosedElement   Array   Sequence   Choice ]
Array = [ EnclosedElement [ <b>Separator</b> EnclosedElement ]* [ <b>Separator</b> StopValue ]
StopValue = SimpleElement
LeadingAlignment = <b>LeadingSkip AlignmentFill</b> TrailingAlignment = <b>TrailingSkip</b> RightPadOrFill = <b>RightPadding</b>   <b>RightFill</b>   <b>RightPadding RightFill</b>

Deleted: **ElementUnused**¶

Deleted: ¶

Deleted:

Deleted: **NilLogicalValue** | **SimpleValue**

Deleted: ComplexContent = Sequence | Choice ¶

**Table 10 DFDL Grammar Productions**

XML Schema and DFDL properties are used to control constraints on the terminals of the above grammar, as well as repetition (the "\*" operator), and alternatives (the "|" operator). For a given set of XML Schema and DFDL properties, and prior data, any terminal may be allowed to be length zero, to contain specific data, or to contain a variety of different admissible data.

Some definitions are needed to cover the range of representations that are possible in the data stream for an element. These definitions are with respect to the grammar above.

### 9.2.1 Nil Representation

An element occurrence has a *nil representation* if the element has XSDL nillable property 'true' and the occurrence either:

- conforms to the grammar for SimpleNilLiteralElementRep or ComplexNilLiteralElementRep. Specifically, the **NilElementInitiator** and **NilElementTerminator** regions must be conformant with dfdl:nilValueDelimiterPolicy. (If non-conformant it is not a processing error and the representation is not nil).
- conforms to the grammar for SimpleNormalRep and its value is **NilLogicalValue**.

The LeadingAlignment, TrailingAlignment, PrefixLength regions may be present.

### 9.2.2 Empty Representation

An element occurrence has an *empty representation* if the occurrence does not have a nil representation and it conforms to the grammar for SimpleEmptyElementRep or ComplexEmptyElementRep. Specifically, the **EmptyElementInitiator** and **EmptyElementTerminator** regions must be conformant with dfdl:emptyValueDelimiterPolicy and the occurrence's SimpleContent or ComplexContent region in the data is of length zero. (If non-conformant it is not a processing error and the representation is not empty). LeadingAlignment, TrailingAlignment, PrefixLength regions may be present.

The empty representation is special in DFDL because when parsing it is used to determine when default values are created in the Infoset. The empty representation can require initiators or terminators be present so as to enable data formats which explicitly distinguish occurrences with empty string/hexBinary values from occurrences that are missing or are absent. See Section 9.4 Element Defaults below about default values.

The empty representation is not possible for fixed-length elements with a non-zero length.

### 9.2.3 Normal Representation

An element occurrence has a normal representation if the occurrence does not have the nil representation or the empty representation and it conforms to the grammar for SimpleNormalRep or ComplexNormalRep.

### 9.2.4 Absent Representation

An element occurrence has an *absent representation* if the occurrence does not have a nil or empty or normal representation, and it conforms to the grammar for AbsentElementRep. Specifically, the occurrence's representation in the data stream is of length zero. Consequently, the Initiator, Terminator, LeadingAlignment, TrailingAlignment, PrefixLength regions must not be present.

Example of an absent representation: During unparsing, if an optional element does not have an item in the infoset then nothing is output. However if a separator of an enclosing structure is subsequently output as the immediate next thing, then a subsequent parse of the element may return a representation of length zero. If this happens, and this zero-length representation does not conform to any of the nil representation, the empty representation, or the normal representation, then it is the absent representation, and it behaves as if the element occurrence is 'missing'. (The term 'missing' is defined below.)

The point of this term 'absent representation', is that often we know the location where an element or group's representation would be in the data based on the delimiters of an enclosing group. (An example: if there are adjacent delimiters of an enclosing sequence.) When this location in the data, which is of zero length, cannot be a nil, empty, or normal representation, then we say it has absent representation, or "the representation is absent".

### 9.2.5 Zero-length Representation

We use the term zero-length representation to describe the situations where any of the above representations turn out to be of length zero due to specific combinations of data type and format properties:

**Deleted:** the occurrence's content in the data stream

**Deleted:** The empty representation is special in DFDL, because when parsing it is this condition that can trigger the creation of a default value for an element occurrence.

After reading this section I do not understand the difference between "absent" and "empty" maybe examples would help?

The nil representation can be a zero-length representation if dfdl:nilValue is '%ES;' or '%WSP\*;' appearing on its own as a literal nil value, and there is no framing or framing is suppressed by dfdl:nilValueDelimiterPolicy.

Deleted: is "%ES;"

The empty representation can be a zero-length representation if there is no framing or framing is suppressed by dfdl:emptyValueDelimiterPolicy.

The normal representation can be a zero-length representation if the type is xs:string or xs:hexBinary and there is no framing.

What is meant by "no framing"?

The absent representation always has a zero-length representation.

If the nil representation may be zero-length, then the absent representation cannot occur because zero-length will be interpreted as nil representation.

If the nil representation may not be zero-length, but the empty representation is zero-length, then the absent representation cannot occur because zero-length will be interpreted as the empty representation.

If the nil and empty representations can not be zero-length, but the normal representation may be zero length then the absent representation cannot occur because zero length will be interpreted as a normal representation.

If the nil representation may not be zero-length, the empty representation is not zero-length, and the normal representation may not be zero-length, then a zero-length representation is the absent representation, or "is absent".

## 9.2.6 Missing

When parsing, an element occurrence is missing if it does not have nil, empty, or normal representations, or it has the absent representation.

When parsing, the term missing really covers two situations. Firstly it subsumes absent representation. Secondly it applies when an element does not have a representation at all in the data stream, that is, when we do not even have the constructs in the data stream to determine the location of the representation of the element; hence, none of the concepts above apply. This will be made clearer in the examples below. If an element occurrence is missing when parsing, no item is ever added to the Infoset.

When unparsing, an element occurrence is missing if there is no item in the infoset. For a required element occurrence, it is this condition that can trigger the creation of a default value in the augmented infoset. See Section 9.4 Element Defaults below about default values. For an optional element occurrence, no item is ever added to the augmented Infoset nor any representation ever output in the data stream.

## 9.2.7 Examples of Missing and Empty Representation

The following examples illustrate missing and empty representation.

```
<xs:sequence dfdl:separator="," dfdl:terminator="@"
  dfdl:separatorSuppressionPolicy="trailingEmpty">
  <xs:element name="A" type="xs:string"
    dfdl:lengthKind="delimited"/>
  <xs:element name="B" type="xs:string" minOccurs="0"
    dfdl:lengthKind="delimited"/>
  <xs:element name="C" type="xs:string" minOccurs="0"
    dfdl:lengthKind="delimited"/>
</xs:sequence>
```

→ extra line  
In data stream aaa,@ element B has the empty representation, and element C does not have a representation so is missing.

```
<xs:sequence dfdl:separator=","
  dfdl:separatorSuppressionPolicy="anyEmpty">
  <xs:element name="A" type="xs:string"
    dfdl:lengthKind="delimited" dfdl:initiator="A:"
```

There are 2 things that I assume are important about these definitions:

- 1) How to syntactically distinguish them
- 2) Semantic differences

Unfortunately neither are clear to me. Let me read the examples and see if they help.

```

dfdl:emptyValueDelimiterPolicy=initiator"/>
<xs:element name="B" type="xs:string" minOccurs="0"
dfdl:lengthKind="delimited" dfdl:initiator="B:"
dfdl:emptyValueDelimiterPolicy="initiator"/>
<xs:element name="C" type="xs:string" minOccurs="0"
dfdl:lengthKind="delimited" dfdl:initiator="C:"
dfdl:emptyValueDelimiterPolicy=initiator"/>
</xs:sequence>

```

*Examples are used to consider putting examples BEFORE defs. Without examples, definitions are incomprehensible.*

In data stream A:aaaa,C:cccc element B does not have a representation so is missing.

In data stream A:aaaa,B:,C:cccc element B has the empty representation.

In the data stream A:aaaa,,C:cccc element B has the absent representation so is missing.

### 9.2.8 Round Trip Ambiguities

The overlapping nature of the possible representations: normal, empty, nil, and absent, creates a number of ambiguities where taking an Infoset, unparsing it, and reparsing it will result in a second Infoset that is not the same as the original. However taking the second Infoset, unparsing it, and reparsing it, will result in a third Infoset which is the same as the second.

When unparsing, if a string Infoset item happens to contain a string that matches either one of the nilValues or the default value, it is not given any special treatment. The string's characters are output, or if the value is the empty string, zero length content is output. (In both cases along with an initiator or terminator if defined.) This creates an ambiguity where one can unparse an Infoset item which has member **[nilled]** true, but when reparsed will produce an Infoset item which has member **[nilled]** false.

These ambiguities are natural and unavoidable. If the nilValue is the 3-character string "nil", then encountering the characters "nil" in the data stream will parse to produce an Infoset item with **[nilled]** true in the Infoset. If you unparsed a string infoset item with contents of the 3 characters "nil", this will be output as the letters "nil", which on parse will not produce a string with the characters "nil", but rather an Infoset item with no data value and member **[nilled]** true.

To avoid this issue, one can use validation, along with a pattern that prevents the string from matching any of the nil values.

### 9.3 Parsing Algorithm

A DFDL parser proceeds by determining the existence of occurrences of schema components. It does this by examining the data and the schema, so as to:

- a) Establish representation
- b) Resolve points of uncertainty

These two activities are defined below. They are mutually recursive in the expected way as a DFDL schema is a recursive nest of schema components.

Establishing the representation of an occurrence of a schema component and resolving points of uncertainty involve the concepts of known-to-exist and known-not-to-exist.

#### 9.3.1 Known-to-exist and Known-not-to-exist

##### 9.3.1.1 Known-to-exist

An occurrence of a schema component is said to be known-to-exist when any of these positive discriminations hold:

*with elaboration. Difference between this and the first example is subtle.*

*Consider a Venn diagram to show overlap/intersection of these concepts*

1. There is a `dfdl:discriminator`<sup>8</sup> applying to the component and its expression evaluates to true or regular expression pattern matches.
2. The component is a direct child of an `xs:sequence` or `xs:choice` with `dfdl:initiatedContent` 'yes' and an initiator defined for the component is found.
3. The component is a direct child of an `xs:choice` with `dfdl:choiceDispatchKey` and the result of the `dfdl:choiceDispatchKey` expression matches one of the `dfdl:choiceBranchKey` property values of the child.

If none of those hold because they are not applicable then the occurrence is still known-to-exist if ALL of the following hold, and no processing error occurs during their determination:

1. When there are `dfdl:asserts` with `failureType` 'processingError' on the component, all their expressions evaluate to true or their regular expression patterns match.
2. It has nil, empty, or normal representation.
3. When it has normal representation the content of the representation is convertible to the element type without error.

Note that validation errors or recoverable errors do not prevent determination that a component is known-to-exist.

### 9.3.1.2 Processing Error After Determining Known-to-exist

Note that it is possible for an occurrence of a schema component to be known-to-exist due to a positive discrimination, but then subsequently a processing error occurs when evaluating a statement annotation such as a `dfdl:assert` or a `dfdl:setVariable`, or a processing error occurs when determining the representation, or in the case of normal representation and `simpleType`, when converting that representation's content into a value of the type. This processing error does not change the fact that the schema component was determined to be known-to-exist. This is important in the discussion of resolving Points of Uncertainty below.

### 9.3.1.3 Known-not-to-exist

An occurrence of a schema component is *known-not-to-exist* when any of these *negative discriminations* holds:

1. There is a `dfdl:discriminator` applying to the component and its expression evaluates to false or regular expression pattern fails to match, or a processing error occurs while processing the `dfdl:discriminator`.
2. The component is a direct child of an `xs:sequence` or `xs:choice` with `dfdl:initiatedContent` 'yes' and an initiator defined for the component is not found.
3. The component is a direct child of an `xs:choice` with `dfdl:choiceDispatchKey` and the result of the `dfdl:choiceDispatchKey` expression does not match any of the `dfdl:choiceBranchKey` property values of the child.

If none of those hold because they are not applicable, then a schema component is known-not-to-exist when any of the following hold:

1. The occurrence is missing
2. There is a `dfdl:assert` with `failureType` 'processingError' on the component and its expression evaluates to false or its regular expression pattern fails to match, or a processing error occurs while processing the `dfdl:assert`.
3. A processing error occurs when parsing the component. Processing errors include, but are not limited to, inability to identify any of nil, empty, normal or absent representations, or failure to convert a value to the built-in logical type.

<sup>8</sup> DFDL discriminators are described in section: 7.4 The `dfdl:discriminator` Statement Annotation Element.

Note that validation errors or recoverable errors do not cause a component to be known-not-to-exist.

Note: based on the above, when processing a sequence for which a separator is defined, the presence of a match in the data for the separator is not sufficient to cause the parser to determine that an associated component is known-to-exist. See Section 14.2 Sequence Groups with Separators for details.

### 9.3.2 Establishing Representation

Unless an element occurrence is known-not-to-exist, it must be established if it has the nil, empty, normal, or absent representation.

The first step is to see if the SimpleContent or ComplexContent region is of length zero as a first approximation. This is dfdl:lengthKind dependent.

- explicit => length is zero (either fixed or from expression evaluation)
- prefixed => length given by the prefix is zero
- implicit (simple) => length is zero<sup>9</sup>
- implicit (complex) => not possible.
- delimited => length is zero (in scope delimiter is immediately encountered)
- pattern => pattern returns zero length match
- endOfParent => already positioned at parent's end so length is zero

#### 9.3.2.1 Simple element

If the result is length zero as described above, the representation is then established by checking, in order for:

1. nil representation (if %ES; or %WSP\* on its own is a literal nil value).
2. empty representation.
3. normal representation (xs:string or xs:hexBinary only)
4. absent representation (if none of the prior representations apply).

If the result is not length zero, the representation is then established by checking, in order, for:

1. nil representation (as a literal nil value)
2. nil representation (as a logical nil value)
3. normal representation

#### 9.3.2.2 Complex element

If the result is length zero as described above, the representation is then established by checking for:

- nil representation (if %ES; is a literal nil value).<sup>10</sup>

To establish any other representations requires that the parser descends into the complex type for the element, and returns successfully (that is, no unsuppressed processing error occurs). If the result is zero bits consumed, the representation is then established by checking, in order, for:

1. empty representation.
2. absent representation (if none of the prior representations apply).

Otherwise the element has normal representation.

<sup>9</sup> This is a corner case that only happens when type is xs:string or xs:hexBinary and the maxLength facet is 0. Such an element can only be of length 0.

<sup>10</sup> It is a Schema Definition Error if a complex element has XSDL nillable 'true' and dfdl:lengthKind 'implicit'.

**Deleted:** The first step is to see if the content is trivially of length zero.



Note: The DFDL parser shall not recursively parse the schema components inside a complex element when it has already established that the element occurrence is missing<sup>11</sup>.

### 9.3.3 Points of Uncertainty - Finally!

A point of uncertainty occurs in the data stream when there is more than one schema component that might occur at that point. Points of uncertainty can be nested.

Any one of the following constructs is a potential point of uncertainty:

- An `xs:choice`
- All `xs:elements` in an unordered `xs:sequence` (`dfdl:sequenceKind` is 'unordered')
- An optional `xs:element`
- An array `xs:element`.
- All `xs:elements` in an `xs:sequence` containing one or more floating `xs:elements`.

The parser resolves these points of uncertainty by way of a set of construct-specific rules given below along with determining whether schema components are known-to-exist or known-not-to-exist. For some of these constructs, whether there is an actual point of uncertainty depends on the representation of the constructs in the data.

An `xs:choice` is always a point of uncertainty. It is resolved sequentially, or by direct dispatch.

Sequential choice resolution occurs by parsing each choice branch in schema definition order until one is known-to-exist. It is a processing error if none of the choice branches are known-to-exist. Direct-dispatch choice resolution occurs by matching the value of the

`dfdl:choiceDispatchKey` property to the value of one of the `dfdl:choiceBranchKey` property values of one of the choice branches. It is a processing error if none of the choice branches have a matching value in their `dfdl:choiceBranchKey` property.

An element in an unordered `xs:sequence` is always a point of uncertainty. It is resolved by parsing for the child components of the sequence in schema definition order at each point in the data stream where a component can exist until the required number of occurrences of each child component is known- to-exist or the sequence is terminated by delimiters or specified length.

An element in a sequence with one or more floating elements is always a point of uncertainty. It is resolved by parsing for the expected element at that point in the data stream. If the expected element is known-not-to-exist then an occurrence of each floating element is parsed in schema definition order.

When parsing an array, points of uncertainty only occur for certain values of `occursCountKind`, as follows:

<code>occursCountKind</code>	Details of Point of Uncertainty
fixed	No point of uncertainty (maxOccurs occurrences expected).
implicit	A point of uncertainty exists after minOccurs occurrences found and until maxOccurs found.
parsed	A point of uncertainty exists for all occurrences
expression	No point of uncertainty ( <code>occursCount</code> occurrences expected)
stopValue	No point of uncertainty ( <code>stopValue</code> must always be present, even when minOccurs is 0).

**Table 11: Points of Uncertainty and `dfdl:occursCountKind`**

<sup>11</sup> The rationale for this is that otherwise this could give rise to misleading error messages where the parser reported that required child elements were missing required occurrences. (This is consistent with XML Schema validation, where if a required element is missing, it gets reported as such, and there is nothing reported about its children).

An optional element point of uncertainty is resolved by parsing the element until it is either known-to-exist or known-not-to-exist. Whether an optional element is an actual point of uncertainty depends on property `dfdl:occursCountKind` as described above. (Property `dfdl:occursCountKind` is defined in Section 16.1 `dfdl:occursCountKind` property.)

For an array element, the point of uncertainty is resolved for each occurrence separately by parsing the occurrence until it is either known-to-exist or known-not-to-exist.

#### 9.3.3.1 Nested Points of Uncertainty

A point of uncertainty can be resolved because a schema component has been determined to be known-to-exist due to positive discrimination. In that case, if a subsequent processing error occurs when completing the parsing of that schema component this will cause the next enclosing schema component surrounding this point of uncertainty to be determined to be known-not-to-exist.

For example, when parsing an element occurrence for an array with a variable number of occurrences, a positive discrimination tells the parser that the currently-being-parsed occurrence is known-to-exist. If a subsequent processing error occurs while completing the parsing of this occurrence, then the entire array is then known-not-to-exist.

Another example is a choice. If a discriminator resolves the choice point of uncertainty to the first of the choice's alternatives, a subsequent processing error causes the entire choice construct to be determined to be known-not-to-exist.

This will cause the next enclosing point of uncertainty to try the next possible alternative, or if there isn't one, will cause an unsuppressed processing error.

The behavior of a DFDL processor on an unsuppressed processing error is not specified, but it is allowable for implementations to abort further parsing. Any other behavior is implementation-defined.

### 9.4 Element Defaults

A DFDL processor can create element defaults in the Infoset for both simple and complex elements. This happens quite differently for parsing and unparsing as will be explained in this section.

#### 9.4.1 Definition 'default value'

A simple element has a default value if any of these are true:

1. The XSDL default property exists. The default value is the property's value.
2. The XSDL fixed property exists. The default value is the property's value.
3. The element has XSDL nillable is 'true' and `dfdl:useNilForDefault` is 'yes'. The corresponding Infoset item will have the `[nilled]` member true, and the `[dataValue]` member will have no value.

*Do I have what this is?*

#### 9.4.2 Element Defaults When Parsing

If *empty* representation is established when parsing, the possibility of applying an element default arises. Essentially, if a required occurrence of an element has empty representation, then an element default will be applied if present, though there are a couple of variations on this rule. Remember that in order to have established empty representation, the occurrence must be compliant with the `dfdl:emptyValueDelimiterPolicy` for the element, and for a complex element the parser must have descended into the type and returned with no unsuppressed processing error.

The rules for applying element defaults are not dependent on `dfdl:occursCountKind`. However, if a required occurrence does not produce an item in the Infoset after the rules have been applied, then whether it is a processing error or a validation error (if validation is enabled) *does* depend on `dfdl:occursCountKind` (see Section 16.1 `dfdl:occursCountKind` property).

The sections below indicate when an item is added to the infoset, and whether it has a default or other value. If there is no processing error then regardless of whether an item is added to the infoset or not, any side-effects due to `dfdl:discriminator` statements evaluating to true, or `dfdl:setVariable` statements, are retained.

Assuming "empty" representation is established

There are three main cases to consider:

#### 9.4.2.1 Simple element (not xs:string and not xs:hexBinary)

Required occurrence: If the element has a default value then an item is added to the Infoset using the default value, otherwise nothing is added to the Infoset.

Optional occurrence: Nothing is added to the Infoset.

#### 9.4.2.2 Simple element (xs:string or xs:hexBinary)

Required occurrence: If the element has a default value then an item is added to the infoset using the default value, otherwise an item is added to the Infoset using empty string (type xs:string) or empty hexBinary (type xs:hexBinary) as the value.

Optional occurrence: if `dfdl:emptyValueDelimiterPolicy` is applicable and is not 'none',<sup>12</sup> then an item is added to the Infoset using empty string (type xs:string) or empty hexBinary (type xs:hexBinary) as the value, otherwise nothing is added to the Infoset.

Note: To prevent unwanted empty strings or empty hexBinary values from being added to the Infoset, use XSD `minLength > '0'` and a `dfdl:assert` that uses the `dfdl:checkConstraints()` function, to raise a processing error.

#### 9.4.2.3 Complex element

Required occurrence: An item is added to the Infoset.

Optional occurrence: if `dfdl:emptyValueDelimiterPolicy` is applicable and is not 'none',<sup>13</sup> then an item is added to the Infoset, otherwise nothing is added to the Infoset.

A complex element can parse, by recursive descent, and construct a complex element in the infoset containing a single child element. This can occur:

1. If the first child element of the complex type is a required simple element, then an empty string (type xs:string), empty hexBinary (type xs:hexBinary), or default value will also be added to the Infoset.
2. If the first child element of the complex type is a required complex element, then an item is added to the Infoset (which may itself have a child via (1))

As an example, consider a sequence S0 with a separator that contains among other content an optional non-nillable non-initiated element E1 of complex type. The content of the type is a sequence S1 with a different separator and the first child is a required non-initiated element E2 of type xs:string. The `dfdl:lengthKind` of both E1 and E2 is 'delimited'. The representation of E1 has zero length, that is, the data contains adjacent S0 separators. On processing E1, the parser will establish a point of uncertainty and descend into E1's complex type and process E2. It scans for in-scope delimiters and immediately encounters S0 separator. E2 has the empty representation, so E1 is added to the Infoset along with a value of empty string for E2. All other content of S1 is missing, so the parser returns from the descent. E1 is therefore known-to-exist. Because the position in the data has not changed, E1 therefore has the empty representation. Because E1 is empty and optional it is not added to the Infoset, and the Infoset items for E1 and E2 are discarded.

#### 9.4.3 Element Defaults When Unparsing

If an element is *missing* from the Infoset when unparsing, the possibility of applying an element default arises. Essentially if a required occurrence of an element is missing, then an element default will be applied if present, and the resulting item is added to the *augmented Infoset*.

The rules for applying element defaults are not dependent on `dfdl:occursCountKind`. However if a required occurrence does not produce an item in the augmented Infoset after the rules have been

<sup>12</sup> If other than 'none', either an initiator, terminator or both must have been found in the data stream.

<sup>13</sup> If other than 'none', either an initiator, terminator or both must have been found in the data stream.

**Deleted:** Optional occurrence: If `dfdl:emptyValueDelimiterPolicy` is not 'none'

**Deleted:**

**Deleted:** Optional occurrence: If `dfdl:emptyValueDelimiterPolicy` is not 'none'

**Deleted:**

**Commented [1]:** Remove 'when' - each of the following bullets starts with 'If'.

**Deleted:** For both required and optional occurrences, the Infoset item may also have a child item.

have to follow. Picture?

applied then whether it is a processing error or a validation error (if enabled) is dependent on dfdl:occursCountKind (see Section 16.1 dfdl:occursCountKind property).

There are two main cases to consider.

#### 9.4.3.1 Simple element

Required occurrence: If an element has a default value then an item is added to the augmented Infoset using the default value, otherwise nothing is added.

Optional occurrence: Nothing is added to the augmented Infoset.

#### 9.4.3.2 Complex element

Required occurrence: An item is added to the augmented Infoset as specified below.

Optional occurrence: Nothing is added to the augmented Infoset.

For a required occurrence, the unparser descends into the complex type:

For a sequence, each child element is examined in schema order and the rules for simple and complex elements applied (recursively). The lack of a default may give rise to a processing error, as described above.

For a choice, each branch is examined in schema order and the above rules applied recursively to the branch. The lack of a default may give rise to a processing error, as described above, and if so the error is suppressed and the next branch is tried, otherwise that branch is selected. It is a processing error if no choice branch is ultimately selected.

### 9.5 Evaluation Order for Statement Annotations

Given a component of a DFDL schema, there is a resolved set of annotations for it.

Of these, some are statement annotations and the order of their evaluation relative to the actual processing of the schema component itself (parsing or unparsing via its format annotations) is as given in the ordered lists below.

For elements and element references:

1. dfdl:discriminator or dfdl:assert(s) with testKind 'pattern' (parsing only)
2. dfdl:element following property scoping rules, which includes establishing representation as described in Section 9.3.2 and conversion to the element type for simple types.
3. dfdl:setVariable(s) - in lexical order, innermost schema component first
4. dfdl:discriminator or dfdl:assert(s) with testKind 'expression' (parsing only)

Deleted: dfdl:element following property scoping rules¶

For sequences, choices and group references:

1. dfdl:discriminator or dfdl:assert(s) with testKind 'pattern' (parsing only)
2. dfdl:newVariableInstance(s) - in lexical order, innermost schema component first
3. dfdl:setVariable(s) - in lexical order, innermost schema component first
4. dfdl:sequence or dfdl:choice or dfdl:group following property scoping rules and evaluating any property expressions (corresponds to ComplexContent grammar region)
5. dfdl:discriminator or dfdl:assert(s) with testKind 'expression' (parsing only)

The dfdl:setVariable annotations at any one annotation point of the schema are always executed in lexical order. However, dfdl:setVariable annotations can also be found in different annotation points that are combined into the resolved set of annotations for one schema component. In this case, the order of execution of the dfdl:setVariable statements from any one annotation point remains lexical. The order of execution of the dfdl:setVariable annotations different annotation points follows the principle of innermost first, meaning that a schema component that references another schema component has its dfdl:setVariable statements executed after those of the referenced schema component. For example, if an element reference and an element declaration both have dfdl:setVariable statements, then those on the element declaration will execute before those on the element reference. Similarly dfdl:setVariable statements on a base simple type execute before those of a simple type derived from it. The dfdl:setVariable statements on a simple type execute before those on an element having that simple type (whether by reference,

or when the simple type is lexically nested within the element declaration). The `dfdl:setVariable` statements on the sequence or choice within a global group definition execute before those on a group reference.

The `dfdl:newVariableInstance` annotations at any one annotation point of the schema are always executed in lexical order. However, `dfdl:newVariableInstance` annotations can also be found in different annotation points that are combined into the resolved set of annotations for one schema component. In this case, the order of execution of the `dfdl:newVariableInstance` statements from any one annotation point remains lexical. The order of execution of the `dfdl:newVariableInstance` annotations different annotation points follows the principle of innermost first, meaning that a schema component that contains or references another schema component has its `dfdl:newVariableInstance` statements executed after those of the contained or referenced schema component. For example, if a group reference and the sequence or choice group of a group definition both have `dfdl:newVariableInstance` statements, then those on the global group definition will execute before those on the group reference.

#### 9.5.1 Asserts and Discriminators with `testKind 'expression'`

Implementations are free to optimize by recognizing and executing discriminators or asserts with `testKind 'expression'` earlier so long as the resulting behavior is consistent with what results from the description above.

#### 9.5.2 Discriminators with `testKind 'expression'`

When parsing, an attempt to evaluate a discriminator must be made even if preceding statements or the parse of the schema component ended in a processing error.

This is because a discriminator's expression could evaluate to true thereby resolving a point of uncertainty even if the complete parsing of the construct ultimately caused a processing error.

Such discriminator evaluation has access to the DFDL Infoset of the attempted parse as it existed immediately before detecting the parse failure. Attempts to reference parts of the DFDL Infoset that do not exist are processing errors.

#### 9.5.3 Elements and `setVariable`

The resolved set of `dfdl:setVariable` statements for an element are executed **after** the parsing of the element. This is in contrast to the resolved set of `dfdl:setVariable` statements for a group which are executed **before** the parsing of the group. *yuck! :)*

For elements, this implies that these variables are set after the evaluation of expressions corresponding to any computed DFDL properties for that element, and so the variables may not be referenced from expressions that compute these DFDL properties.

That is, if an expression is used to provide the value of a property (such as `dfdl:terminator` or `dfdl:byteOrder`), the evaluation of that property expression occurs before any `dfdl:setVariable` annotation from the resolved set of annotations for that element are executed; hence, the expression providing the value of the property may not reference the variable. Schema authors can insert sequences to provide more precise control over when variables are set.

## 10. Core Representation Properties and their Format Semantics

The next sections specify the core set of DFDL v1.0 properties that may be used in DFDL annotations in DFDL Schemas to describe data formats.

It is a Schema Definition Error when a DFDL schema does *not* contain a definition for a representation property that is needed to interpret the data. For example, a DFDL schema containing any textual data must provide a definition of the character set encoding property (dfdl:encoding) for that textual data, and if it is not part of the format properties context for that data, then it is a Schema Definition Error.

Furthermore, no default values are provided for representation properties as built-in definitions by any DFDL processor. This requires DFDL schemas to be explicit about the representation properties of the data they describe, and avoids any possibility of DFDL schemas that are meaningful for some DFDL processors but not others.

The properties are organized as follows:

- Common to both Content and Framing (see 11)
- Common Framing, Position, and Length (see 12)
- Simple Type Content (see 13 )
- Sequence Groups (see 14 )
- Choice Groups (see 15 )
- Array elements and optional elements (see 16 )
- Calculated Values (see 17 )

*useful explicit structure*

Where properties are specific to a physical representation, the property name may choose to reflect this. Where properties are related to a specific logical type grouping (defined below), the property name may choose to reflect this.

A limited number of properties can take a DFDL expression which must return a value of the proper type for the property. Those properties that take an expression explicitly state in the description. Other properties do not take an expression.

The property description defines which schema component that the property may be specified on. In addition most DFDL properties may be specified on a dfdl:format annotation.

## 11. Properties Common to both Content and Framing

*— should you say when/if a property is required?*

Property Name	Description
byteOrder	<p>Enum or DFDL Expression</p> <p>Valid values 'bigEndian', 'littleEndian'.</p> <p>This property can be computed by way of an expression which returns the string 'bigEndian' or 'littleEndian'. The expression must not contain forward references to elements which have not yet been processed.</p> <p>Note that there is, intentionally, no such thing as 'native' endian<sup>14</sup>.</p> <p>This property applies to all Number, Calendar, and Boolean types with representation binary. Specifically that is binary integers, binary booleans, all packed decimals, binary floats, binary seconds and binary milliseconds.</p> <p>This property is never used to establish the byte order for text /strings with Unicode fixed-width encodings that do not specify the byte order (UTF-16 and UTF-32).</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
bitOrder	<p>Enum</p> <p>Valid values 'mostSignificantBitFirst', 'leastSignificantBitFirst'.</p> <p>The bits of a byte each have a place value or significance of <math>2^n</math>, for <math>n</math> from 0 to 7. Hence, the byte value <math>255 = 2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0</math>. A bit can always be unambiguously identified as the <math>2^n</math>-bit.</p> <p>The bit order is the correspondence of a bit's numeric significance to the bit position (1 to 8) within the byte.</p> <p>Value 'mostSignificantBitFirst' means:</p> <ul style="list-style-type: none"> <li>• The <math>2^7</math> bit is first, i.e., has bit position 1.</li> <li>• In general the <math>2^n</math> bit has position <math>8 - n</math>.</li> <li>• The least significant bits of byte N are considered to be adjacent to the most significant bits of byte N+1.</li> </ul> <p>Value 'leastSignificantBitFirst' means:</p> <ul style="list-style-type: none"> <li>• The <math>2^0</math> bit is first, i.e., has bit position 1.</li> <li>• In general the <math>2^n</math> bit has position <math>n + 1</math>.</li> <li>• The most significant bits of byte N are considered to be adjacent to the least significant bits of byte N+1.</li> </ul> <p>This property applies to all content and framing since it determines which bits of a byte occupy what bit positions. Content and framing are defined in terms of regions of the data stream, and these regions are defined in terms of the starting bit position and ending bit position; hence, dfdl:bitOrder is relevant to determining the specific bits of any grammar region (see Section 9.2) when the region's starting bit position or ending bit position are not on a byte boundary.</p>

*to be specified(?)*

**Deleted:** See Section 11.1 Unicode Byte Order Marks (BOM) for details.¶

<sup>14</sup> The concept of native-endian is avoided in DFDL since a DFDL schema containing such a property binding would not contain a complete description of data, but rather an incomplete one which would behave differently based on characteristics of the machine and implementation where the DFDL processor is executed. In DFDL this same behavior is achieved through the use of explicit parameterization using DFDL variables to set dfdl:byteOrder. See Section 7.7.2 Predefined Variables.

	<p>The bit order can only change on byte boundaries, and alignment of up to 7 bits will be <u>skipped (parsing) or inserted (unparsing)</u> to ensure byte-alignment whenever the bit order changes.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
encoding	<p>Enum or DFDL Expression</p> <p>Values are one of:</p> <ul style="list-style-type: none"> <li>IANA charset name</li> <li>CCSID<sup>15</sup></li> <li>DFDL standard encoding name</li> <li>Implementation-specific encoding name</li> </ul> <p>This property can be computed by way of an expression which returns an appropriate string value. The expression must not contain forward references to elements which have not yet been processed.</p> <p>Note that there is, deliberately, no concept of 'native' encoding<sup>16</sup>.</p> <p>Conforming DFDL v1.0 processors must accept at least 'UTF-8', 'UTF-16', 'UTF-16BE', 'UTF-16LE', 'ASCII', and 'ISO-8859-1' as encoding names.</p> <p>Encoding names are case-insensitive, so 'utf-8' and 'UTF-8' are equivalent.</p> <p>Unicode character set encodings that do not specify a byte order (such as UTF-16 or UTF-32) can have their byte-order controlled by a document-level byte-order-mark (BOM). See Section 11.1 <u>Unicode Byte Order Mark (BOM)</u> for details.</p> <p>The encoding name 'UTF-8' is interpreted strictly and does not include variants such as CESU-8.</p> <p>DFDL standard encoding names are defined in Section 34 Appendix D: DFDL Standard Encodings. When supported, a conforming DFDL implementation must implement them in a uniform manner so that they are portable across all DFDL implementations that implement them.</p> <p>Additional implementation-defined encoding names may be provided only for character set encodings for which there is no IANA name standard nor CCSID standard nor DFDL standard encoding. These implementation-defined encodings must have "X-" as a prefix to their name, as they are subject to being superseded by IANA or DFDL standard encoding names.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
utf16Width	<p>Enum</p> <p>Valid values are 'fixed', 'variable'.</p>

Deleted: inserted

← Feels like if you call it an enum you should list the values explicitly. Appendix?

← Does this apply to all property values? What about 'utf8', 'UTF8' are they valid?

Deleted: Unicode Byte Order Marks (BOM)

<sup>15</sup> CCSID stands for Coded Character Set ID, a decimal number syntax for a coded character set specifier. [CCSID].

<sup>16</sup> The concept of native character encoding is avoided in DFDL since a DFDL schema containing such a property binding would not contain a complete description of data, but rather an incomplete one which would behave differently based on characteristics of the operating environment where the DFDL processor executes. In DFDL this same behavior is achieved through the use of explicit parameterization using DFDL variables to set dfdl:encoding. See Section 7.7.2 Predefined Variables.



	<p>Applies only when encoding is 'UTF-16', 'UTF-16BE', 'UTF-16LE' or their CCSID equivalents.</p> <p>Specifies whether the encoding 'UTF-16' should be treated as a fixed or variable width encoding. 'UTF-16' can contain characters which require two codepoints (called a surrogate pair) to represent. When utf16Width is 'fixed', these surrogate code points are treated as separate characters. When utf16Width is 'variable', then surrogate pairs are converted into a single character on parsing, and such a character is split into two characters on unparsing.</p> <p>When utf16Width is 'variable', then on parsing an un-paired surrogate codepoint causes a decode error, which can be controlled via dfdl:encodingErrorPolicy described below.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
ignoreCase	<p>Enum</p> <p>Valid values are 'yes', 'no'.</p> <p>Whether mixed case data is accepted when matching delimiters and data values on input.</p> <p>This affects the behavior of matching for these properties: dfdl:initiator, dfdl:terminator, dfdl:separator, dfdl:nilValue, dfdl:textStandardExponentRep, dfdl:textStandardInfinityRep, dfdl:textStandardNaNRep, dfdl:textStandardZeroRep, dfdl:textBooleanTrueRep, and dfdl:textBooleanFalseRep.</p> <p>Property ignoreCase plays no part when comparing an element value with an XSDL enum facet, matching an element value to an XSDL pattern facet, or comparing an element value with the XSDL fixed property. It is therefore not used by validation (when validation is enabled), nor by the dfdl:checkConstraints function.</p> <p>On unparsing always use the delimiters or value as specified.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
encodingErrorPolicy	<p>Enum</p> <p>Valid values are 'error' or 'replace'.</p> <p>This property applies whenever dfdl:encoding is applicable.</p> <p>This property provides control of how decoding and encoding errors are handled when converting the data to text, or text to data. This includes converting when scanning for delimiters, matching regular expression length or test patterns, matching textual data type representation patterns against the data, and of course isolating the text content that will become the value of an element (parsing) or constructing the content from the value (unparsing).</p> <p>When parsing, an error can occur when decoding characters from their encoded form into the DFDL Infoset character set (ISO10646). This can occur due to invalid byte sequences, or not enough bytes found to make up the full encoding of a character.</p> <p>If 'replace', then the Unicode replacement character (U+FFFD) is substituted for the offending errors, one replacement character for any incorrect fragment of an encoding.</p> <p>If 'error' then a processing error occurs.</p>

	<p>When unparsing, the errors that can occur when encoding characters from Unicode/ISO 10646 into the specified encoding include when no mapping is provided by the encoding character set specification and when there is not enough space to output the entire encoding of the character (e.g., need 2 bytes for a 2-byte character codepoint, but only 1 byte remains in the available length.)</p> <p>If 'replace' then encoding-specific replacement/substitution character is output. It is a processing error if no such character is defined, and it is a processing error if there is any error when attempting to output the replacement (such as not enough room for the representation of the entire encoding of the replacement character).</p> <p>If 'error' then a processing error occurs.</p> <p>See Section 11.2 Character Encoding and Decoding Errors for further details.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
--	---

Table 12 Properties Common to both Content and Framing

11.1 Unicode Byte Order Mark (BOM)

DFDL does not provide any special treatment of Unicode Byte-Order Marks. They are treated as a Unicode ZWNBS character.

11.2 Character Encoding and Decoding Errors

When parsing, these are the errors that can occur when decoding characters into Unicode/ISO 10646.

1. The data is broken - invalid bit/byte sequences are found which do not match the definition of a character for the encoding.
2. Not enough data is found to make up the entire encoding of a character. That is, a fragment of a valid encoding is found.

When unparsing, these are the errors that can occur when encoding characters from Unicode/ISO 10646 into the specified encoding.

1. No mapping provided by the encoding specification.
2. Not enough room to output the entire encoding of the character (e.g., need 3 bytes for a character encoding that uses 3-bytes for that character, but only 1 byte remains in the available length.

The subsections below describe how these errors are handled.

11.2.1 Property dfdl:encodingErrorPolicy

The property dfdl:encodingErrorPolicy has two possible values: 'error' and 'replace'.

11.2.1.1 dfdl:encodingErrorPolicy 'error'

If 'error', then any error when decoding characters while parsing causes a processing error. For unparsing, any error when encoding characters causes a processing error.

When parsing, it does not matter if this happens when scanning for delimiters, matching a regular expression, matching a literal nil value, or constructing the value of a textual element.

There is one exception. When dfdl:lengthUnits is 'bytes', the 'not enough data' decoding error is ignored, and the data making up the fragment character is skipped over. Symmetrically, when unparsing the 'not enough room' encoding error is ignored and the left-over bytes are filled with the dfdl:fillByte.

Detection of character set decoding errors is often implementation-dependent because DFDL Implementations are free to optimize processing speed by skipping character decoding or encoding whenever possible. For example: when character set encodings are fixed-width, it is

**Deleted:** <#>provides automatic detection and generation of a Unicode BOM at the document level and saves (for parsing), or retrieves (for unparsing) the BOM information from the DFDL Infoset [unicodeByteOrderMark] member. ¶

<#>Parsing behaviour: When the dfdl:encoding property of the root element is specified, and is exactly one of UTF-8, UTF-16, or UTF-32 (or CCSID equivalents), then a DFDL parser will look for the appropriate BOM as the very first bytes in the data stream. ¶

<#>UTF-8. If a BOM is found<sup>17</sup> then this is used to set the document information item [unicodeByteOrderMark] member. If no BOM is found the parser takes no action. There is no need to model the BOM explicitly. ¶

<#>UTF-16. If a BOM is found then this is used to set the document information item [unicodeByteOrderMark] member, and all data with dfdl:encoding UTF-16 throughout the rest of the stream are assumed to have the implied byte order. If no BOM is found then all data with dfdl:encoding UTF-16 throughout the rest of the stream are assumed to have big-endian byte order. There is no need to model the BOM explicitly. ¶

<#>UTF-32. If a BOM is found then this is used to set the document information item [unicodeByteOrderMark] member, and all data with dfdl:encoding UTF-32 throughout the rest of the stream are assumed to have the implied byte order. If no BOM is found then all data with dfdl:encoding UTF-32 throughout the rest of the stream are assumed to have big-endian byte order. There is no need to model the BOM explicitly. ¶

<#>When the dfdl:encoding property of the root element is specified, and is exactly one of UTF-16LE, UTF-16BE, UTF-32LE or UTF-32BE (or CCSID equivalents), then a DFDL parser will not look for the appropriate BOM. The byte order to use is implicit in the encoding. If a BOM does appear at the start of the data stream, then it will simply be treated as a Unicode Zero-Width Non-Breaking Space (ZWNBS) character, because this shares the same codepoint as a BOM. ¶

<#>The dfdl:byteOrder property is never used to establish the byte order for Unicode encodings. ¶

<#>The parser never looks for a BOM at any other point in the data stream, so if a BOM appears elsewhere it will be treated as a Unicode ZWNBS character as described above<sup>18</sup>. ¶

<#>Unparsing behaviour: When the dfdl:encoding property of the root element is specified, and is exactly one of UTF-8, UTF-16 or UTF-32 (or CCSID equivalents), then a DFDL unparsing will look in the infoset document information item for a BOM. ¶

<#>UTF-8. If the document information item [unicodeByteOrderMark] member is 'UTF-8', the UTF-8 BOM is output as the very first bytes in the data stream. If the property is empty then no BOM is output. If the property has any other value, it is a processing error. There is no need to model the BOM explicitly. ¶

<#>UTF-16. If the document information item [unicodeByteOrderMark] member is 'UTF-16LE' or 'UTF-16BE', the corresponding UTF-16 BOM is output as the very first bytes in the data stream, and all data with dfdl:encoding UTF-16 throughout the rest of the document will be output with the implied byte order. If the property is empty then no BOM is output, and all data with dfdl:encoding UTF-16 throughout the rest of the document are assumed to have big-endian byte order. If the property has any other value, it is a processing error. There is no need to model the BOM explicitly. ¶

<#>UTF-32. If the document information item [unicodeByteOrderMark] member is 'UTF-32LE' or 'UTF-32BE', the corresponding UTF-32 BOM is output as the very first bytes in the data stream, and all data with dfdl:encoding UTF-32 throughout the rest of the document will be output with the implied byte order. If the property is empty then no BOM is output, and all data with dfdl:encoding UTF-32 throughout the rest of the document are assumed to have big-endian byte order. If the property has any other value, it is a processing error. There is no need to model the BOM explicitly. ¶

possible to determine lengths in bytes or bits from the length in characters by multiplying the length value by the character width, without having to decode any characters.

When parsing, character decoding errors MUST be detected when

- a) the decoding results in a character being placed into the DFDL Infoset
- b) the decoding is necessary to identify a delimiter
- c) the decoding is necessary to determine a match or non-match of a regular expression in a dfdl:assert or dfdl:discriminator with testKind='pattern'.

When unparsing, character encoding errors MUST be detected when

- a) an unmapped character appears in the infoset value of an element.

In all other cases, character set decoding and encoding errors MAY NOT be detected.

#### 11.2.1.2 dfdl:encodingErrorPolicy 'replace' for parsing

If 'replace' then any error when decoding characters results in the insertion of the Unicode Replacement Character (U+FFFD) as the replacement for that error.

It does not matter if this error and replacement happens when scanning for delimiters, matching a regular expression, matching a literal nil value, or constructing the value of a textual element.

There is one exception. When dfdl:lengthUnits is 'bytes', the 'not enough data' decoding error is ignored, no replacement character is created. The data making up the fragment character is skipped over. (It will be filled with the dfdl:fillByte when unparsing.)

Note that the "." wildcard in regular expressions will match the Unicode Replacement Character, so ".\*" and ".+" regular expressions can potentially cause very large matches (up to the entire data stream) to occur when data contains errors and dfdl:encodingErrorPolicy 'replace'. DFDL Schema authors are advised that bounded length negated regular expressions can help in this case. E.g., "[^\uFFFD]{0,50}" says to match any character (excluding the Unicode Replacement Character), but only up to length 50.

It is also worth noting that the Unicode Replacement Character can appear in data as an ordinary character, and this cannot be distinguished from the insertion of the Unicode Replacement Character due to a decoding error. This is likely to happen for data that is (a) initially parsed by a DFDL parser with dfdl:encodingErrorPolicy 'replace', and (b) which contains some decoding errors, but (c) is nevertheless successfully parsed, (d) is written back out to a file or other data repository, and (e) is parsed again. The written data will have replaced data errors with the Unicode Replacement Character, and so if the data is parsed again, it will no longer have errors, but will have the Unicode Replacement Character as a regular character in the data.

If dfdl:lengthUnits is 'characters', then a Unicode Replacement Character counts as contributing a single character to the length.

If the data contains more than one adjacent decode error, then the specific number of Unicode Replacement Characters that are inserted as the replacement of these errors is implementation-dependent. That is, some implementations may view, for example, three consecutive erroneous bytes as three separate decode errors, others may view them as a single or two decode errors. All implementations MUST, however, insert some number of Unicode Replacement Characters, and then continue to decode characters following the erroneous data.

The trimming of pad characters always happens after Unicode Replacement Characters have been inserted into the data.

#### 11.2.1.3 dfdl:encodingErrorPolicy 'replace' for unparsing

For unparsing, each encoding has a replacement/substitution character specified by the ICU. This character is substituted for the unmapped character or the character that has too large an encoding to fit in the available space.

There is one exception. When dfdl:lengthUnits is 'bytes', the 'not enough room' encoding error is ignored. The left-over bytes are filled with the dfdl:fillByte (they are skipped when parsing.)

The definitions of these substitution characters can be conveniently found for many encodings in the ICU Converter Explorer (<http://demo.icu-project.org/icu-bin/convexp>).

An encoding error is a processing error if the encoding does not provide a substitution/replacement character definition. (This would be rare, but could occur if a DFDL implementation allows many encodings beyond the minimum set.)

#### 11.2.2 Unicode UTF-16 Decoding/Encoding Non-Errors

The following specific situations involving encodings UTF-16, UTF-16LE, and UTF-16BE when `dfdl:utf16Width` "fixed", and they do not cause a decoding or encoding error.

- unpaired surrogate codepoint
- out-of-order surrogate codepoint pair
- surrogate codepoint pair is encountered

In all these cases the code-point(s) becomes a character code in the DFDL Information Item for the string.

#### 11.2.3 Preserving Data Containing Decoding Errors

There can be situations where data wants to be preserved exactly even if it contains errors.

It is suggested that if a DFDL schema author wants to preserve information containing data where the encodings have these kinds of errors, that they model such data as `xs:hexBinary`, or as `xs:string` but using an encoding such as `iso-8859-1` which preserves all bytes.

### 11.3 Byte Order and Bit Order

Byte order and bit order are separate concepts. However, of the possible combinations, only the following are allowed:

1. 'bigEndian' with 'mostSignificantBitFirst'
2. 'littleEndian' with 'mostSignificantBitFirst'
3. 'littleEndian' with 'leastSignificantBitFirst' <sup>19</sup>

Other combinations must produce Schema Definition Errors.

#### 11.4 dfdl:bitOrder Example

Consider a structure of 4 logical elements. The total length is 16 bits. Assume `dfdl:lengthUnits` is 'bits', `dfdl:representation` is 'binary', `dfdl:binaryNumberRep` is 'binary':

```
<element name="A" type="xs:int" dfdl:length="3"/> <!-- having value 3 -->
<element name="B" type="xs:int" dfdl:length="7"/> <!-- having value 9 -->
<element name="C" type="xs:int" dfdl:length="4"/> <!-- having value 5 -->
<element name="D" type="xs:int" dfdl:length="2"/> <!-- having value 1 -->
```

The above are colored so as to highlight the corresponding bits in the data below.

In a format where `dfdl:bitOrder` is 'mostSignificantBitFirst':

	01100010	01010101	
	AAABBBBB	BBCCCCDD	
Significance	M	L M	L
Bit Position	12345678	12345678	
Byte Position	----1----	----2---	

As presented here, the bits corresponding to each element appear left to right, and all bits for an individual element are adjacent. Within the bits of an individual element the most significant bit is on the left, least significant on the right, consistent with the way the bytes themselves are presented.

In contrast, in a format where `dfdl:bitOrder` is 'leastSignificantBitFirst':

	01001011	01010100	
	BBBBBAAA	DDCCCB	
Significance	M	L M	L
Bit Position	87654321	87654321	

*Very useful/interesting example.*

<sup>19</sup> Used by data format MIL-STD-2045

Byte Position ----1--- ----2---

In the above presentation note how the bits of the element 'B' do not appear adjacent to each other. The most significant bits of byte N are adjacent to the least significant bits of byte N+1.

11.4.1 Example Using Right-to-Left Display for 'leastSignificantBitFirst'

When working exclusively with data having dfdl:bitOrder 'leastSignificantBitFirst', it is useful to present data with bytes Right to Left. That is, with the bytes starting at byte 1 on the right, and increasing to the left.

	01010100	01001011
	DDCCCCBB	BBBBBAAA
Significance	M	L M L
Bit Position	87654321	87654321
Byte Position	----2---	----1---

With this reorientation, the bits of the element 'B' are once again displayed adjacently. Within the bits of an individual element the most significant bit is on the left, least significant on the right, consistent with the way the bytes themselves are presented.

Often the specification documents for data formats that with least-significant-bit-first bit order will describe data using this Right-to-Left presentation style.

11.4.2 dfdl:bitOrder and Grammar Regions

When any grammar region appears before (to the left of) or after (to the right of) another grammar region in the grammar rules of Section 9.2, and the boundary between the two falls within a byte rather than on a byte boundary, then the dfdl:bitOrder determines which bits are occupied by the regions.

In general, the notion of *before* means occupying lower-numbered bit positions, and the bit positions are numbered according to dfdl:bitOrder. Hence, when dfdl:bitOrder is 'mostSignificantBitFirst', grammar regions that are before, will occupy more-significant bits, and when dfdl:bitOrder is 'leastSignificantBitFirst', grammar regions that are before will occupy less-significant bits.

## 12. Framing

Several properties are common across the various framing styles or are used to distinguish them. Generally these have to do with position and length for text, bit fields, or opaque data.

### 12.1 Aligned Data *— Picture would be useful to illustrate how the properties interact*

Alignment properties control the leading alignment and trailing alignment regions.

When the alignment properties are applied to an array element, the properties are applied to each occurrence of the element; that is, not only to the first occurrence.

The following properties are used to define alignment rules.

Property Name	Description
alignment	<p>Non-negative Integer or 'implicit'</p> <p>A non-negative number that gives the alignment required for the beginning of the item. If alignment is needed then the size of the <b>AlignmentFill</b> grammar region will be non-zero if the item must be aligned to a boundary.</p> <p>'implicit' specifies that the natural alignment for the representation type is used. See the table of implicit alignments Table 14 Implicit Alignment in bits for simple elements. The 'implicit' alignment of a complex element is the alignment of its model group. The 'implicit' alignment of a model group is always 1. If alignment is 'implicit' then dfdl:alignmentUnits is ignored.</p> <p>For textual data, minimum alignment is mandated by the character-set encoding, and this property must be 'implicit' or set to a multiple of the character-set's mandatory alignment. See Section 12.1.2.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
alignmentUnits	<p>Enum</p> <p>Valid values are 'bits' or 'bytes'</p> <p>Scales the alignment so alignment can be specified in either units of bits or units of bytes.</p> <p>Only used when dfdl:alignment not 'implicit'</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
fillByte	<p>DFDL String Literal</p> <p>A single byte specified as a DFDL byte value entity or a single character. If a character is specified, it must be a single-byte character in the applicable encoding.</p> <p>Used on unparsing to fill empty space such as between two aligned elements.</p> <p>Used to fill these regions specified in the grammar: <b>RightFill</b>, <b>ElementUnused</b>, <b>ChoiceUnused</b>, <b>LeadingSkip</b>, <b>AlignmentFill</b>, and <b>TrailingSkip</b>.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
leadingSkip	<p>Non-negative Integer</p> <p>A non-negative number of bytes or bits, depending on dfdl:alignmentUnits, to skip before alignment is applied. Gives the size of the grammar region having the same name.</p>

	Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group
trailingSkip	<p>Non-negative Integer</p> <p>A non-negative number of bytes or bits, depending on dfdl:alignmentUnits, to skip after the element, but before considering the alignment of the next element. Gives the size of the grammar region having the same name.</p> <p>If dfdl:trailingSkip is specified when dfdl:lengthKind is 'delimited' then a dfdl:terminator must be specified.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>

**Table 13 Aligned Data Properties**

There are two properties which control the data alignment by controlling the length of the **AlignmentFill** region

- alignment - an integer 1 or greater
- alignmentUnits - bits or bytes

An element's representation is aligned to N units if P is the first position in the representation and  $P \bmod N = 1$ . When parsing, the position of the first unit of the data stream is 1.

For example, if dfdl:alignment is 4, and dfdl:alignmentUnits is 'bytes', then the element's representation must begin at 1 or 1 plus a multiple of 4 bytes. That is, 1, 5, 9, 13, 17 and so on.

The length of the **AlignmentFill** region is measured in bits. If alignmentUnits is 'bytes' then we multiply the alignment value by 8 to get the bit alignment. If the position in the data stream of the start of the **AlignmentFill** region is bit position N, then the length of the **AlignmentFill** region is the smallest non-negative integer L such that  $(L + N) \bmod B = 1$ . The position of the first bit of the aligned component is  $P = L + N$ .

The **LeadingSkip** and **TrailingSkip** regions length are controlled by two properties of corresponding names and the dfdl:alignmentUnits property.

#### 12.1.1 Implicit Alignment

When dfdl:alignment is 'implicit' the following alignment values are applied for each logical type.

Type	Alignment	
	text	binary

String	Encoding Specific (usually 8 bits, with exceptions: See Section 12.1.2)	Not applicable	
Float		32	
Double		64	
Decimal, Integer, nonNegativeInteger		Packed decimals: 8	binary: 8
Long, UnsignedLong			binary: 64
Int, UnsignedInt			binary: 32
Short, UnsignedShort			binary: 16
Byte, UnsignedByte			binary: 8
DateTime			binarySeconds: 32, binaryMilliseconds:64
Date			binarySeconds: 32, binaryMilliseconds:64
Time			binarySeconds: 32, binaryMilliseconds:64
Boolean		32	
HexBinary	Not applicable	8	

**Table 14 Implicit Alignment in bits**

Note: The above table specifies the implicit alignment in bits, but this does not imply that `dfdl:alignmentUnits 'bits'` can be specified for all simple types. Rather, `dfdl:alignmentUnits` and `dfdl:lengthUnits` are independent and have their own rules for when they are applicable.

#### 12.1.2 Mandatory Alignment for Textual Data

The term *textual data* is used to describe data of type `xs:string`, data with `dfdl:representation "text"`, as well as data being matched to delimiters (parsing) or output as delimiters (unparsing), and data being matched to regular expressions (parsing only - as in a `dfdl:assert` with `testKind 'pattern'`, or an element with `dfdl:lengthKind 'pattern'`).

Textual data has mandatory alignment that is character-set-encoding dependent. That is, these mandates come from the character set encoding specified by the `dfdl:encoding` property.

When processing textual data, it is a Schema Definition Error if the `dfdl:alignment` and `dfdl:alignmentUnits` properties are used to specify alignment that is not a multiple of the encoding-specified mandatory alignment.

If the data is not aligned to the proper boundary for the encoding when textual data is processed, then bits are skipped (parsing) or filled from `dfdl:fillByte` (unparsing) to achieve the mandatory alignment.

All required character set encodings in DFDL have 8-bit/1-byte alignment.

DFDL standard encodings specify their alignment. See Section 34 Appendix D: DFDL Standard Encodings.

Some implementations may include additional implementation-defined encodings which have other alignments.

Note the 16-bit and 32-bit Unicode character set encodings UTF-16, UTF-16BE, UTF-16LE, UTF-32, UTF-32BE, UTF-32LE, all have 8-bit/1-byte alignment.



### 12.1.3 Mandatory Alignment for Packed Decimal Data

→ not sure it's clear what this means in terms of properties

Packed decimal data must have a multiple of 4-bit alignment. It is a Schema Definition Error otherwise.

#### 12.1.4 Example: AlignmentFill

When `dfdl:alignmentUnits` is 'bits', and the `dfdl:alignment` is not a multiple of 8, then the `dfdl:bitOrder` property affects the alignment by controlling which bits are skipped as part of the grammar **AlignmentFill** region.

In general, the **AlignmentFill** region is *before* the regions it is aligning, and within a byte, the meaning of '*before*' is interpreted with respect to the `dfdl:bitOrder`.

When `dfdl:bitOrder` is 'mostSignificantBitFirst', then bits with more significance are before bits with less significance, so the **AlignmentFill** region occupies the most significant bits of the byte.

When `dfdl:bitOrder` is 'leastSignificantBitFirst', then bits with less significance are before bits with more significance, so the **AlignmentFill** region occupies the least significant bits of the byte.

Consider a structure of 2 logical elements. Assume `dfdl:lengthUnits='bits'`, `dfdl:representation='binary'`, `dfdl:binaryNumberRep='binary'` `dfdl:alignmentUnits='bits'`, and assume the data is at the beginning of the data stream.

```
<element name="A" type="xs:int" dfdl:length="2" dfdl:alignment='8' />
<!-- having value 1 -->
<element name="B" type="xs:int" dfdl:length="4" dfdl:alignment='4' />
<!-- having value 5 -->
```

The above are colorized so as to highlight the corresponding bits in the data below. The total length due to the alignment region appearing before element 'B' will be 8 bits.

In a format where `dfdl:bitOrder` is 'mostSignificantBitFirst' the data can be visualized as:

```
01000101
AAxxBBBB
Significance M L
Bit Position 12345678
```

In the above, the **AlignmentFill** region is marked with 'x' characters, and contains all 0 bit values.

In a format where `dfdl:bitOrder` is 'leastSignificantBitFirst' the presentation is different:

```
01010001
BBBBxxAA
Significance M L
Bit Position 87654321
```

In the above the **AlignmentFill** region still appears before element 'B', and in this case that is in less significant bits of the byte than the bits of content of element 'B', and these bits are displayed to the right of the bits of element 'B'.

## 12.2 Properties for Specifying Delimiters

The following properties apply to all objects that use text delimiters to delimit, that is, to initiate and/or terminate data. Delimiters can apply to binary data; however they are most often called 'text' delimiters because the concept is much more commonly used for textual data formats.

Property Name	Description
initiator	<p>List of DFDL String Literals or DFDL Expression</p> <p>Specifies a whitespace separated list of alternative DFDL String Literals one of which marks the beginning of the element or group of elements.</p> <p>This property can be computed by way of an expression which returns a string containing a whitespace separated list of DFDL String Literals. The expression must not contain forward references to elements which have not yet been processed.</p> <p>Each string literal in the list, whether apparent in the schema, or returned as the value of an expression, is restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed.</li> <li>• DFDL Byte Value entities ( <code>%#rXX;</code> ) are allowed.</li> <li>• DFDL Character Classes NL, WSP, WSP+, WSP*, and ES are allowed.</li> <li>• <b>If the ES entity or the WSP* entity appear alone as one of the string literals in the list, then <code>dfdl:initiatedContent</code> must be "no"</b>.</li> </ul> <p>If the above rules are not followed it is a Schema Definition Error.</p> <p>The <b>Initiator</b> region contains one of the initiator strings defined by <code>dfdl:initiator</code>.</p> <p>When parsing, the list of values is processed in a greedy manner, meaning it takes all the initiators, that is, each of the string literals in the whitespace separated list, and matches them each against the data. The initiator with the <b>longest match</b> is the one that is selected as having been 'found'. Once a matching initiator is found, no other matches will be subsequently attempted (ie, there is no backtracking).</p> <p>When an initiator is specified, it is a processing error if the component is required and one of the values is not found.</p> <p>If <code>dfdl:initiator</code> is "" (the empty string), then the <b>Initiator</b> region is of length zero, and no initiator is expected. It is not permitted for an expression to return an empty string. That is a Schema Definition Error.</p> <p>On unparsing the first initiator in the list is automatically inserted into the <b>Initiator</b> region.</p> <p>If <code>dfdl:ignoreCase</code> is 'yes' then the case of the string is ignored by the parser.</p> <p><b>Annotation:</b> <code>dfdl:element</code>, <code>dfdl:simpleType</code>, <code>dfdl:sequence</code>, <code>dfdl:choice</code>, <code>dfdl:group</code></p>

**Deleted:** `<#>ES` must not appear as the only DFDL string literal in the property. It can only appear as a member of a list.¶

This feels unclear:

- 1/ longest matching initiator
- 2/ longest matching data element following the initiator.

I assume (2) BUT why is this an appropriate criteria?

Why is this meaningful?  
Does it need to be supported?

I assume this means these are legal locations for this property.  
I cannot remember if you have said that explicitly.

terminator	<p>List of DFDL String Literals or DFDL Expression</p> <p>Specifies a whitespace separated list of alternative text strings that one of which marks the end of an element or group of elements. The strings <b>MUST</b> be searched for in the <b>longest first order</b>.</p> <p>This property can be computed by way of an expression which returns a string of whitespace separated list of values. The expression <del>must not</del> <sup>may</sup> contain forward references to elements which have not yet been processed.</p> <p>This property can be used to determine the length of an element as described in Section <a href="#">12.3.2</a> dfd:lengthKind 'delimited'.</p> <p>Each string literal in the list, whether apparent in the schema, or returned as the value of an expression, is restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed.</li> <li>• DFDL Byte Value entities ( %rXX; ) are allowed.</li> <li>• DFDL Character Classes NL, WSP, WSP+, WSP*, and ES are allowed.</li> <li>• Neither the ES entity nor the WSP* entity may appear on their own as one of the string literals in the list when the parser is determining the length of a component by scanning for delimiters.</li> </ul> <p>If the above rules are not followed it is a Schema Definition Error.</p> <p>The <b>Terminator</b> region contains the terminator string.</p> <p>If dfd:terminator is "" (the empty string), then the terminator region is of length zero, and no terminator is expected. It is not permitted for an expression to return an empty string, that is a Schema Definition Error.</p> <p>When parsing, the list of values is processed in a greedy manner, meaning it takes all the terminators, that is, each of the string literals in the whitespace separated list, and matches them each against the data. The terminator with the longest match is the one that is selected as having been 'found'. Once a matching terminator is found, no other matches will be subsequently attempted (ie, there is no backtracking).</p> <p>When a terminator is expected it is a processing error if no matching terminator is found. However, if dfd:documentFinalTerminatorCanBeMissing is specified then it is not an error if the last terminator in the data stream is not found.</p>
------------	---

Ah so it is (i) above.  
I think this is a better way to express.

**Deleted:** <#>ES must not appear as the only DFDL string literal in the property. It can only appear as a member of a list.¶

I don't understand

Again how does this make any sense?

Fuller explanation but seems repetitive/redundant with above. Something could be cut

*How is list order determined lexically?*

	<p>On unparsing the first terminator in the list is automatically inserted in the Terminator region.</p> <p>If dfdl:ignoreCase is 'yes' then the case of the string is ignored by the parser.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>
--	--

emptyValueDelimiterPolicy	<p>Enum</p> <p>Valid values are 'none', 'initiator', 'terminator' or 'both'</p> <p>Indicates that when an element in the data stream is empty, an initiator (if one is defined), a terminator (if one is defined), both an initiator and a terminator (if defined) or neither must be present.</p> <p>Ignored if both dfdl:initiator and dfdl:terminator are "" (empty string).</p> <p>'initiator' indicates that, on parsing, if the content region (which can be either the SimpleContent region or the ComplexContent region defined in Section 9.2) is empty then the dfdl:initiator must be present. It also indicates that on unparsing when the content region is empty that the dfdl:initiator will be output.</p> <p>'terminator' indicates that, on parsing, if the content region is empty then the dfdl:terminator must be present. It also indicates that on unparsing when the content region is empty the dfdl:terminator will be output.</p> <p>'both' indicates that, on parsing, if the content region is empty both the dfdl:initiator and dfdl:terminator must be present. On unparsing when the content region is empty the dfdl:initiator followed by the dfdl:terminator will be output.</p> <p>'none' indicates that if the content region is empty neither the dfdl:initiator or dfdl:terminator must be present. On unparsing when the content region is empty nothing will be output.</p> <p>It is a Schema Definition Error if dfdl:emptyValueDelimiterPolicy set to 'none' or 'terminator' when the parent group has dfdl:initiatedContent 'yes'.</p> <p>This property plays an important role in establishing empty representation. See 9.2.2 Empty Representation for details.</p> <p><u>This property is ignored if the element is fixed-length and length is not zero (as no empty representation is possible).</u></p> <p><u>The value of dfdl:emptyValueDelimiterPolicy should only be checked if there is a dfdl:initiator or dfdl:terminator in scope. If so, and dfdl:emptyValueDelimiterPolicy is not set, it is a Schema Definition Error. If dfdl:initiator is not "" and dfdl:terminator is "" and dfdl:emptyValueDelimiterPolicy is 'terminator' it is a Schema Definition Error. If dfdl:terminator is not "" and dfdl:initiator is "" and dfdl:emptyValueDelimiterPolicy is 'initiator' it is a Schema Definition Error. It is not an SDE if dfdl:emptyValueDelimiterPolicy is 'both' and one or both of dfdl:initiator and dfdl:terminator is "".</u> This is</p>
---------------------------	---

hard to parse.

Random thought:  
Is there a property like  
Decideability for DFDL Schema?  
It must be possible to define Schema that  
are highly ambiguous. Does DFDL  
guarantee to create one and only one  
input set for data conforming to such a  
schema?

Deleted: ¶

new paragraph  
full

	<p><u>to accommodate the common use of setting 'both' as a schema-wide setting.</u></p> <p><u>It is a Schema Definition Error if dfdl:emptyValueDelimiterPolicy is in effect and is set to 'none' or 'terminator' when the parent xs:sequence has dfdl:initiatedContent 'yes'.</u></p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
documentFinalTerminatorCanBeMissing	<p>Enum</p> <p>Valid values are 'yes', 'no'</p> <p>When the dfdl:documentFinalTerminatorCanBeMissing property is true, then when an element is the last element in the data stream, then on parsing, it is not an error if the terminator is not found.</p> <p>For example, if the data are in a file, and the format specifies lines terminated by the newline character (typically LF or CRLF), then if the last line is missing its newline, then this would normally be an error, but if dfdl:documentFinalTerminatorCanBeMissing is true, then this is not a processing error.</p> <p>On unparsing the terminator is always written out regardless of the state of this property.</p> <p>Annotation: dfdl:format (but applies to elements only)</p>
outputNewLine	<p>DFDL String Literal or DFDL Expression</p> <p>Specifies the character or characters that will be used to replace the %NL; character class entity during unparse</p> <p>It is a Schema Definition Error if any of the characters are not in the set of characters allowed by the DFDL entity %NL; Only individual characters or the %CR;%LF; combination are allowed.</p> <p>It is a Schema Definition Error if the DFDL entity %NL; is specified</p> <p>This property can be computed by way of an expression which returns a DFDL string literal. The expression must not contain forward references to elements which have not yet been processed.</p> <p>Annotation: dfdl:element, dfdl:simpleType, dfdl:sequence, dfdl:choice, dfdl:group</p>

→ Only used on unparse?  
How do we know when to create %NL during parse

**Table 15 Properties for Specifying Delimiters**

### 12.3 Properties for Specifying Lengths

These properties are used to determine the content length of an element and apply to elements of all types (simple and complex).

Property Name	Description
lengthKind	<p>Enum</p> <p>Controls how the content length of the component is determined.</p>

	<p>Valid values are: 'explicit', 'delimited', 'prefixed', 'implicit', 'pattern', 'endOfParent'</p> <p>A full description of each enumeration is given in the later sections.</p> <p>'explicit' means the length of the element is given by the dfdl:length property.</p> <p>'delimited' means the element length is determined by scanning for a terminator or separator.</p> <p>'prefixed' means the length of the element is given by an immediately preceding PrefixLength data region the format of which is specified using dfdl:prefixLengthType.</p> <p>'implicit' means the length is to be determined in terms of the type of the element and its schema-specified properties if any.</p> <p>'pattern' means the length of the element is given by scanning for a regular expression specified using the dfdl:lengthPattern property.</p> <p>'endOfParent' means that the length extends to the end of the containing (parent) construct.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
lengthUnits	<p>Enum</p> <p>Valid values 'bytes', 'characters', 'bits'.</p> <p>Specifies the units to be used whenever a length is being used to extract or write data. Applicable when dfdl:lengthKind is 'explicit', 'implicit' (for xs:string and xs:hexBinary) or 'prefixed'.</p> <p>Usage is restricted as follows:</p> <ul style="list-style-type: none"> <li>'characters' may only be used for complex elements and simple elements with text representation.</li> <li>'bits' may only be used for xs:boolean, xs:byte, xs:short, xs:int, xs:long, xs:unsignedByte, xs:unsignedShort, xs:unsignedInt, and xs:unsignedLong simple types with binary representation, <u>and for calendar simple types with binary packed representation.</u></li> <li>'bytes' must be used for type xs:hexBinary.</li> <li>'bytes' must be used for types xs:float and xs:double with binary representation.</li> </ul> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 16 Properties for Specifying Length****12.3.1 dfdl:lengthKind 'explicit'**

When dfdl:lengthKind is 'explicit' the length of the item is given by the dfdl:length property.

When the value of the dfdl:length property is a constant, it is used both when parsing and unparsing.

When unparsing an element with dfdl:lengthKind 'explicit' and where dfdl:length is an expression, then the data in the Infoset is treated as fixed-length and the dfdl:length property, whether literal constant or expression, is evaluated to provide the length to use.

When parsing and dfdl:lengthKind is 'explicit', delimiter scanning is turned off and in-scope delimiters are not looked for within or between elements.

Property Name	Description
---------------	-------------

→ Numbers?

Deleted: .

Can bytes be used for string? Rules do not prohibit as defined.

Deleted: variable length and not

Deleted: . The behaviour is the same as dfdl:lengthKind 'prefixed'. See Section 12.3.4.

length	<p>Non-negative Integer or DFDL Expression.</p> <p>Only used when lengthKind is 'explicit'.</p> <p>Specifies the length of this element in units that are specified by the dfdl:lengthUnits property.</p> <p>This property can be computed by way of an expression which returns a non-negative integer. The expression must not contain forward references to elements which have not yet been processed.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
--------	---

**Table 17 The dfdl:length Property**

When dfdl:lengthKind 'explicit', the method of extracting data is described in section: 12.3.7 Elements of Specified Length

#### 12.3.2 dfdl:lengthKind 'delimited'

On parsing, the length of an element with dfdl:lengthKind 'delimited' is determined by scanning the datastream for the delimiter.

The data stream is scanned for any of

- the element's terminator (if specified)
- an enclosing construct's separator or terminator
- the end of an enclosing element designated by its known length
- the end of the data stream

dfdl:lengthKind 'delimited' may be specified for

- elements of simple type with text representation
- elements of number or calendar simple type with dfdl:representation 'binary' that have a packed decimal representation
- elements of type xs:hexBinary
- elements of complex type.

The rules for resolving ambiguity between delimiters are:

1. When two delimiters have a common prefix, the longest delimiter is tried first.
2. When two delimiters have exactly the same length, but on different schema components, the innermost (most deeply nested) delimiter is tried first.
3. When the separator and terminator on a group have the same value, then at a point in the data where either the separator or terminator could be found, the separator is tried first. (Speculative execution may try the terminator subsequently).
4. If the length of the delimiters cannot be determined because character class entities ~~(which are variable length)~~ are being used then the delimiters must each be matched against the data, and the longest matching delimiter is taken as the match for the delimiter.
5. Ties (same matched length) are broken by giving a separator priority over a terminator of a sequence, or by choosing the innermost, or first in schema order.

When unparsing a simple element with text representation, the length in the data stream is the length of the content region, padded to dfdl:textOutputMinLength or the XSD minLength facet if dfdl:textPadKind is 'padChar'.

When unparsing a simple element with binary representation, then for hexBinary the length is the number of bytes in the Infoset value padded to the XSD minLength facet value using dfdl:fillByte, and for the other types the length is the minimum number of bytes to represent the value and any sign.

When unparsing a complex element, the length is that of the ComplexContent region.



### 12.3.2.1 Non-Delimited Elements within Delimited Constructs

When a simple or complex element has a specified length, `dfdl:lengthKind` 'pattern', or `dfdl:lengthKind` 'endOfParent' then delimiter scanning is suspended for the duration of the processing of that element.

This allows formats to be parsed which are delimited, but have nested elements which contain non-character data so long as that nested data can be isolated from the delimited data context surrounding it.

### 12.3.2.2 Delimited Binary Data

Formats involving binary data, most notably packed decimals, can use delimiter scanning but care must be taken that the delimiters cannot match data represented in these formats. In particular, the delimiters must be chosen with knowledge that BCD data can contain any byte both of whose nibbles are 0 to 9 (that is, excluding A to F). Packed data adds bytes with a sign indicator, that is, a nibble in the range A to F.

General binary data can contain any bit pattern whatsoever, so delimiter scanning for numbers and calendars with `dfdl:representation` 'binary' are disallowed, with the specific exception of packed decimals. Delimiter scanning is also allowed for type `xs:hexBinary`.

Implementation Note: Scanning for delimiters when data is binary, or when using byte-value (aka raw byte) entities in delimiters, means that a simple character based delimiter scanner is not sufficient, as the delimiter may not be representable as characters.

### 12.3.3 `dfdl:lengthKind` 'implicit'

When `dfdl:lengthKind` is 'implicit', the length is determined in terms of the type of the element and its schema-specified properties.

For complex elements, 'implicit' means the length is determined by the combined lengths of the contained children, that is the ComplexValue region, and the ElementUnused region is of size 0. However, note that alignment regions inside the contained children within the ComplexValue region may be of different lengths depending on the ComplexValue's starting position alignment.

For simple elements the length is fixed and is given in [Table 18 Length in Bits for SimpleTypes](#) when `dfdl:lengthKind` is 'implicit'.

Type	Length	
	text	binary
String	The XSD <code>maxLength</code> facet gives length in characters, but this is also the length in bytes. (See note below: character set encoding must be single-byte.) Multiply by 8 to get number of bits.	Not applicable
Float	Not allowed	32 bits
Double	Not allowed	64 bits
Decimal, Integer, nonNegativeInteger	Not allowed	packed decimal: Not allowed binary: Not allowed
Long, UnsignedLong	Not allowed	binary: 64 bits
Int, UnsignedInt	Not allowed	binary: 32 bits

*I don't think you have discussed which properties are required/optional in which locations.*

**Deleted:** For complex elements, 'implicit' means the length is determined by the combined lengths of the contained children, that is the ComplexContent region. However, note that alignment regions inside the ComplexContent region may be of different lengths depending on the ComplexContent's starting position alignment.¶

**Deleted:** Table 18 Length in Bits for SimpleTypes when `dfdl:lengthKind` is 'implicit'

Short, UnsignedShort	Not allowed		binary: 16 bits
Byte, UnsignedByte	Not allowed		binary: 8 bits
DateTime	Not allowed		binarySeconds: 32 bits, binaryMilliseconds: 64 bits.
Date	Not allowed		binarySeconds: <del>Not allowed</del> , binaryMilliseconds: <del>Not allowed</del> ,
Time	Not allowed		binarySeconds: <del>Not allowed</del> , binaryMilliseconds: <del>Not allowed</del> ,
Boolean	Length of longest of dfdl:textBooleanTrueRep and dfdl:textBooleanFalseRep values	32 bits	
HexBinary	Not applicable	The XSD maxLength facet gives the length in bytes. Multiply by 8 to convert to number of bits.	

**Deleted:** 32 bits**Deleted:** 64 bits**Deleted:** 32 bits**Deleted:** 64 bits**Table 18 Length in Bits for SimpleTypes when dfdl:lengthKind is 'implicit'**

- 'Not Allowed' means that there is no implicit length for the combination of simple type and representation and it is a Schema Definition Error if dfdl:lengthKind 'implicit' is specified.
- packed decimal means dfdl:binaryNumberRep is 'packed', 'bcd', or 'ibm4690Packed'
- binary means dfdl:binaryNumberRep is 'binary'
- binarySeconds means dfdl:binaryCalendarRep is 'binarySeconds'
- binaryMilliseconds means dfdl:binaryCalendarRep is 'binaryMilliseconds'.

When dfdl:lengthKind is 'implicit', the method of extracting data is described in section: 12.3.7 Elements of Specified Length.

It is a Schema Definition Error if type is xs:string and dfdl:lengthKind is 'implicit' and dfdl:lengthUnits is 'bytes' and encoding is not an SBCS (exactly 1 byte per character code) encoding. This prevents a scenario where validation against the XSD maxLength facet is in characters but parsing and unparsing using the XSD maxLength facet is in bytes.

#### 12.3.4 dfdl:lengthKind 'prefixed'

When dfdl:lengthKind is 'prefixed' the length of the element is given by the integer value of the PrefixLength region specified using dfdl:prefixLengthType. The property dfdl:prefixIncludesPrefixLength also can be used to adjust the length appropriately.

When dfdl:lengthKind is 'prefixed' the method of extracting data is described in section: 12.3.7 Elements of Specified Length

When dfdl:lengthKind is 'prefixed', delimiter scanning is turned off and in-scope delimiters are not looked for within or between elements.

Property Name	Description
---------------	-------------

prefixIncludesPrefixLength	<p>Enum</p> <p>Valid values are 'yes', 'no'</p> <p>Whether the length given by a prefix includes the length of the prefix as well as the length of the content region (which can be either the SimpleContent region or the ComplexContent region defined in Section 9.2 DFDL Data Syntax Grammar.)</p> <p>Used only when dfdl:lengthKind 'prefixed'.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
prefixLengthType	<p>QName</p> <p>Name of a simple type derived from xs:integer or any subtype of it.</p> <p>This type, with its DFDL annotations specifies the representation of the length prefix, which is in the PrefixLength region.</p> <p>It is a Schema Definition Error if the xs:simpleType specifies any of:</p> <ul style="list-style-type: none"> <li>• dfdl:lengthKind 'delimited', 'endOfParent', or 'pattern'</li> <li>• dfdl:lengthKind 'explicit' where length is an expression</li> <li>• dfdl:outputValueCalc</li> <li>• dfdl:initiator or dfdl:terminator other than empty string</li> <li>• dfdl:alignment other than '1'</li> <li>• dfdl:leadingSkip or dfdl:trailingSkip other than '0'.</li> </ul> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 19 Properties for dfdl:lengthKind 'prefixed'**

The representation of the element is in two parts.

1. The 'prefix length' is an integer which specifies the length of the element's content. The representation of the length prefix is described by a simple type which is identified using the dfdl:prefixLengthType property.
2. The content of the element.

When parsing, the length of the element's content is obtained by parsing the simple type specified by dfdl:prefixLengthType to obtain an integer value. Note that all required properties must be present on the specified simple type or defaulted because there is no element declaration to supply any missing required properties.

If the dfdl:prefixIncludesPrefixLength property is 'yes' then the length of the element's content is the value of the prefix length minus the length of the content of the prefix length.

If the prefix type is dfdl:lengthKind 'implicit' or 'explicit' then the dfdl:lengthUnits properties of both the prefix type and the element must be the same.

The DFDL properties that specify the format of the prefix come from annotations directly on the dfdl:prefixLengthType's type definition, and from the default format annotation for the schema document containing the definition of that type. If the using-element resides in a separate schema, the simple type does not pick up values from the element's schema's default dfdl:format annotation.

When unparsing, the length of the element's content region must be determined first as described below. Then the value of the prefix length must be adjusted using dfdl:prefixIncludesPrefixLength. *can be adjusted using based on the value of*

Then the prefix length can be written to the data stream using the properties on the dfdl:prefixLengthType, and finally the element's content can be written to the data stream.

Consider this example:

```
<xs:element name="myString" type="xs:string"
```

```

        dfdl:lengthKind="prefixed"
        dfdl:prefixIncludesPrefixLength="no"
        dfdl:prefixLengthType="packed3"/>

<xs:simpleType name="packed3"
  dfdl:representation="binary"
  dfdl:binaryNumberRep="packed"
  dfdl:lengthKind="explicit"
  dfdl:length="2" >
  <xs:restriction base="integer" />
</xs:simpleType>
```

In the above, the string has a prefix length of type 'packed3' containing 3 packed decimal digits. The property dfdl:prefixIncludesPrefixLength is an enumeration which allows the length computation to be varied to include or exclude the length of the prefix element itself. The prefix length's value contains the length measured in units given by dfdl:lengthUnits. When parsing, if the dfdl:lengthUnits are bits, then any number of bits can be in the representation. However, the same is not true when unparsing. The DFDL Infoset does not store the number of bits in a number, so the number of bits will always be a multiple of 8 bits. When unparsing, the value of the prefix is computed automatically by obtaining the length of the element's content.

For a simple element with text representation, the length is computed as for dfdl:lengthKind 'delimited'.

For a simple element with binary representation, the length is given in the table below.

For a complex element, the length is that of the ComplexContent region.

Type	Length	
String	Not applicable	
Float	32	
Double	64	
Decimal, Integer, NonNegativeInteger	Compute the minimum number of bytes to represent the value (per dfdl:binaryNumberRep) and sign (if applicable). Multiply by 8 for number of bits.	
Long, UnsignedLong	packed decimal: as Decimal	binary: 64
Int, UnsignedInt		binary: 32
Short, UnsignedShort		binary: 16
Byte, UnsignedByte		binary: 8
DateTime		binarySeconds: 32, binaryMilliseconds: 64
Date		binarySeconds: <del>Not allowed</del> , binaryMilliseconds: <del>Not allowed</del>
Time		binarySeconds: <del>Not allowed</del> , binaryMilliseconds: <del>Not allowed</del>
Boolean	32	
HexBinary	Compute the number of bytes in the infoset value padded to the value of the XSD minLength facet (which gives minimum length in bytes) using dfdl:fillByte if necessary. This gives the unparse length in bytes. Multiply by 8 for the number of bits.	

Deleted: 32

Deleted: 64

Deleted: 32

Deleted: 64

**Table 20 Unparse Lengths (in Bits) for Binary Data with dfdl:lengthKind 'prefixed'****12.3.4.1 Nested Prefix Lengths<sup>20</sup>**

It is possible for a prefix length, as specified by dfdl:prefixLengthType, to itself have a prefix length

It is a Schema Definition Error if this nesting exceeds 1 deep. That is, an element can have a prefix length, which defines a PrefixLength region (see Section 9.2 DFDL Data Syntax Grammar). The PrefixLength region can itself have a type which also specifies a prefix length, thereby defining a PrefixPrefixLength region. It is a Schema Definition Error unless the type associated with the PrefixPrefixLength is different from the type associated with the PrefixLength.

**12.3.5 dfdl:lengthKind 'pattern'**

The dfdl:lengthKind 'pattern' means the length of the element is given by a regular expression specified using the dfdl:lengthPattern property. The DFDL processor scans the data stream to determine a string value that is the match to a regular expression. The pattern is only used on parsing.

When dfdl:lengthKind is 'pattern', delimiter scanning is turned off and in-scope delimiters are not looked for within or between elements.

Property Name	Description
lengthPattern	<p>DFDL Regular Expression.</p> <p>Only used when lengthKind is 'pattern'.</p> <p>Specifies a regular expression that, on parsing, is executed against the datastream to determine the length of the element.</p> <p>The data stream beginning at the starting offset of the content region (which can be either the SimpleContent region or the ComplexContent region defined in Section 9.2 DFDL Data Syntax Grammar) of the element is interpreted as a stream of characters in the encoding of the element, and the regular expression contained in the dfdl:lengthPattern property is executed against that stream of characters. When the element is complex the encoding used is the dfdl:encoding of the complex element itself.</p> <p>It is a Schema Definition Error if there is no value for the dfdl:encoding property in scope.</p> <p>DFDL Escape Schemes (per dfdl:escapeSchemeRef) are not used when executing the regular expression.</p> <p>If the pattern matching of the regular expression reads data that cannot be decoded into characters of the current encoding, then the behavior is controlled by the dfdl:encodingErrorPolicy property. See dfdl:encodingErrorPolicy in Section 11 Properties Common to both Content and Framing.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 21 The dfdl:lengthPattern Property**

On unparsing the behavior is the same as for dfdl:lengthKind 'prefixed'.

When the DFDL regular expression is matched against data:

- The data is considered to be text in the character set encoding specified by the dfdl:encoding property, regardless of the actual representation of the element.

<sup>20</sup>This feature allows DFDL to describe the needed "one more level" of prefix that is needed for modeling an ASN.1 format, but without the complexities of general recursion.

- The data is decoded from the specified encoding into Unicode before the actual matching takes place.
- If there is no match (ie, the length of the data found to match the pattern is zero) it is not a processing error but instead it means the length is zero.

### 12.3.6 dfdl:lengthKind 'endOfParent'

The dfdl:lengthKind 'endOfParent' means that the element is terminated either by the end of the data stream, or the end of an enclosing complex element with dfdl:lengthKind 'explicit', 'pattern', 'prefixed' or 'endOfParent', or the end of an enclosing choice with dfdl:choiceLengthKind 'explicit'. The 'parent' element or choice does not have to be the immediate enclosing component of the element, but there must be no other components defined between the element specifying dfdl:lengthKind 'endOfParent' and the end of the parent.

A convenient way of describing the parent is as a 'box', being defined as a portion of the data stream that has an established content length prior to the parsing of its children. If the parent is such a 'box' then the element specifying dfdl:lengthKind 'endOfParent' is the last element in the 'box' and its content extends to the end of the 'box'.

A dfdl:lengthKind of 'endOfParent' can only be used on simple and complex elements in the following locations:

- When the immediate containing model group is a sequence, on the final element in the sequence
- When the immediate containing model group is a choice, on any element that is a branch of the choice
- A simple type or global element declaration referenced by one of the above.
- A global element declaration that is the document root.

It is a Schema Definition Error if:

- the element has a terminator.
- the element has dfdl:trailingSkip not equal to 0.
- the element has maxOccurs > 1.
- any other model-group is defined between this element and the end of the enclosing component.
- any other represented element is defined between this element and the end of the enclosing component.
- the parent is an element with dfdl:lengthKind 'implicit' or 'delimited'.
- the element has text representation, does not have a single-byte character set encoding, and the *effective length units* of the parent is not 'characters'.

The effective length units of the parent are:

- dfdl:lengthUnits if parent is an element with dfdl:lengthKind 'explicit' or 'prefixed';
- 'characters' if parent is an element with dfdl:lengthKind 'pattern';
- 'bytes' if parent is a choice with dfdl:choiceLengthKind 'explicit';
- 'characters' if the element is the document root;
- the effective length units of the parent's parent if parent is an element with dfdl:lengthKind 'endOfParent'

If the element is in a sequence then it is a Schema Definition Error if:

- the dfdl:separatorPosition of the sequence is 'postfix'
- the dfdl:sequenceKind of the sequence is not 'ordered'
- the sequence has a terminator
- there are floating elements in the sequence

- the sequence has a non-zero dfdl:trailingSkip

If the element is in a choice where dfdl:choiceLengthKind is 'implicit' then it is a Schema Definition Error if:

- the choice has a terminator
- the choice has a non-zero dfdl:trailingSkip

A simple element must have one of:

- type xs:string
- dfdl:representation 'text'
- type xs:hexBinary
- dfdl:representation 'binary' and a packed decimal representation

A complex element can have dfdl:lengthKind 'endOfParent'. If so then its last child element can be any dfdl:lengthKind including 'endOfParent'.

The dfdl:lengthKind 'endOfParent' can also be used on the document root to allow the last element to consume the data up to the end of the data stream.

The use of dfdl:lengthKind 'endOfParent' is distinct from the situation where the length of the last element in the parent is known but is not sufficient to fill the parent. In the latter case the remaining data are ignored on parsing and filled with dfdl:fillByte on unparsing.

When parsing an element with dfdl:lengthKind 'endOfParent', delimiter scanning is turned off and in-scope terminating delimiters are not looked for within the element.

When unparsing an element with dfdl:lengthKind 'endOfParent', if the parent is a complex element with dfdl:lengthKind 'explicit' where dfdl:length is not an expression, or a choice with dfdl:choiceLengthKind 'explicit', then the element with dfdl:lengthKind 'endOfParent' is padded or filled in the usual manner to the required length, by completing the **LeftPadding**, **RightPadOrFill**, **ElementUnused**, or **ChoiceUnused** regions of the data as appropriate.

### 12.3.7 Elements of Specified Length

An element has a specified length when dfdl:lengthKind is 'explicit', 'implicit' (simple type only) or 'prefixed'. The units that the length represents are specified by the dfdl:lengthUnits property except where noted in Section 12.3.3.

Using specified length, it is possible for an element to have content length longer than needed to represent just the data value. For example, a simple text element may be padded in the **RightPadding** region if the data is not long enough.

When an element has specified length, but appears inside a complex type element having delimited length kind, delimiter scanning is turned off and in-scope delimiters are not looked for within or between elements.

An element of specified length with dfdl:lengthKind 'implicit' or 'explicit' where dfdl:length is not an expression has a known length when unparsing.

An element of specified length with dfdl:lengthKind 'prefixed' is considered to have a variable length when unparsing. Specifically, the processor automatically determines the value to store in the prefix, based on the length of the info set element, and the properties which modify the interpretation of the prefix length value, such as dfdl:prefixIncludesPrefixLength.

For dfdl:lengthKind 'explicit' (expression), whether parsing or unparsing the expression is evaluated to obtain the length. When unparsing, the processor cannot automatically determine in what way the length information is to be stored as it comes from an expression that may access one or more elements and perform any calculation. Hence, normally the value of the element or elements involved in the length calculation would be computed using dfdl:outputValueCalc, using an expression that measures the length of the element by way of functions such as dfdl:contentLength or dfdl:valueLength.

*I don't think I know what these are yet. bold face?*

*on these are tokens from the formal grammar you defined. It would be useful to identify these as such*

**Deleted:** However, an

**Deleted:** or 'explicit' where dfdl:length is an expression

**Commented [MJB2]:** This language was moved from below. There is no change to the behavior for prefixed unparsing.

**Deleted:** For

**Deleted:**

**Deleted:** . N

**Deleted:** the

**Deleted:** using

When parsing, if the data stream ends without enough data to parse an element, that is, N bits are needed based on the dfdl:length, but only M < N bits are available, then it is a processing error.

If dfdl:lengthUnits is 'characters' then the length (in bits) of the content region (i.e., SimpleContent or ComplexContent defined in Section 9.2 DFDL Data Syntax Grammar) will depend on the encoding of the characters.

- If the dfdl:encoding property specifies a fixed-width encoding then the content length is the character width (in bits) multiplied by the length.
- If the dfdl:encoding property specifies a variable-width encoding then the length will depend on the actual characters in the element's value. The characters must be decoded one by one, adding up their widths (in bits), while counting up to the specified length value.

For a simple element, dfdl:lengthUnits 'characters' may only be used for textual elements, it is a Schema Definition Error otherwise.

Some DFDL implementations may support character set encodings where the characters are not a multiple of 8-bits wide. Encodings which are 5, 6, 7, and 9 bits wide are rare, but do exist, so the overall length of the content region may not be a multiple of 8-bits wide.

#### 12.3.7.1 Length of Simple Elements with Textual Representation

*Textual data* is defined to mean either data of type string or data where the dfdl:representation property is 'text'.

For a textual element, the dfdl:lengthUnits property can be either 'bytes' or 'characters'.

##### 12.3.7.1.1 Text Length Specified in Bytes

If a textual element has dfdl:lengthUnits of 'bytes', and the dfdl:encoding is not SBCS, then it is possible for a partial character encoding to appear after the code units of the characters. In this case, the following rules apply:

- When parsing, as many characters as possible are extracted from the bytes of the simple content region. Any left over bytes are skipped. (They are considered part of the grammar *RightFill* region).
- When unparsing, if the simple content region is larger than the encoded length of the element (as padded when dfdl:textPadKind is not 'none') then the remaining bytes, which are insufficient to hold another character code, are filled with dfdl:fillByte (Again, this is the grammar *RightFill* region.)

It is a Schema Definition Error if type is xs:string and dfdl:textPadKind is not 'none' and dfdl:lengthUnits is 'bytes' and dfdl:encoding is not an SBCS encoding and the XSD minLength facet is not zero. This prevents a scenario where validation against the XSD minLength facet is in characters but padding would be performed in bytes.

##### 12.3.7.2 Length of Simple Elements with Binary Representation

This section discusses the dfdl:lengthKind 'explicit' and 'prefixed' specified lengths for the different binary representations. When dfdl:lengthKind is 'implicit', see Section 12.3.3 dfdl:lengthKind 'implicit'.

The dfdl:lengthUnits can be 'bytes' or 'bits' unless otherwise stated. It is Schema Definition Error if dfdl:lengthUnits is 'characters'.

It is a Schema Definition Error if the specified dfdl:length for an element of dfdl:lengthKind 'explicit' is a string literal integer such that the length of the data exceeds the capacity of the simple type.

It is a processing error if the specified length for an element of dfdl:lengthKind 'prefixed' or 'explicit' (with dfdl:length an expression) is an integer such that the length of the data exceeds the capacity of the simple type.

**Deleted:** <#>For dfdl:lengthKind 'prefixed' the processor automatically determines the value to store in the prefix, based on the length of the infoset element, and the properties which modify the interpretation of the prefix length value, such as dfdl:prefixIncludesPrefixLength.¶¶



12.3.7.2.1 Length of Base-2 Binary Number Elements

Non-floating point numbers with binary representation and dfdl:binaryNumberRep 'binary' are represented as a bit string which contains a base-2 representation.  
The value of the specified length is constrained per the table below. The lengths are expressed in bits and are inclusive.

Type	Minimum value of length	Maximum value of length
xs:byte	2	8
xs:short	2	16
xs:int	2	32
xs:long	2	64
xs:unsignedByte	1	8
xs:unsignedShort	1	16
xs:unsignedInt	1	32
xs:unsignedLong	1	64
xs:nonNegativeInteger	1	Implementation-dependent (but not less than 64)
xs:integer	2	Implementation-dependent (but not less than 64)
xs:decimal	8 <sup>21</sup>	Implementation-dependent (but not less than 64)

Deleted: 2

Table 22: Allowable Specified Lengths in Bits for Base-2 Binary Number Elements

See Section 13.7.1.1 Converting Base-2 Binary Numbers for details of the conversion to/from numeric values.

12.3.7.2.2 Length of Floating Point Binary Number Elements

For binary elements of types xs:float or xs:double, a specified length must be either exactly 4 bytes or exactly 8 bytes respectively.  
The dfdl:lengthUnits property must be 'bytes'. It is a Schema Definition Error otherwise.  
See Section 13.8 Properties Specific to Float/Double with Binary Representation.

12.3.7.2.3 Length of Packed Decimal Number Elements

Non-floating point numbers with binary representation and dfdl:binaryNumberRep 'packed', 'bcd', or 'ibm4690Packed', are represented as a bit string of 4 bit nibbles. The term *packed decimal* is used to describe such numbers.  
It is a Schema Definition Error if the specified length is not a multiple of 4 bits.  
The maximum specified length of a packed decimal number is implementation-defined.  
See Section 13.7 [Properties Specific to Number with Binary Representation](#), for details of the conversion of the packed decimal bit string to/from a numeric value.

Deleted: Properties Specific to Numbers with Binary Representation...

12.3.7.2.4 Length of Binary Boolean Elements

<sup>21</sup> Type decimal must be a minimum of 8 bits because lengthUnits 'bits' is not allowed for xs:decimal.

The specified length of a binary element of type `xs:boolean` is as for type `xs:unsignedInt` described in 12.3.7.2.1 [Length of Base-2 Binary Number Elements](#).

See also Section 13.10 Properties Specific to Boolean with Binary Representation for details of how the data is converted to/from a Boolean value.

**Deleted:** Length of Base-2 Binary Number Elements

#### 12.3.7.2.5 Length of Base-2 Binary Calendar Elements

Calendars with binary representation and `dfdl:binaryCalendarRep` 'binarySeconds' or 'binaryMilliseconds' are represented as a bit string which contains a base-2 representation. The specified length must be either exactly 4 bytes or exactly 8 bytes respectively.

The `dfdl:lengthUnits` property must be 'bytes'. It is a Schema Definition Error otherwise.

See Section 13.13 Properties Specific to Calendar with Binary Representation for details of how the data is converted to/from the calendar type.

#### 12.3.7.2.6 Length of Packed Decimal Calendar Elements

Calendars with binary representation and `dfdl:binaryCalendarRep` 'packed', 'bcd', or 'ibm4690Packed', are represented as a bit string of 4 bit nibbles. The term *packed decimal* is used to describe such calendars.

It is a Schema Definition Error if the specified length is not a multiple of 4 bits.

The maximum specified length of a packed decimal calendar is implementation-~~defined~~ (but not less than 9 bytes, which corresponds to calendar pattern 'yyyyMMddhhmmssSSS')<sup>22</sup>.

See Section 13.13 Properties Specific to Calendar with Binary Representation for details of how the data is converted to/from the calendar type.

**Deleted:** dependent

#### 12.3.7.2.7 Length of Binary Opaque Elements

The `dfdl:lengthUnits` property must be 'bytes'. It is a Schema Definition Error otherwise.

When unparsing a specified length element of type `xs:hexBinary`, and the simple content region is larger than the length of the element in the Infoset, then the remaining bytes are filled using the `dfdl:fillByte` property.

The `dfdl:fillByte` is **not** used to trim an element of type `xs:hexBinary` when parsing.

#### 12.3.7.3 Length of Complex Elements

A complex element of specified length is defining a 'box' in which its child elements exist. An example of this would be a fixed-length record element with a variable number of children elements. The `dfdl:lengthUnits` may be 'bytes' or 'characters' and it is a Schema Definition Error otherwise.

It is possible that the children may not entirely fill the full length of the complex element. An example is a complex element with a specified length of 100 characters, which contains a sequence of child elements that use up less than 100 characters of data, perhaps because an optional element is not present. In this case the remaining unused data is called the `ElementUnused` region in the data syntax grammar of section 9.2. Another example is a complex element with a specified length of 100 bytes, which contains a sequence of child elements the last of which has `dfdl:lengthKind` 'endOfParent', `dfdl:representation` 'text' and a multi-byte `dfdl:encoding` such that the element does not use up all the bytes of data. In this case the remaining unused bytes comprise the child element's *RightFill* region in the data syntax grammar of section 9.2. In both examples, the unused area is skipped when parsing, and is filled with the `dfdl:fillByte` on unparsing.

<sup>22</sup> This is the smallest pattern that contains all the digit-only symbols. SSS is the minimum precision that must be supported for fractional seconds, but in can be more, hence why 'not less than 9 bytes'.

Note that a poorly chosen value for `dfdl:fillByte` may fill the region with data that cannot be decoded in the character set encoding, resulting in a decode error when this data is subsequently parsed again. When `dfdl:lengthUnits` is 'characters' the value for `dfdl:fillByte` should be chosen so as to avoid this error.

### 13. Simple Types

The 'representation' property identifies the physical representation of the element. The DFDL logical types are grouped to illustrate which physical representations apply to each logical type. These properties provide the correct interpretation of the data found in the SimpleContent grammar region.

The allowable physical representations for each logical type grouping are also shown, where the logical type groupings are defined as:

Logical Type Group	Types
Number	xs:double, xs:float, xs:decimal, xs:integer and its restrictions (xs:int, xs:unsignedLong, etc.)
String	xs:string
Calendar	xs:dateTime, xs:date, xs:time
Opaque	xs:hexBinary
Boolean	xs:boolean

Table 23 Logical type groups

#### 13.1 Properties Common to All Simple Types

Property Name	Description
representation	Enum Valid values are dependent on logical type. <b>Number:</b> 'text', 'binary' <b>String:</b> representation is assumed to be 'text' and the dfdl:representation property is not examined <b>Calendar:</b> 'text', 'binary' <b>Boolean:</b> 'text', 'binary' <b>Opaque:</b> representation is assumed to be 'binary' and the dfdl:representation property is not examined. Annotation: dfdl:element, dfdl:simpleType

Table 24 Properties Common to All Simple Types

The permitted representation properties for each logical type are shown in Table 25: Logical Type to Representation properties

Logical type	dfdl:representation	Additional representation property
String	Assumed to be text	
Float, Double	text	<b>dfdl:textNumberRep:</b> standard
	binary	<b>dfdl:binaryFloatRep:</b> ieee, ibm390Hex
Decimal, Integer, nonNegativeInteger	text	<b>dfdl:textNumberRep:</b> standard, zoned
	binary	<b>dfdl:binaryNumberRep:</b> packed, bcd, ibm4690Packed, binary

Long, Int, Short, Byte, UnsignedLong, UnsignedInt, UnsignedShort, UnsignedByte	text	<b>dfdl:textNumberRep:</b> standard, zoned
	binary	<b>dfdl:binaryNumberRep:</b> packed, bcd, ibm4690Packed, binary
DateTime, Date, Time	text	
	binary	<b>dfdl:binaryCalendarRep:</b> packed, bcd, ibm4690Packed, binarySeconds, binaryMilliseconds
Boolean	text	
	binary	
HexBinary	Assumed to be binary	

Table 25: Logical Type to Representation properties

## 13.2 Properties Common to All Simple Types with Text representation

Property Name	Description
textPadKind	<p>Enum</p> <p>Valid values 'none', 'padChar'.</p> <p>Indicates whether to pad the data value on unparsing. This controls the contents of the <b>LeftPadding</b> and <b>RightPadding</b> regions of the data syntax grammar in section 9.2.</p> <p>'none': No padding occurs. When dfdl:lengthKind is 'implicit' or 'explicit' (and dfdl:length is not an expression) the unparsed data value must match the expected length otherwise it is a processing error.</p> <p>'padChar': The data value is padded using the dfdl:textStringPadCharacter, dfdl:textNumberPadCharacter, dfdl:textBooleanPadCharacter or dfdl:textCalendarPadCharacter depending on the type of the element. The padding characters populate the <b>LeftPadding</b> and/or <b>RightPadding</b> regions depending on <b>dfdl:textStringJustification</b>, <b>dfdl:textNumberJustification</b>, <b>dfdl:textBooleanJustification</b>, or dfdl:textCalendarJustification, depending on the type of the element.</p> <p>When dfdl:lengthKind is 'implicit' the data value is padded to the implicit length for the type.</p> <p>When dfdl:lengthKind is 'explicit' (and dfdl:length is not an expression) the data value is padded to the length given by the dfdl:length property.</p> <p>When dfdl:lengthKind is 'explicit' (and dfdl:length is an expression), 'delimited', 'prefixed', 'pattern' the data value is padded to the length given by the XSD minLength facet for type 'xs:string' or dfdl:textOutputMinLength property for other types.</p> <p>When dfdl:lengthKind is 'endOfParent' the data value is padded to the available length.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textTrimKind	<p>Enum</p> <p>Valid values 'none', 'padChar'</p>

*why are these properties not discussed in this table?*

*I thought it was because they are only relevant when textPadKind = 'padChar' but that is also true of textOutputMinLength which is in this table.*

	<p>Indicates whether to trim data on parsing. This controls the expected contents of the <b>LeftPadding</b> and <b>RightPadding</b> regions of the data syntax grammar in section 9.2.</p> <p>When 'none' no trimming takes place.</p> <p>When 'padChar' the element is trimmed of the dfdl:textStringPadCharacter, dfdl:textNumberPadCharacter, dfdl:textBooleanPadCharacter or dfdl:textCalendarPadCharacter depending on the type of the element. The padding characters populate the <b>LeftPadding</b> and/or <b>RightPadding</b> regions depending on dfdl:textStringJustification, dfdl:textNumberJustification, or dfdl:textCalendarJustification, depending on the type of the element.</p> <p>Annotation: dfdl:element , dfdl:simpleType</p>
textOutputMinLength	<p>Non-negative Integer.</p> <p>Only used when dfdl:textPadKind is 'padChar' and dfdl:lengthKind is 'delimited', 'prefixed', 'pattern', 'explicit' (when dfdl:length is an expression) or 'endOfParent', and type is not xs:string</p> <p>Specifies the minimum content length during unparsing for simple types that do not allow the XSD minLength facet to be specified.</p> <p>For dfdl:lengthKind 'delimited', 'pattern' and 'endOfParent' the length units are always characters, for other dfdl:lengthKinds the length units are specified by the dfdl:lengthUnits property.</p> <p>If dfdl:textOutputMinLength is zero or less than the length of the representation text then no padding occurs.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
escapeSchemeRef	<p>QName or empty String</p> <p>The name of the dfdl:defineEscapeScheme annotation that provides the additional properties used to describe the escape scheme. If the value is the empty string then escaping is explicitly turned off.</p> <p>See: Section 7.6 The dfdl:escapeScheme Annotation Element, and Section 7.5 The dfdl:defineEscapeScheme Defining Annotation Element.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 26 Properties Common to All Simple Types with Text Representation****13.2.1 The dfdl:escapeScheme Properties**

The dfdl:escapeScheme annotation is used within a dfdl:defineEscapeScheme annotation to group the properties of an escape scheme and allows a common set of properties to be defined that can be reused.

An escape scheme is needed when the content of a text element contains sequences of characters that are the same as an in-scope separator or terminator. If the characters are not escaped, a parser scanning for a separator or terminator would erroneously find the character sequence in the content.

An escape scheme defines the properties that describe the text escaping rules. There are two variants on such schemes:

- The use of a single escape character to cause the next character to be interpreted literally. The escape character itself is escaped by the escape escape character.
- The use of a pair of escape strings to cause the enclosed group of characters to be interpreted literally. The ending escape string is escaped by the escape escape character.

On parsing, the escape scheme is applied after pad characters are trimmed and on unparsing before pad characters are added. A pad character is not escaped by an escape character. When

parsing, pad characters are trimmed without reference to an escape scheme. When unparsing, pad characters are added without reference to an escape scheme.

On unparsing, the application of escape scheme processing takes place before the application of the dfdl:emptyValueDelimiterPolicy property.

Property Name	Description
escapeKind	<p>Enum</p> <p>Valid values 'escapeCharacter', 'escapeBlock'</p> <p>The type of escape mechanism defined in the escape scheme</p> <p>When 'escapeCharacter': On unparsing a single character of the data is escaped by adding an dfdl:escapeCharacter <u>or dfdl:escapeEscapeCharacter immediately</u> before it. <u>The characters to escape are determined by property dfdl:escapeCharacterPolicy.</u></p> <p>On parsing any in-scope terminating delimiter encountered in the data is not interpreted as such when it is immediately preceded by the dfdl:escapeCharacter (when not itself preceded by the dfdl:escapeEscapeCharacter). Occurrences of the dfdl:escapeCharacter and dfdl:escapeEscapeCharacter are removed from the data <u>as determined by property dfdl:escapeCharacterPolicy</u>, unless the dfdl:escapeCharacter is preceded by the dfdl:escapeEscapeCharacter, or the dfdl:escapeEscapeCharacter does not precede the dfdl:escapeCharacter, <u>respectively.</u></p> <p>When 'escapeBlock': On unparsing the entire data are escaped by adding dfdl:escapeBlockStart to the beginning and dfdl:escapeBlockEnd to the end of the data. The data is either always escaped or escaped when needed as specified by dfdl:generateEscapeBlock. If the data is escaped and contains the dfdl:escapeBlockEnd then first character of each appearance of the dfdl:escapeBlockEnd is escaped by the dfdl:escapeEscapeCharacter.</p> <p>On parsing the dfdl:escapeBlockStart string must be the first characters in the (trimmed) data in order to activate the escape scheme. The dfdl:escapeBlockStart string is removed from the beginning of the data. Until a matching dfdl:escapeBlockEnd string (that is, one not preceded by the dfdl:escapeEscapeCharacter) is found in the data, any in-scope terminating delimiter encountered in the data is not interpreted as such, and any dfdl:escapeEscapeCharacters are removed when they precede an dfdl:escapeBlockEnd string. The matching dfdl:escapeBlockEnd string is removed from the data.. The matching dfdl:escapeBlockEnd does not have to be the last character(s) in the (trimmed) data in order to de-activate the escape scheme. A dfdl:escapeBlockStart occurring anywhere in the data other than the first characters has no significance.</p> <p>Annotation: dfdl:escapeScheme</p>
escapeCharacter	<p>DFDL String Literal or DFDL Expression</p> <p>Specifies one character that escapes the subsequent character.</p> <p>Used when dfdl:escapeKind is 'escapeCharacter'</p> <p>It is a Schema Definition Error if dfdl:escapeCharacter is empty when dfdl:escapeKind is 'escapeCharacter'</p>

**Deleted:** The following are escaped if they are in the data¶  
Any in-scope terminating delimiter by escaping its first character. ¶  
dfdl:escapeCharacter (escaped by dfdl:escapeEscapeCharacter)¶  
any dfdl:extraEscapedCharacters

	<p>This property can be computed by way of an expression which returns a character. The expression must not contain forward references to elements which have not yet been processed.</p> <p><i>Escape and Quoting Character Restrictions:</i> The string literal is restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed</li> <li>• The DFDL byte value entity ( <code>%#rXX;</code> ) is not allowed</li> <li>• DFDL Character classes NL, WSP, WSP+, WSP*, and ES are not allowed</li> </ul> <p>It is a Schema Definition Error if the string literal contains any of the disallowed constructs.</p> <p>Escape characters contribute to the <a href="#">simple value region (SimpleLogicalValue or NilLiteralValue)</a> of the field</p> <p>Annotation: <code>dfdl:escapeScheme</code></p>
escapeBlockStart	<p>DFDL String Literal</p> <p>The string of characters that denotes the beginning of a sequence of characters escaped by a pair of escape strings.</p> <p>Used when <code>dfdl:escapeKind</code> is 'escapeBlock'</p> <p>It is a Schema Definition Error if <code>dfdl:escapeBlockStart</code> is empty when <code>dfdl:escapeKind</code> is 'escapeBlock'</p> <p>The string literal value is restricted in the same way as described in "Escape and Quoting Character Restrictions" in the description of the <code>dfdl:escapeCharacter</code> property.</p> <p>An <code>dfdl:escapeBlockStart</code> string contributes to the <a href="#">simple value region (SimpleLogicalValue or NilLiteralValue)</a> of the field</p> <p>Annotation: <code>dfdl:escapeScheme</code></p>
escapeBlockEnd	<p>DFDL String Literal</p> <p>The string of characters that denotes the end of a sequence of characters escaped by a pair of escape strings.</p> <p>Used when <code>dfdl:escapeKind</code> is 'escapeBlock'.</p> <p>It is a Schema Definition Error if <code>dfdl:escapeBlockEnd</code> is empty when <code>dfdl:escapeKind</code> is 'escapeBlock'.</p> <p><u>When parsing, it is a processing error if the end of the data for the element is reached and the escapeBlockEnd is not found in the data.</u></p> <p>The string literal value is restricted in the same way as described in "Escape and Quoting Character Restrictions" in the description of the <code>escapeCharacter</code> property.</p> <p>A <code>dfdl:escapeBlockEnd</code> string contributes to the <a href="#">simple value region (SimpleLogicalValue or NilLiteralValue)</a> of the field</p> <p>Annotation: <code>dfdl:escapeScheme</code></p>
escapeEscapeCharacter	<p>DFDL String Literal or DFDL Expression</p> <p>Specifies one character that escapes an immediately following <code>dfdl:escapeCharacter</code> or first character of <code>dfdl:escapeBlockEnd</code>.</p> <p>Used when <code>dfdl:escapeKind</code> is 'escapeCharacter' or 'escapeBlock'.</p> <p>This property can be computed by way of an expression which returns a character. The expression must not contain forward references to elements which have not yet been processed.</p>

Deleted: content

Deleted: length

Deleted: content length

Deleted: content length



	<p>The string literal value is restricted in the same way as described in "Escape and Quoting Character Restrictions" in the description of the escapeCharacter property.</p> <p>If the empty string is specified then no escaping of escape characters occurs.</p> <p>It is explicitly allowed for both the dfdl:escapeCharacter and the dfdl:escapeEscapeCharacter to be the same character. In that case processing functions as if the dfdl:escapeCharacter escapes itself.</p> <p><a href="#">Escape-escape characters contribute to the simple value region (SimpleLogicalValue or NilLiteralValue) of the field.</a></p> <p>Annotation: dfdl:escapeScheme</p>
extraEscapedCharacters	<p>List of DFDL String Literals</p> <p>A whitespace separated list of single characters that must be escaped in addition to the in-scope delimiters. If there are no extra characters to escape the property should be set to "".</p> <p>The string literal values are restricted in the same way as described in "Escape and Quoting Character Restrictions" in the description of the dfdl:escapeCharacter property.</p> <p>This property only applies on unparsing.</p> <p><a href="#">Extra escaped characters contribute to the simple value region (SimpleLogicalValue or NilLiteralValue) of the field.</a></p> <p>Annotation: dfdl:escapeScheme</p>
generateEscapeBlock	<p>Enum</p> <p>Valid values 'always', 'whenNeeded'</p> <p>Controls when escaping is used on unparsing when dfdl:escapeKind is 'escapeBlock'.</p> <p>If 'always' then escaping is always occurs as described in dfdl:escapeKind.</p> <p>If 'whenNeeded' then escaping occurs as described in dfdl:escapeKind when the data contains any of the following:</p> <ul style="list-style-type: none"> <li>any in-scope terminating delimiter</li> <li>dfdl:escapeBlockStart at the start of the data</li> <li>any dfdl:extraEscapedCharacters</li> </ul> <p>Annotation: dfdl:escapeScheme</p>
<a href="#">escapeCharacterPolicy</a>	<p><a href="#">Enum</a></p> <p><a href="#">Valid values are 'all', 'delimiters'.</a></p> <p><a href="#">Controls when escape characters are removed during parsing, and output during unparsing, when dfdl:escapeKind is 'escapeCharacter'.</a></p> <p><a href="#">When 'all':</a></p> <p><a href="#">During unparsing the following are escaped as described in dfdl:escapeKind when they are in the data.</a></p> <ul style="list-style-type: none"> <li><a href="#">Any in-scope terminating delimiter by escaping its first character.</a></li> </ul>

	<ul style="list-style-type: none"><li>• <u>dfdl:escapeCharacter (escaped by dfdl:escapeEscapeCharacter)</u></li><li>• <u>any dfdl:extraEscapedCharacters</u></li></ul> <p><u>During parsing, occurrences of dfdl:escapeCharacter and dfdl:escapeEscapeCharacter are interpreted and removed from the data as described in dfdl:escapeKind.</u></p> <p><u>When 'delimiters':</u></p> <p><u>During unparsing the following are escaped as described in dfdl:escapeKind when they are in the data.</u></p> <ul style="list-style-type: none"><li>• <u>Any in-scope terminating delimiter by escaping its first character.</u></li><li>• <u>dfdl:escapeCharacter (escaped by dfdl:escapeEscapeCharacter)</u></li></ul> <p><u>During parsing, occurrences of dfdl:escapeCharacter and dfdl:escapeEscapeCharacter are interpreted and removed from the data as described in dfdl:escapeKind, except that dfdl:escapeCharacter is only removed when it immediately precedes an in-scope terminating delimiter.</u></p> <p><u>Annotation: dfdl:escapeScheme</u></p>
--	---

Table 27 Escape Scheme Properties

13.2.1.1 Escape Scheme Examples

Consider a dfdl:escapeScheme annotation with the following properties:

- dfdl:escapeBlockStart="start"
- dfdl:escapeBlockEnd="end"
- dfdl:escapeEscapeCharacter="#"

If this is used to serialize a DFDL Infoset element of type xs:string with value "A hash is a #", then the value is wrapped with the dfdl:escapeBlockStart and dfdl:escapeBlockEnd, giving simple content "startA hash is a #end". If this data was parsed, the "#end" will be treated as an escaped escape block end and the parse will fail, reporting that there is no escape block end in the data.

In this scenario, the data is not compliant with the escape scheme, and the DFDL serializer must issue a processing error.

*More examples?  
This one is interesting  
but limited*

### 13.3 Properties for Bidirectional support for All Simple Types with Text representation

Bidirectional text is a feature expected in a future revision of the DFDL standard.

Property name	Description
textBidi	Enum Valid <u>value is, 'no'</u> <u>This property exists in anticipation of future DFDL features that will enable bidirectional text processing.</u> Annotation: dfdl:element, dfdl:simpleType (representation text)

**Deleted:** consists of mainly right-to-left text with some left-to-right nested segments (such as an Arabic text with some information in English), or vice versa (such as an English letter with a Hebrew address nested within it.)  
Note: the bidirectional properties apply to the content of the element and not to the initiator, terminator or separator if defined....

**Deleted:** values are 'yes'

**Deleted:** Indicates the text content of the element is bidirectional....

**Deleted:** textBidiOrdering

**Table 28 Properties for Bidirectional support for All Simple Types with Text representation**

#### 13.4 Properties Specific to String

Property Name	Description
textStringJustification	Enum Valid values 'left', 'right', 'center' Unparsing: 'left': Justifies to the left and adds padding chars to the string contents if the string is too short, to the length determined by the dfdl:textPadKind property. 'right': Justifies to the right and adds padding chars to the string contents if the string is too short, to the length determined by the dfdl:textPadKind property. 'center': Adds equal padding chars left and right of the string contents if the string is too short, to the length determined by the dfdl:textPadKind property. It adds one extra padding char on the left if needed. Parsing: 'left': Trims any pad characters from the right of the string, according to dfdl:textTrimKind property. 'right': Trims any pad characters from the left of the string, according to dfdl:textTrimKind property. 'center' Trims any pad characters from the left and right of the string, according to dfdl:textTrimKind property. Annotation: dfdl:element, dfdl:simpleType
textStringPadCharacter	DFDL String Literal The value that is used when padding or trimming string elements. The value can be a single character or a single byte. If a character, then it can be specified using a literal character or using DFDL entities. If a byte, then it must be specified using a single byte value entity otherwise it is a Schema Definition Error If a pad character is specified when dfdl:lengthUnits is 'bytes' then the pad character must be a single-byte character. If a pad byte is specified when dfdl:lengthUnits is 'characters' then <ul style="list-style-type: none"> <li>the encoding must be a fixed-width encoding</li> </ul>

	<ul style="list-style-type: none"> <li>padding and trimming must be applied using a sequence of N pad bytes, where N is the width of a character in the fixed-width encoding.</li> </ul> <p><i>Padding Character Restrictions:</i> The string literal is restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>DFDL character entities are allowed</li> <li>The DFDL byte value entity ( <code>%#rXX;</code> ) is allowed.</li> <li>DFDL Character classes NL, WSP, WSP+, WSP*, and ES are not allowed</li> </ul> <p>It is a Schema Definition Error if the string literal contains any of the disallowed syntax.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
<truncatespecifiedlengthstring </truncatespecifiedlengthstring  <i>why needed? (seems a little arbitrary)</i>	<p>Enum</p> <p>Valid values are 'yes', 'no'</p> <p>Used on unparsing only</p> <p>'yes' means if the logical type is <code>xs:string</code> and the value is longer than the specified length, the string is truncated to this length. (See section 12.3.7 Elements of Specified Length.) No processing error is raised.</p> <p>The position from which data is truncated is determined by the value of the <code>dfdl:textStringJustification</code> property. If the value of the <code>dfdl:textStringJustification</code> property is 'left', data is truncated from the right; if the value of the <code>dfdl:textStringJustification</code> property is 'right', data is truncated from the left. However if the value of the <code>dfdl:textStringJustification</code> property is 'center', truncation does not occur and a processing error occurs if the value is too long.</p> <p>When unparsing, validation errors cannot be prevented by truncation as validation takes place on the augmented infoset, before any truncation has occurred.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>

Table 29 13.4 Properties Specific to String

**13.5 Properties Specific to Number with Text or Binary Representation**

Property Name	Description
decimalSigned	<p>Enum</p> <p>Valid values are 'yes', 'no'</p> <p>Indicates whether an xs:decimal element is signed. See 13.6.2 Converting logical numbers to/from text representation and 13.7.1 Converting Logical Numbers to/from Binary to see how this affects the presence of the sign in the data stream.</p> <p>'yes' means that the xs:decimal element is signed</p> <p>'no' means that the xs:decimal element is not signed</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 30 Properties Specific to Number with Text or Binary Representation****13.6 Properties Specific to Number with Text Representation**

Property Name	Description
textNumberRep	<p>Enum</p> <p>Valid values are 'standard', 'zoned'</p> <p>'standard' means represented as characters in the character set encoding specified by the dfdl:encoding property.</p> <p>'zoned' means represented as a zoned decimal in the character set encoding specified by the dfdl:encoding property. Zoned is not supported for float and double numbers. Base 10 is assumed, and the encoding must be for an EBCDIC or ASCII compatible encoding. It is a Schema Definition Error if any of these requirements are not met.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textNumberJustification	<p>Enum</p> <p>Valid values 'left', 'right', 'center'</p> <p>Controls how the data is padded or trimmed on parsing and unparsing.</p> <p>Behavior as for dfdl:textStringJustification.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textNumberPadCharacter	<p>DFDL String Literal</p> <p>The value that is used when padding or trimming number elements.</p> <p>The value can be a single character or a single byte.</p> <p>If a character, then it can be specified using a literal character or using DFDL entities.</p> <p>If a byte, then it must be specified using a single byte value entity</p> <p>If a pad character is specified when dfdl:lengthUnits is 'bytes' then the pad character must be a single-byte character.</p> <p>If a pad byte is specified when dfdl:lengthUnits is 'characters' then</p>

	<ul style="list-style-type: none"> <li>the encoding must be a fixed-width encoding</li> <li>padding and trimming must be applied using a sequence of N pad bytes, where N is the width of a character in the fixed-width encoding.</li> </ul> <p><u>When parsing, if the pad character is '0' and dfdl:textTrimKind is 'padChar' then the SimpleContent region is trimmed of the '0' characters as defined by the trimming rules. If at least one '0' character is removed and the trimmed text causes a processing error when parsed, a single '0' character is re-instated and the text is parsed again. This is to handle the case when '0' characters are trimmed away leaving no digits. This rule also applies when the pad character is a DFDL character entity equivalent to '0'. This rule does not apply when the pad character is any other character nor when a pad byte is specified.</u></p> <p>The string literal value is restricted in the same way as described in "Pad Character Restrictions" in the description of the dfdl:textStringPadCharacter property.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textNumberPattern	<p>String</p> <p>Defines the ICU-like pattern that describes the format of the text number. The pattern defines where grouping separators, decimal separators, implied decimal points, exponents, positive signs and negative signs appear. It permits definition by either digits/fractions or significant digits. Allows rounding.</p> <p>When dfdl:textNumberRep is 'standard' this property only applies when dfdl:textStandardBase is 10. When dfdl:textNumberRep is 'standard' and dfdl:textStandardBase is not 10 the number is represented as the minimum number of characters to represent the digits. There is no sign or virtual decimal point.</p> <p>The syntax of dfdl:textNumberPattern is described in section 13.6.1 The dfdl:textNumberPattern Property</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textNumberRounding	<p>Enum</p> <p>Specifies how rounding is controlled during unparsing.</p> <p>Valid values 'pattern', 'explicit'</p> <p>When dfdl:textNumberRep is 'standard' this property only applies when dfdl:textStandardBase is 10.</p> <p>If 'pattern' then rounding takes place according to the pattern. A rounding increment may be specified in the dfdl:textNumberPattern using digits '1' through '9', otherwise rounding is to the width of the pattern. The rounding mode is always 'roundHalfEven'.</p> <p>If 'explicit' then the rounding increment is specified by the dfdl:textNumberRoundingIncrement property, and any digits '1' through '9' in the dfdl:textNumberPattern are</p>

**Deleted:** When parsing, if the pad character is '0' and the SimpleContent region consists entirely of '0' characters, then the last remaining '0' is not trimmed and a single '0' is the result of the trimming. This rule also applies when the pad character is a DFDL character entity equivalent to '0'. This rule does not apply when the pad character is any other character nor when a pad byte is specified.

<p><i>What exactly is this?</i></p>	<p>treated as digit '0'. The rounding mode is specified by the <code>dfdl:textRoundingMode</code> property.</p> <p>To disable rounding, use 'explicit' in conjunction with 'roundUnnecessary' for the <code>dfdl:textNumberRoundingMode</code>. If rounding is disabled then any <b>excess precision</b> is treated as a processing error.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
<p><i>With mathematical definitions of these</i></p>	<p>Enum</p> <p>Specifies how rounding occurs during unparsing, when <code>dfdl:textNumberRounding</code> is 'explicit'.</p> <p>When <code>dfdl:textNumberRep</code> is 'standard' this property only applies when <code>dfdl:textStandardBase</code> is 10.</p> <p>To switch off rounding, use 'roundUnnecessary'.</p> <p>Valid values 'roundCeiling', 'roundFloor', 'roundDown', 'roundUp', 'roundHalfEven', 'roundHalfDown', 'roundHalfUp', 'roundUnnecessary'</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
<p><code>textNumberRoundingIncrement</code></p>	<p>Double</p> <p>Specifies the rounding increment to use during unparsing, when <code>dfdl:textNumberRounding</code> is 'explicit'.</p> <p>When <code>dfdl:textNumberRep</code> is 'standard' this property only applies when <code>dfdl:textStandardBase</code> is 10.</p> <p>A negative value is a Schema Definition Error.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
<p><code>textNumberCheckPolicy</code></p>	<p>Enum</p> <p>Values are 'strict' and 'lax'.</p> <p>Indicates how lenient to be when parsing against the pattern.</p> <p>When <code>dfdl:textNumberRep</code> is 'standard' this property only applies when <code>dfdl:textStandardBase</code> is 10.</p> <p>If 'lax' and <code>dfdl:textNumberRep</code> is 'standard' then grouping separators are ignored, leading and trailing whitespace is ignored, leading zeros are ignored and quoted characters may be omitted.</p> <p>If 'lax' and <code>dfdl:textNumberRep</code> is 'zoned' then positive punched data is accepted when parsing an unsigned type, and unpunched data is accepted when parsing a signed type</p> <p>If 'strict' and <code>dfdl:textNumberRep</code> is 'standard' then the data must follow the pattern with the exceptions that digits 0-9, decimal separator and exponent separator are always recognised and parsed.</p> <p>If 'strict' and <code>dfdl:textNumberRep</code> is 'zoned' then the data must follow the pattern.</p> <p>On unparsing the pattern is always followed and follow the rules in 13.6.2 Converting logical numbers to/from text representation.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>

textStandardDecimalSeparator	<p>List of DFDL String Literals or DFDL Expression</p> <p>Defines the a whitespace separated list of single characters that will appear (individually) in the data as the decimal separator.</p> <p>This property is applicable, when dfdl:textNumberRep is 'standard' and dfdl:textStandardBase is 10. It must be set if dfdl:textNumberPattern contains a decimal separator symbol ("."), or the E or @ symbols. (it is a Schema Definition Error otherwise.) Empty string is not an allowable value.</p> <p>This property can be computed by way of an expression which returns a character. The expression must not contain forward references to elements which have not yet been processed.</p> <p><i>Text Number Character Restrictions:</i> The the string literal is restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed</li> <li>• The DFDL byte value entity ( %rXX; ) is not allowed.</li> <li>• DFDL Character classes NL, WSP, WSP+, WSP*, and ES are not allowed</li> </ul> <p>It is a Schema Definition Error if the string literal contains any of the disallowed syntax constructs.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textStandardGroupingSeparator	<p>DFDL String Literal or DFDL Expression</p> <p>Defines the single character that will appear in the data as the grouping separator.</p> <p>This property is applicable when dfdl:textNumberRep is 'standard' and dfdl:textStandardBase is 10. It must be set if dfdl:textNumberPattern contains a grouping separator symbol (it is a Schema Definition Error otherwise.) Empty string is not an allowable value.</p> <p>This property can be computed by way of an expression which returns a character. The expression must not contain forward references to elements which have not yet been processed.</p> <p>The string literal value is restricted in the same way as described in "Text Number Character Restrictions" in the description of the dfdl:textStandardDecimalSeparator property.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textStandardExponentRep	<p>DFDL String Literal or DFDL Expression</p> <p>Defines the actual character(s) that will appear in the data as the exponent indicator. If the empty string is specified then no exponent character will be used.</p> <p>This property is applicable when dfdl:textNumberRep is 'standard' and dfdl:textStandardBase is 10. Empty string is an allowable value, so that formats like NNN+M (meaning NNN x 10 with MM exponent) can be expressed.</p>



	<p>This property must be set even if the <code>dfdl:textNumberPattern</code> does not contain an 'E' (exponent) character. It is a Schema Definition Error if this property is not set or in scope for any number with <code>dfdl:representation</code> 'text'.</p> <p>This property can be computed by way of an expression which returns a DFDL String Literal character. The expression must not contain forward references to elements which have not yet been processed.</p> <p>The string literal value is restricted in the same way as described in "Text Number Character Restrictions" in the description of the <code>dfdl:textStandardDecimalSeparator</code> property.</p> <p>If <code>dfdl:ignoreCase</code> is 'yes' then the case of the string is ignored by the parser.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
<code>textStandardInfinityRep</code>	<p>DFDL String Literal</p> <p>The value used to represent infinity.</p> <p>Infinity is represented as a string with the positive or negative prefixes and suffixes from the <code>dfdl:textNumberPattern</code> applied.</p> <p>This property is applicable when <code>dfdl:textNumberRep</code> is 'standard', <code>dfdl:textStandardBase</code> is 10 and the simple type is float or double.</p> <p>If <code>dfdl:ignoreCase</code> is 'yes' then the case of the string is ignored by the parser.</p> <p>The string literal value is restricted in the same way as described in "Text Number Character Restrictions" in the description of the <code>dfdl:textStandardDecimalSeparator</code> property.</p> <p><u>It is a Schema Definition Error if empty string found as the property value.</u></p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
<code>textStandardNaNRep</code>	<p>DFDL String Literal</p> <p>The value used to represent NaN.</p> <p>NaN is represented as a string and the positive or negative prefixes and suffixes from the <code>dfdl:textNumberPattern</code> are not used.</p> <p>This property is applicable when <code>dfdl:textNumberRep</code> is 'standard', <code>dfdl:textStandardBase</code> is 10 and the simple type is float or double.</p> <p>If <code>dfdl:ignoreCase</code> is 'yes' then the case of the string is ignored by the parser.</p> <p>The string literal value is restricted in the same way as described in "Text Number Character Restrictions" in the description of the <code>dfdl:textStandardDecimalSeparator</code> property.</p> <p><u>It is a Schema Definition Error if empty string found as the property value.</u></p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>

textStandardZeroRep	<p>List of DFDL String Literals</p> <p>Valid values: empty string, any character string</p> <p>The whitespace separated list of alternative DFDL String Literals that are equivalent to zero, for example the characters 'zero'.</p> <p>The representation is examined for a match to one of the values of this property after padding has been trimmed away.</p> <p>On unparsing the first value is used.</p> <p>If dfdl:ignoreCase is 'yes' then the case of the string is ignored by the parser.</p> <p>The empty string means that there is no special literal string for zero.</p> <p>This property is applicable when dfdl:textNumberRep is 'standard' and dfdl:textStandardBase is 10.</p> <p>Each string literal in the list is restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed.</li> <li>• DFDL Byte Value entities ( %rXX; ) are not allowed.</li> <li>• DFDL Character class entities NL and ES are not allowed.</li> <li>• DFDL Character class entities WSP, WSP+, and WSP* are allowed.</li> </ul> <p>However, the WSP* entity cannot appear on its own as one of the string literals in the list. It must be used in combination with other text characters or entities so as to describe a representation that cannot ever be an empty string.</p> <p>It is a Schema Definition Error if the string literal contains any of the disallowed syntax constructs.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textStandardBase	<p>Non-negative Integer</p> <p>Valid Values 2, 8, 10, 16</p> <p>Indicates the number base.</p> <p>Only used when dfdl:textNumberRep is 'standard'.</p> <p>When base is not 10, xs:decimal, xs:float and xs:double are not supported.</p> <p>When dfdl:textNumberRep is 'zoned' dfdl:<del>textStandardBase</del> is not used and base 10 is assumed.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textZonedSignStyle	<p>Enum</p> <p>Specifies the code points that are used to overpunch the sign nibble when the dfdl:encoding is an ASCII-derived character set encoding. The location of this sign nibble is indicated in the dfdl:textNumberPattern.</p> <p>This property is applicable when dfdl:textNumberRep is 'zoned'.</p>

Deleted: textNumberBase

Deleted: 10

	<p>Used only when dfdl:encoding is an ASCII-derived character set encoding. The encoding must provide the character to single byte code point mapping used by the specified value of dfdl:textZonedSignStyle, as stated below.</p> <p>Valid values 'asciiStandard', 'asciiTranslatedEBCDIC', 'asciiCARealiaModified', and 'asciiTandemModified'</p> <p>Which characters are used to represent 'overpunched' (included) positive and negative signs, varies by encoding, Cobol compiler and system. The code points are fixed for EBCDIC systems but not for ASCII.</p> <p>In EBCDIC-based encodings, code points 0xC0 to 0xC9 or 0xF0 to 0xF9 represent a positive sign and digits 0 to 9 (typically characters '{ABCDEFGHI' or '0123456789'), and code points 0xD0 to 0xD9 or 0xB0 to 0xB9 represent a negative sign and digits 0 to 9 (typically characters 'JKLMNOPQR' or '£¥©\$¶¼½¾'). On parsing both ranges will be accepted. On unparsing the range 0xC0 to 0xC9 will be produced for positive signs and the range 0xD0 to 0xD9 will be produced for negative signs.</p> <p>asciiStandard: ASCII characters '0123456789' represent a positive sign and the corresponding digit. (Sign nibble for '+' is 0x3, which is the high nibble of these code points unmodified.) ASCII characters 'pqrstuvwxyz' represent negative sign and digits 0 to 9. (Code points 0x70 to 0x79)</p> <p>asciiTranslatedEBCDIC: The overpunched character is the ASCII equivalent of the typical EBCDIC above. So the characters '{ABCDEFGHI' still represent a positive sign and digits 0 to 9. (These are code points 0x7B, 0x41 through 0x49). The characters 'JKLMNOPQR' still represent negative sign and digits 0 to 9. (These are code points 0x7D, 0x4A through 0x52). This case comes up if EBCDIC zoned decimal data is translated to ASCII as if it were textual data.)</p> <p>asciiCARealiaModified<sup>23</sup>: In this style, the ASCII characters '0123456789' represent positive sign and digits 0 to 9 as in standard. However, ASCII characters from code points 0x20 to 0x29 are used for negative sign and the corresponding decimal digit. This doesn't translate well into printing characters. These characters include the space (' ') for zero, characters '!"#%&amp;' for 1 through 6, the single quote character "'" for 7, and the parenthesis '()' for 8 and 9.</p> <p>asciiTandemModified: In this style the ASCII characters '0123456789' represent positive sign and digits 0 to 9, but code points 0x80 to 0x89 are used to represent negative</p>
--	--

<sup>23</sup> Reference for this CA Realia 0x20 overpunch for negative sign is the article: "EBCDIC to ASCII Conversion of Signed Fields" [CARealia] where it says:

COBOL compilers that run on ASCII platforms have a "signed" data type that operates in a similar manner to the EBCDIC Signed field -- that is, they over punch the sign on the LSD. However, this is not standardized in ASCII, and different compilers use different overpunch codes. For example, Computer Associates' Realia compiler uses a 30 hex for positive values and a 20 hex for negative values, but Micro Focus® and Microsoft® use 30 hex for positive values and 70 hex for negative values.

	<p>sign and a digit. There are no corresponding code points in the standard ASCII encoding since these values are all above 128 (decimal). This means the resultant bytes are not code points in standard ASCII, so the modeller must specify an encoding like ISO-8859-1 in order for such zoned decimals to parse without an encoding error. (Note that neither ISO-8859-1 encoding nor Unicode have assigned glyphs for these code points. They are considered control characters.)</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
--	---

Table 31 Properties Specific to Number with Text Representation

The dfdl:textStandardDecimalSeparator, dfdl:textStandardGroupingSeparator, dfdl:textStandardExponentRep, dfdl:textStandardInfinityRep, dfdl:textStandardNaNRep, and dfdl:textStandardZeroRep must all be distinct, and it is a Schema Definition Error otherwise. Note that if dfdl:textStandardDecimalSeparator, dfdl:textStandardGroupingSeparator, or dfdl:textStandardExponentRep are expressions, this checking can only be carried out during processing (parsing or unparsing.)

Implementation note: This rule is in the interests of clarity, and is an extra constraint compared to ICU.

13.6.1 The dfdl:textNumberPattern Property

The dfdl:textNumberPattern describes how to parse and unparse text representations of number logical types with base 10.

The length of the representation of the number is determined first, and the number pattern is used only for conversion of the content text to and from a numeric logical infoset value.

The pattern described below is derived from the ICU DecimalFormat class described here: [\[ICUDecimal\]](#)

The pattern is an ICU-like syntax that defines where grouping separators, decimal separators, implied decimal points, exponents, positive signs and negative signs appear. It permits definition by either digits/fractions or significant digits.

13.6.1.1 dfdl:textNumberPattern for dfdl:textNumberRep 'standard'

When dfdl:textNumberRep is 'standard' this property only applies when dfdl:textStandardBase is 10.

The pattern comes in two parts separated by a semi-colon. The first is mandatory and applies to positive numbers, the second is optional and applies to negative numbers.

Examples: The first shows digits/fractions and positive/negative signs, the second shows exponent, the third shows virtual decimal point, the fourth shows scaling position.

<pre>+###,##0.00; (###,##0.00)  ##0.0#E0  000V00  PPP0000</pre>
---

The 'V' symbol is used to indicate the location of an implied decimal point for fixed point number representations. (This is an extension to the ICU pattern language.)

The 'P' symbol is used to indicate that a decimal scaling factor needs to be applied. (This is an extension to the ICU pattern language.)

The actual grouping separator, decimal separator and exponent characters are defined independently of the pattern.

The actual positive sign and negative sign are defined within the pattern itself.

Many characters in a pattern are taken literally; they are matched during parsing and output unchanged during unparsing. Special characters, on the other hand, stand for other characters, strings, or classes of characters. For example, the '#' character is replaced by a digit.

To insert a special character in a pattern as a literal, that is, without any special meaning, the character must be quoted. There are some exceptions to this which are noted below.

Symbol	Location	Meaning
0	Number	Digit
1-9	Number	'1' through '9' indicates rounding.
#	Number	Digit, zero shows as absent
.	Number	Decimal separator or monetary decimal separator
-	Number	Minus sign
,	Number	Grouping separator
E	Number	Separates mantissa and exponent in scientific notation. Need not be quoted in prefix or suffix.
+	Exponent	Prefix positive exponents with plus sign. Need not be quoted in prefix or suffix.
;	Subpattern boundary	Separates positive and negative subpatterns
'	Prefix or suffix	Used to quote special characters in a prefix or suffix, for example, "###" formats 123 to "#123". To create a single quote itself, use two in a row: "# o'clock".
*	Prefix or suffix boundary	Pad escape, precedes pad character
V	Number	Virtual decimal point marker. Only used with decimal, float and double simple types.
P	Number	Decimal scaling position. Only used with decimal, float and double simple types.
@	Number	Significant digits specifier. Only used with decimal simple type. Controls number of significant digits when used alone or in conjunction with the # character.

**Table 32 dfdl:textNumberPattern Special Characters**

A pattern contains a positive and negative subpattern, for example, "#,##0.00;(##0.00)". Each subpattern has a prefix, a numeric part, and a suffix. If there is no explicit negative subpattern, the negative subpattern is the minus sign prefixed to the positive subpattern. That is, "0.00" alone is equivalent to "0.00;-0.00". If there is an explicit negative subpattern, it serves only to specify the negative prefix and suffix; the number of digits, minimal digits, and other characteristics are ignored in the negative subpattern. That means that "#,##0.0#;(##0.0#)" has precisely the same result as "#,##0.0#;(##0.0#)".

The prefixes, suffixes, and various symbols used for infinity, digits, grouping separators, decimal separators, etc. may be set to arbitrary values, and they will appear properly during unparsing. However, care must be taken that the symbols and strings do not conflict, or parsing will be unreliable. For example, either the positive and negative prefixes or the suffixes must be distinct for parse to be able to distinguish positive from negative values.

The *grouping separator* is a character that separates clusters of integer digits to make large numbers more legible. It commonly used for thousands, but in some locales it separates ten-thousands. The *grouping size* is the number of digits between the grouping separators, such as 3 for "100,000,000" or 4 for "1 0000 0000". There are actually two different grouping sizes: One used for the least significant integer digits, the *primary grouping size*, and one used for all others, the *secondary grouping size*. In most locales these are the same, but sometimes they are different. For example, if the primary grouping interval is 3, and the secondary is 2, then this corresponds to the pattern "#,##,###0", and the number 123456789 is formatted as "12,34,56,789". If a pattern contains multiple grouping separators, the interval between the last one and the end of the integer defines the primary grouping size, and the interval between the last two defines the secondary grouping size. All others are ignored, so "#,##,###,####" == "###,###,####" == "#,##,###,####".

The P symbol is used to derive the location of an assumed decimal point when the point is not within the number that appears in the data. It acts as a decimal scaling factor.

The symbol P can be specified only as a continuous string of Ps in the leftmost or rightmost digit positions in the vpinteger region of the pattern.

It is a Schema Definition Error if any symbols other than "0", "1" through "9" or # are used in the vpinteger region of the pattern.

#### Examples

Data representation	Pattern	Value
123	PP000	0.00123
123	000PP	12300

**Table 33 Examples of P Symbol in the dfdl:textNumberPattern Property**

```

pattern      := subpattern (';' subpattern)?
subpattern   := prefix? ((number exponent?) | vpinteger) suffix?
number       := (integer ('.' fraction)?) | sigdigits

vpinteger    := pinteger | (vinteger exponent?)
pinteger     := ('P'* integer) | (integer 'P'*)
vinteger     := ('V'? integer) |
               ('*' 'V'? integer) |
               ('*' '0'* 'V'? '0'* '0') |
               (integer 'V'?)

prefix       := '\u0000'..' \uFFFD' - specialCharacters
suffix       := '\u0000'..' \uFFFD' - specialCharacters
integer      := '#'* '0'* '0'
fraction     := '0'* '#'*
sigDigits    := '#'* '@' '@'* '#'*
exponent     := 'E'? '+'? '0'* '0'
padSpec      := '*' padChar
padChar      := '\u0000'..' \uFFFD' - quote

Notation:
X*           0 or more instances of X
X?           0 or 1 instances of X
X|Y          either X or Y
C..D         any character from C up to D, inclusive
S-T          characters in S, except those in T

```

**Figure 4 dfdl:textNumberPattern BNF syntax**

The first subpattern is for positive numbers. The second (optional) subpattern is for negative numbers.

Not indicated in the BNF syntax above:

The grouping separator ',' can occur inside the integer region, between any two pattern characters of that region, as long as the number region is not followed by an exponent region.

Two grouping intervals are recognized: That between the decimal point and the first grouping symbol, and that between the first and second grouping symbols. These intervals are identical in most locales, but in some locales they differ. For example, the pattern "#,##,###" formats the number 123456789 as "12,34,56,789".

The pad specifier `padSpec` may appear before the prefix, after the prefix, before the suffix, after the suffix, or not at all.

In place of '0', the digits '1' through '9' in the number or `vpinteger` region may be used to indicate a rounding increment.

The term *maximum fraction digits* is the total number of '0' and '#' characters in the fraction sub-pattern above.

The term *minimum fraction digits* is the total number of '0' characters (only) in the fraction sub-pattern above.

The term *maximum integer digits* is a limit that is implementation-dependent, but must be at least 20 (which is the number of digits in a base 10 unsigned long).<sup>24</sup>

The term *minimum integer digits* is the total number of '0' characters (only) in the integer sub-pattern above.

- A pattern with a V symbol must not have # symbols to the right of the V symbol.
- A pattern with P symbols at the left end must not have # symbols .
- A pattern with P symbols at the right end can have # symbols.
- A pattern with a V symbol must not have @ or \* symbols.
- A pattern with P symbols must not have @ or E or \* symbols.

### Parsing

During parsing, grouping separators are removed from the data.

### Unparsing

Unparsing is guided by several parameters all of which can be specified using a pattern. The following description applies to formats that do not use scientific notation.

If the number of actual integer digits exceeds the *maximum integer digits*, then only the least significant digits are shown. For example, 1997 is formatted as "97" if the maximum integer digits is 2.

If the number of actual integer digits is less than the *minimum integer digits*, then leading zeros are added. For example, 1997 is formatted as "01997" if the minimum integer digits is 5.

If the number of actual fraction digits exceeds the *maximum fraction digits*, then half-even rounding is performed to the maximum fraction digits. For example, 0.125 is formatted as "0.12" if the maximum fraction digits is 2. This behavior can be changed by specifying a rounding increment and a rounding mode.

If the number of actual fraction digits is less than the *minimum fraction digits*, then trailing zeros are added. For example, 0.125 is formatted as "0.1250" if the minimum fraction digits is 4.

Trailing fractional zeros are not displayed if they occur *j* positions after the decimal, where *j* is less than the maximum fraction digits. For example, 0.10004 is formatted as "0.1" if the maximum fraction digits is four or less.

### Special Values

<sup>24</sup> Implementations which use current versions of the popular ICU library will allow 309 digits as *maximum integer digits*.

NaN is represented as a string determined by the `dfdl:textStandardNaNRep` property. This is the only value for which the prefixes and suffixes are not used.

Infinity is represented as a string with the positive or negative prefixes and suffixes applied. The infinity string is determined by the `dfdl:textStandardInfinityRep` property.

### Scientific Notation

Numbers in scientific notation are expressed as the product of a mantissa and a power of ten, for example, 1234 can be expressed as  $1.234 \times 10^3$ . The mantissa is typically in the half-open interval  $[1.0, 10.0)$  or sometimes  $[0.0, 1.0)$ , but it need not be. In a pattern, the exponent character immediately followed by one or more digit characters indicates scientific notation. Example: "0.###E0" formats the number 1234 as "1.234E3".

The number of digit characters after the exponent character gives the minimum exponent digit count. There is no maximum. Negative exponents are formatted using the minus sign, *not* the prefix and suffix from the pattern. This allows patterns such as "0.###E0 m/s". To prefix positive exponents with a plus sign, specify '+' between the exponent and the digits: "0.###E+0" will produce formats "1E+1", "1E+0", "1E-1", etc.

The minimum number of integer digits is achieved by adjusting the exponent. Example: 0.00123 formatted with "00.###E0" yields "12.3E-4". This only happens if there is no maximum number of integer digits. If there is a maximum, then the minimum number of integer digits is fixed at one.

The maximum number of integer digits, if present, specifies the exponent grouping. The most common use of this is to generate *engineering notation*, in which the exponent is a multiple of three, e.g., "###0.###E0". The number 12345 is formatted using "###0.###E0" as "12.345E3".

When using scientific notation, the formatter controls the digit counts using significant digits logic. The maximum number of significant digits limits the total number of integer and fraction digits that will be shown in the mantissa; it does not affect parsing. For example, 12345 formatted with "###0.###E0" is "12.3E3".

Exponential patterns may not contain grouping separators.

### Significant Digits

The '@' pattern character can be used with the '#' to control how many integer and fraction digits are needed to display the specified number of significant digits. The '@' only affects unparsing behavior. Examples:

Pattern	Minimum significant digits	Maximum significant digits	Number	Formatted Output
@@@	3	3	12345	12300
@@@	3	3	0.12345	0.123
@@##	2	4	3.14159	3.142
@@##	2	4	1.23004	1.23

**Table 34 Significant Digits '@' Symbol in the `dfdl:textNumberPattern` Property**

- Significant digit counts may be expressed using patterns that specify a minimum and maximum number of significant digits. These are indicated by the '@' and '#' characters. The minimum number of significant digits is the number of '@' characters. The maximum number of significant digits is the number of '@' characters plus the number of '#' characters following on the right. For example, the pattern "@@@" indicates exactly 3 significant digits. The pattern "@##" indicates from 1 to 3 significant digits. Trailing zero digits to the right of the decimal separator are suppressed after the minimum number of significant digits have been shown. For example, the pattern "@##" formats the number 0.1203 as "0.12".
- If a pattern uses significant digits, it may not contain a decimal separator, nor the '0' pattern character. Patterns such as "@00" or "@.###" are disallowed.



- Any number of '#' characters may be prepended to the left of the leftmost '@' character. These have no effect on the minimum and maximum significant digits counts, but may be used to position grouping separators. For example, "#,#@#" indicates a minimum of one significant digits, a maximum of two significant digits, and a grouping size of three.
- The number of significant digits has no effect on parsing.
- Significant digits may be used together with exponential notation. For example, the pattern "@@###E0" is equivalent to "0.0###E0".
- The '@' pattern character can be used only in 'standard' textNumberRep (not 'zoned'), and excludes the 'P' and 'V' pattern characters. It is a Schema Definition Error if the '@' pattern character appears in 'zoned' textNumberRep, or in conjunction with the 'P' or 'V' pattern characters.

### Padding

Padding may be specified through the pattern syntax. In a pattern the pad escape character, followed by a single pad character, causes padding to be parsed and formatted. The pad escape character is ". For example, "x#.##0.00" formats 123 to "x123.00", and 1234 to "1,234.00".

When padding is in effect, the width of the positive subpattern, including prefix and suffix, determines the format width. For example, in the pattern "x #0 o'clock", the format width is 10.

The width is counted in 16-bit code units.

Some parameters which usually do not matter have meaning when padding is used, because the pattern width is significant with padding. In the pattern "x ###.###.###0.##", the format width is 14. The initial characters "x###.##" do not affect the grouping size or maximum integer digits, but they do affect the format width.

Padding may be inserted at one of four locations: before the prefix, after the prefix, before the suffix, or after the suffix. If there is no prefix, before the prefix and after the prefix are equivalent, likewise for the suffix.

When specified in a pattern, the 32-bit codepoint immediately following the pad escape is the pad character. This may be any character, including a special pattern character. That is, the pad escape escapes the following character. If there is no character after the pad escape, then the pattern is illegal.

Note: Padding specified through the pattern syntax is distinct from, and in addition to, padding specified using dfdl:textPadKind.

### Rounding

How rounding is controlled is given by dfdl:textNumberRounding. The rounding increment may be specified in the dfdl:textNumberPattern itself using digits '1' through '9' or using an explicit increment in dfdl:textNumberRoundingIncrement. For example, 1230 rounded to the nearest 50 is 1250. 1.234 rounded to the nearest 0.65 is 1.3.

- Rounding only affects the string produced by unparsing. It does not affect parsing or change any numerical values.
- In a pattern, digits '1' through '9' specify rounding, but otherwise behave identically to digit '0'. For example, "#.##50" specifies a rounding increment of 50.
- Using digits in a pattern, rounding is always 'half even', meaning rounds towards the nearest integer, or towards the nearest even integer if equidistant.

Using an explicit rounding increment, dfdl:textNumberRoundingMode determines how values are rounded.

#### 13.6.1.2 dfdl:textNumberPattern for dfdl:textNumberRep 'zoned'

When dfdl:textNumberRep is 'zoned' a subset of the number pattern language described in Section 13.6.1.1 dfdl:textNumberPattern for dfdl:textNumberRep 'standard' is used.

Only the pattern for positive numbers is used. It is a Schema Definition Error if the negative pattern is specified.

In addition, only the following pattern characters may be used:

- '+' MUST BE present at the beginning or end of the pattern to indicate whether the leading or trailing digit carries the overpunched sign, if the logical type is signed
- '-' MAY BE present at the beginning or end of the pattern to indicate whether the leading or trailing digit carries the overpunched sign, if the logical type is unsigned. If logical type is unsigned and dfdl:textNumberPolicy 'lax' specified it is a Schema Definition Error if no '-' is present.
- 'V' MAY BE used to indicate the location of an implied decimal point
- 'P' MAY BE used to indicate the decimal scaling
- '0-9' indicates the number of needed digits (including overpunched).
- '#' indicates the number of optional digits.

Rounding occurs as described under Rounding in 13.6.1.1 dfdl:textNumberPattern for dfdl:textNumberRep 'standard'

#### 13.6.2 Converting logical numbers to/from text representation

- Signed numbers with dfdl:textNumberRep 'standard' and dfdl:textStandardBase 10 are mapped using the dfdl:textNumberPattern.
- Signed numbers with dfdl:textNumberRep 'standard' and dfdl:textStandardBase not 10 are mapped to an unsigned representation. On unparsing the minimum number of characters to represent the digits is output and it is a processing error if the value is negative.
- Signed numbers with dfdl:textNumberRep 'zoned' are mapped using the dfdl:textNumberPattern to indicate the position of the sign and virtual decimal point. On parsing if the sign is not overpunched, that is it does not have a sign, it is treated as positive. On unparsing the sign is always overpunched.
- Unsigned numbers with dfdl:textNumberRep 'standard' and dfdl:textStandardBase 10 are mapped using the dfdl:textNumberPattern. On parsing it is a processing error if the data are negative.
- Unsigned numbers with dfdl:textNumberRep 'standard' and dfdl:textStandardBase not 10 are mapped to an unsigned representation. On unparsing the minimum number of characters to represent the digits is output. .
- Unsigned numbers with dfdl:textNumberRep 'zoned' are mapped using the dfdl:textNumberPattern to indicate the position of the sign and virtual decimal point. On parsing it is a processing error if the data are negative. On unparsing the data are not overpunched with a sign.

### 13.7 Properties Specific to Number with Binary Representation

These properties are applicable to simple type `xs:decimal` and its derived types which include all the signed and unsigned integer types. These properties are not applicable to types `xs:float` and `xs:double`. See section 13.8. Note that simple types derived from `xs:decimal` do not imply base-10 representations in the data stream.

Property Name	Description								
binaryNumberRep	<p>Enum</p> <p>Valid values are 'packed', 'bcd', 'binary', 'ibm4690Packed'</p> <p>Allowable values for each number type are:</p> <table> <tr> <th>Logical Type</th><th>Permitted Value</th></tr> <tr> <td>Decimal, Integer, NonNegativeInteger</td><td>packed, bcd, binary, ibm4690Packed</td></tr> <tr> <td>Long, Int, Short, Byte,</td><td>packed, binary, ibm4690Packed (but not bcd)</td></tr> <tr> <td>UnsignedLong, Unsignedint, UnsignedShort, UnsignedByte</td><td>packed, bcd, binary, ibm4690Packed</td></tr> </table> <p>'packed' means represented as an IBM 390 packed decimal. Each byte contains two decimal digits, except for the least significant byte, which contains a sign in the least significant nibble.</p> <p>'bcd' means represented as a binary coded decimal with two digits per byte.</p> <p>'binary' means represented as twos complement for signed types and unsigned base-2 binary for unsigned types.</p> <p>Note that the maximum allowed value for twos-complement and unsigned base-2 binary integers is implementation-dependent, but must be at least that of a <code>xs:long</code> type, which is the equivalent of an 8 byte/64-bit signed integer.</p> <p>'ibm4690Packed' is a variant of a packed decimal having the following characteristics:</p> <ul style="list-style-type: none"> <li>• Nibbles represent digits 0 - 9 in the usual BCD manner.</li> <li>• A positive value is simply indicated by digits.</li> <li>• A negative number is indicated by digits with the most significant nibble being xD.</li> <li>• If a positive or negative value packs to an odd number of nibbles, an extra xF nibble is added as the most significant nibble.</li> </ul> <p>For all values, the <code>dfdl:byteOrder</code> property is used to determine the numeric significance of the bytes making up the representation, and the <code>dfdl:bitOrder</code> property is used to determine the numeric significance of the bits within a byte.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>	Logical Type	Permitted Value	Decimal, Integer, NonNegativeInteger	packed, bcd, binary, ibm4690Packed	Long, Int, Short, Byte,	packed, binary, ibm4690Packed (but not bcd)	UnsignedLong, Unsignedint, UnsignedShort, UnsignedByte	packed, bcd, binary, ibm4690Packed
Logical Type	Permitted Value								
Decimal, Integer, NonNegativeInteger	packed, bcd, binary, ibm4690Packed								
Long, Int, Short, Byte,	packed, binary, ibm4690Packed (but not bcd)								
UnsignedLong, Unsignedint, UnsignedShort, UnsignedByte	packed, bcd, binary, ibm4690Packed								
binaryDecimalVirtualPoint	<p>Integer.</p> <p>Used when base <code>simpleType</code> is <code>xs:decimal</code>.</p> <p>An integer that represents the position of an implied decimal point within a number, or specify 0.</p> <p>If you specify 0 then there is no virtual decimal point</p>								

	<p>If you specify a positive integer, the position of the decimal point is moved from the least-significant side of the number toward the most-significant side of the number. For example, if 3 is specified then, the integer value 1234 represents 1.234. This is equivalent to dividing by <math>10^3</math>.</p> <p>If you specify a negative integer, the position of the decimal point is moved from the least significant side of the number further in the less-significant direction. For example, if you specify -3, the integer value 1234 represents 1 234 000. This is equivalent to multiplying by <math>10^3</math>.</p> <p><u>When unparsing, if the property value is not sufficient to remove the decimal point from the infoSet value, it is a processing error. This is true even if the resultant number can be converted into an integer (that is, all digits after the decimal point are zero) because it is an example of excess precision.</u></p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
binaryPackedSignCodes	<p>List of Characters</p> <p>Used only when dfdl:binaryNumberRep or dfdl:binaryCalendarRep is 'packed'</p> <p>A whitespace separated string giving the hex sign nibbles to use for a positive value, a negative value, an unsigned value, and zero.</p> <p>Valid values for positive nibble: A, C, E, F</p> <p>Valid values for negative nibble: B, D</p> <p>Valid values for unsigned nibble: F</p> <p>Valid values for zero sign: A C E F 0</p> <p>Example: 'C D F C' – typical S/390 usage</p> <p>Example: 'C D F 0' – handle special case for zero</p> <p>On parsing, whether to accept all valid values for a positive, negative or unsigned number, and for zero, is governed by the dfdl:binaryNumberCheckPolicy property. On unparsing, the specified values are always used.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
binaryNumberCheckPolicy	<p>Enum</p> <p>Values are 'strict' and 'lax'.</p> <p>Indicates how lenient to be when parsing binary numbers.</p> <p>If 'lax' then the parser tolerates all valid alternatives where such alternatives exist. Specifically, for dfdl:binaryNumberRep 'packed' the sign nibble for positive, negative, unsigned and zero is allowed to be any of the valid respective values.</p> <p>On unparsing, the specified value is always used</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 35 Properties Specific to Number with Binary Representation****13.7.1 Converting Logical Numbers to/from Binary Representation**

When unparsing a binary number (packed decimal or twos-complement) and excess precision is supplied in the InfoSet no rounding occurs. It is a processing error.

### 13.7.1.1 Converting Base-2 Binary Numbers

For both parsing and unparsing, the bit string that represents the content region for a base-2 binary number is converted to/from an Infoset value by a calculation that involves the length and the `dfdl:byteOrder` and `dfdl:bitOrder` properties.

When parsing, DFDL specifies how an unsigned integer of unbounded magnitude is computed from a bit string based on its length, and the `dfdl:byteOrder` and `dfdl:bitOrder` properties. For signed types, this unbounded integer is converted into a signed value by way of the well-known twos-complement scheme, and for the `xs:decimal` type, the `dfdl:binaryDecimalVirtualPoint` property can be used to convert this integer into a decimal value with an integer and a fractional component.

A DFDL implementation can use any conversion technique consistent with this description.

**Deleted:** y

**Deleted:** For unparsing, the `dfdl:fillByte` property can also be involved. ¶

**Deleted:** y

**Deleted:** , and for both `xs:decimal` and the integer types the `dfdl:binaryDecimalVirtualPoint` or to scale up the integer by some scale factor...

### 13.7.1.2 Bit strings, Alignment, and `dfdl:fillByte`

The `dfdl:alignmentUnits` of 'bits', and `dfdl:alignment` of '1' can be used to position a bit string anywhere in the data stream without regard for any other grouping of bits into bytes.

The numeric value of the unsigned integer represented by a bit string is unaffected by alignment.

When unparsing a bit string, alignment may cause the bits within the bit string to occupy only some of the bits within a byte of the data stream. The bits of data in the alignment fill region are unspecified by the elements of the DFDL schema, and when parsing, neither they, nor any data computed from them are put into the DFDL infoset. During unparsing, such unspecified bits are filled in using the value of the `dfdl:fillByte` property. Corresponding bits from the `dfdl:fillByte` value are used to fill in unspecified bits of the data stream. That is, if bit *K* (*K* will be 1 or greater, but less than or equal to 8) of a data stream byte is unspecified, its value will be taken from bit *K* of the `dfdl:fillByte` property value.

Since the value of any bit string element is unaffected by alignment, the logical unsigned integer value for a bit-string is always computed as if the first bit were at position 1 of the bit stream. If the `dfdl:length` for the bit-string evaluates to *M*, then the bit-string conceptually occupies bits 1 to *M* of a data stream for purposes of computing its value.

### 13.7.1.3 Bits within Bit Strings of Length ≤ 8

Any time the length in bits is < 8, then when set, the bit at position *Z* supplies value  $2^{(M-Z)}$ , and the value of the bit string as an integer is the sum of these values for each of its bits.

### 13.7.1.4 Bits within Bit Strings of Length > 8

Call *M* the length of the bit string element in bits. In general, when *M* > 8 the contribution of a bit in position *i* to the numeric value of a bit string is given by a formula specific to the `dfdl:byteOrder`.

For `dfdl:byteOrder` of 'bigEndian' the value of bit *i* is given by  $2^{(M-i)}$ .

For `dfdl:byteOrder` of 'littleEndian' the value of bit *i* is given by a more complex formula. The following pseudo code computes the value of a bit in a littleEndian bit string. It is just a very big expression, but is spread out over many local variables to illustrate the various sub-calculations clearly. DFDL implementations may use any way of converting bit strings to the corresponding integer values that is consistent with this:

In the pseudo code below:

- '%' is modular division (division where remainder is returned)
- '/' is regular division (quotient is returned)
- the expression '*a* ? *b* : *c*' means 'if *a* is true, then the value is *b*, otherwise the value is *c*'

```
littleEndianBitValue(bitPosition, bitStringLength)
    assert bitPosition >= 1;
    assert bitStringLength >= 1;
    assert bitStringLength >= bitPosition;
    numBitsInFinalPartialByte = bitStringLength % 8;
    numBitsInWholeBytes = bitStringLength -
```

```

        numBitsInFinalPartialByte;
    bitPosInByte = ((bitPosition - 1) % 8) + 1;
    widthOfActiveBitsInByte = (bitPosition <= numBitsInWholeBytes)
        ? 8 : numBitsInFinalPartialByte;
    placeValueExponentOfBitInByte = widthOfActiveBitsInByte -
        bitPosInByte;
    bitValueInByte = 2^placeValueExponentOfBitInByte;
    byteNumZeroBased = (bitPosition - 1)/8;
    scaleFactorForBytePosition = 2^(8 * byteNumZeroBased);
    bitValue = bitValueInByte * scaleFactorForBytePosition;
    return bitValue;

```

**Figure 5 Little Endian bit position and value**

#### 13.7.1.4.1 Examples of Unsigned Integer Conversion

Consider the first three bytes of the data stream. Imagine their numeric values as 0x5A 0x92 0x00.

```

Positions:
00000000 01111111 11122222
12345678 90123456 78901234
Bits:
01011010 10010010 00000000
Hex values
    5    A    9    2    0    0

```

Beginning at bit position 1, (the very first bit) if we consider the first two bytes as a bigEndian short, the value will be 0x5A92.

```

<xs:element name="num" type="unsignedShort"
  dfdl:alignment="1"
  dfdl:alignmentUnits="bytes"
  dfdl:byteOrder="bigEndian"
  dfdl:representation="binary"
  dfdl:binaryNumberRep="binary"/>

```

As a littleEndian short, the value will be 0x925A.

```

<xs:element name="num" type="unsignedShort"
  dfdl:alignment="1"
  dfdl:alignmentUnits="bytes"
  dfdl:byteOrder="littleEndian"
  dfdl:representation="binary"
  dfdl:binaryNumberRep="binary"/>

```

Now let us examine a bit string of length 13, beginning at position 2.

```

<xs:sequence>
  <xs:element name="ignored" type="unsignedByte"
    dfdl:alignment="1"
    dfdl:alignmentUnits="bits"
    dfdl:lengthUnits="bits"
    dfdl:length="1"
    dfdl:representation="binary"
    dfdl:binaryNumberRep="binary"/>
  <xs:element name="x" type="unsignedShort"
    dfdl:alignment="1"
    dfdl:alignmentUnits="bits"
    dfdl:byteOrder="bigEndian"
    dfdl:lengthUnits="bits"
    dfdl:length="13"
    dfdl:representation="binary"
    dfdl:binaryNumberRep="binary"/>
  ...
</xs:sequence>

```

Let's examine the same data stream and consider the bit positions that make up element 'x', which are the bits at positions 2 through 14 inclusive.

```
Positions:
00000000 01111111 11122222
12345678 90123456 78901234
Bits:
1011010 100100
```

Since alignment does not affect logical value, we will obtain the same logical value as if we realigned the bits. That is, the value is the same as if we began the bits of the element's representation with bit position 1.

```
Realigned Positions:
00000000 01111111 11122222
12345678 90123456 78901234
Bits:
10110101 00100
```

The DFDL schema fragment above gives element 'x' the dfdl:byteOrder 'bigEndian' property. In this case the place value of each position is given by  $2^{(M-i)}$

PlaceValue positions  $2^{(M-i)}$

```
...11110 00000000
...21098 76543210
Bit values
...10110 10100100
Hex values
1 6 A 4
```

The value of element 'x' is 0x16A4. Notice how it is the most-significant byte -- which is the first byte when big endian -- that becomes the partial byte (having fewer than 8 bits) in the case where the length of the bit string is not a multiple of 8 bits.

For dfdl:byteOrder of 'littleEndian'. The place values of the individual bits are not as easily visualized. However there is still a basic formula (given in the pseudo code in 13.7.1.4 Bits within Bit Strings of Length > 8) and value.

Looking again at our realigned positions:

```
Realigned Positions:
00000000 01111111 11122222
12345678 90123456 78901234
Bits:
10110101 00100
```

The place values of each of these bits, for little endian byte order can be seen to be:

```
PlaceValue positions
00000000 ...111100
76543210 ...21098
Bit values
10110101 ...00100
Hex values
B 5 0 4
```

We must reorder the bytes for little endian byte order. The value of element 'x' is 0x04B5. In little endian form, the first 8 bits make up the first byte, and that contains the least-significant byte of the logical numeric unsignedShort value. The additional bits of the partial byte are once again the most significant byte; however, for little endian form, this is the second byte. The second byte contains only 5 bits, those make up the least significant 5 bits of that byte, but that logical 5-bit value makes up the most-significant byte of the unsignedShort integer.

Now let us examine the 13 bits beginning at position 2, in the context where dfdl:byteOrder is 'littleEndian' and dfdl:bitOrder is 'leastSignificantBitFirst' and dfdl:binaryNumberRep is 'binary'. In this case, the bit positions are assigned differently. Below the bytes are shown left-to-right:

```
Positions:
00000000 11111110 22222111
87654321 65432109 43210987
Bits:
01011010 10010010 00000000
Hex values
5 A 9 2 0 0
```

The bits of interest are highlighted above. If we redisplay this same data, but reversing the order of the bytes to right-to-left, then we get:

```
Positions:
22222111 11111110 00000000
43210987 65432109 87654321
Bits:
00000000 10010010 01011010
Hex values
0 0 9 2 5 A
```

The above shows more clearly that we are looking at a contiguous region of bits containing

```
0 1001 0010 1101
```

or the value 0x092D.

#### 13.7.1.5 Converting Packed Decimal Numbers

Signed numbers with dfdl:binaryNumberRep 'packed' are parsed using a nibble to indicate the sign. The unsigned nibble is treated as positive. On unparsing the sign nibble is written according to dfdl:binaryPackedSignCodes. The unsigned nibble is never written.

Signed numbers with dfdl:binaryNumberRep 'bcd' are always positive. On unparsing it is a processing error if the InfoSet data is negative.

Signed numbers with dfdl:binaryNumberRep 'ibm4690Packed' are parsed using the sign nibble to identify negative values. There is no sign nibble for positive values. On unparsing the nibble 0xD is written for negative values.

Unsigned numbers with dfdl:binaryNumberRep 'packed' are parsed if the nibble is positive or unsigned. It is a processing error if the data is negative. On unparsing the unsigned nibble is used.

Unsigned numbers with dfdl:binaryNumberRep 'bcd' are readily parsed as BCD data is always positive.

Unsigned numbers with dfdl:binaryNumberRep 'ibm4690Packed' are parsed if there is no sign nibble of 0xD to identify a negative value. It is a processing error if the data is negative. On unparsing no sign nibble is written.



### 13.8 Properties Specific to Float/Double with Binary Representation

Property Name	Description
binaryFloatRep	<p>Enum or DFDL Expression</p> <p>This specifies the encoding method for the float and double.</p> <p>Valid values are 'ieee', 'ibm390Hex'. This property can be computed by way of an expression which returns a string of 'ieee' or 'ibm390Hex'. The expression must not contain forward references to elements which have not yet been processed.</p> <p>The enumeration value 'ieee' refers to the IEEE 754-1985 specification.</p> <p>For both 'ieee' and 'ibm390hex', an xs:float must have a physical length of 4 bytes. It is a Schema Definition Error if there is a specified length not equivalent to 4 bytes.</p> <p>Similarly, for both 'ieee' and 'ibm390hex', an xs:double must have a physical length of 8 bytes. It is a Schema Definition Error if there is a specified length not equivalent to 8 bytes.</p> <p>The dfdl:byteOrder property is used to construct a value from the bytes in the binary representation.</p> <p>Note: The DFDL Infoset float and double data types match the precision of the IEEE specification. There may be precision/rounding issues when converting IBM float/double to/from the DFDL infoset float/double types.</p> <p>Half-precision IEEE and quad-precision IEEE/IBM are not supported.<sup>25</sup></p> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 36 Properties Specific to Float/Double with Binary Representation**

### 13.9 Properties Specific to Boolean with Text Representation

Property Name	Description
textBooleanTrueRep	<p>List of DFDL String Literals or DFDL Expression</p> <p>A whitespace separated list of representations to be used for 'true'. These are compared after trimming when parsing, and before padding when unparsing.</p> <p>If dfdl:lengthKind is 'explicit' or 'implicit' and either dfdl:textPadKind or dfdl:textTrimKind is 'none' then both dfdl:textBooleanTrueRep and dfdl:textBooleanFalseRep must have the same length else it is a Schema Definition Error.</p> <p>This property can be computed by way of an expression which returns a string of whitespace separated list of values. The expression must not contain forward references to elements which have not yet been processed.</p> <p>On unparsing the first value is used</p>

<sup>25</sup> Note that XSDL 1.1 moved to IEEE 754-2008 only because of new decimal support, and not for enhanced float support. That's why in XSDL 1.1 there are still just the xs:float and xs:double built-in types. Any future support for half-precision and quad-precision in XSDL would very likely be implemented by adding new built-in types that derive from xs:anySimpleType. It is likely therefore that future DFDL support for half-precision and quad-precision will build on XSDL.

	<p>If dfdl:ignoreCase is 'yes' then the case of the string is ignored by the parser.</p> <p><i>Text Boolean Character Restrictions:</i> The string literal is restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed</li> <li>• The DFDL byte value entity ( %rXX; ) is not allowed.</li> <li>• DFDL Character classes NL, WSP, WSP+, WSP*, and ES are not allowed</li> </ul> <p>It is a Schema Definition Error if the string literal <u>is the empty string, or</u> contains any of the disallowed constructs.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textBooleanFalseRep	<p>List of DFDL String Literals or DFDL Expression</p> <p>A whitespace separated list of representations to be used for 'false' These are compared after trimming when parsing, and before padding when unparsing.</p> <p>If dfdl:lengthKind is 'explicit' or 'implicit' and either dfdl:textPadKind or dfdl:textTrimKind is 'none' then both dfdl:textBooleanTrueRep and dfdl:textBooleanFalseRep must have the same length else it is a Schema Definition Error.</p> <p>This property can be computed by way of an expression which returns a string of whitespace separated list of values. The expression must not contain forward references to elements which have not yet been processed.</p> <p>On unparsing the first value is used</p> <p>If dfdl:ignoreCase is 'yes' then the case of the string is ignored by the parser.</p> <p>The string literal value is restricted in the same way as described in "Text Boolean Character Restrictions" in the description of the dfdl:textBooleanTrueRep property.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textBooleanJustification	<p>Enum</p> <p>Valid values 'left', 'right', 'center'</p> <p>Controls how the data is padded or trimmed on parsing and unparsing.</p> <p>Behavior as for dfdl:textStringJustification.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textBooleanPadCharacter	<p>DFDL String Literal</p> <p>The value that is used when padding or trimming boolean elements. The value can be a single character or a single byte. If a character, then it can be specified using a literal character or using DFDL entities.</p> <p>If a byte, then it must be specified using a single byte value entity.</p> <p>If a pad character is specified when lengthUnits is 'bytes' then the pad character must be a single-byte character.</p> <p>If a pad byte is specified when lengthUnits is 'characters' then</p> <ul style="list-style-type: none"> <li>• the dfdl:encoding must be a fixed-width encoding</li> </ul>

	<ul style="list-style-type: none"> <li>padding and trimming must be applied using a sequence of N pad bytes, where N is the width of a character in the fixed-width encoding.</li> </ul> <p>The string literal value is restricted in the same way as described in "Pad Character Restrictions" in the description of the <code>dfdl:textStringPadCharacter</code> property.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
--	---

**Table 37 Properties Specific to Boolean with Text Representation****13.10 Properties Specific to Boolean with Binary Representation**

Property Name	Description
<code>binaryBooleanTrueRep</code>	<p>Non-negative Integer</p> <p>This value, treated as a binary <code>xs:unsignedInt</code> (See Section 13.7.1 Converting Logical Numbers to/from Binary Representation ), gives the representation to be used for 'true'</p> <p>If this property value is the empty string, when parsing it means <code>dfdl:binaryBooleanTrueRep</code> is any value other than <code>dfdl:binaryBooleanFalseRep</code>; when unparsing, the one's complement of the <code>dfdl:binaryBooleanFalseRep</code> will be used.</p> <p>The length of the data value of the element must be between 1 bit and 32 bits (4 bytes) as described in Section 12.3.7.2. It is a Schema Definition Error if the value (when provided) of <code>dfdl:binaryBooleanTrueRep</code> cannot fit as an unsigned binary integer in the specified length.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>
<code>binaryBooleanFalseRep</code>	<p>Non-negative Integer</p> <p>This value, treated as a binary <code>xs:unsignedInt</code> (See Section 13.7.1 Converting Logical Numbers to/from Binary Representation ), gives the representation to be used for 'false'</p> <p>The length of the data value of the element must be between 1 bit and 32 bits (4 bytes) as described in Section 12.3.7.2. It is a Schema Definition Error if the value of <code>dfdl:binaryBooleanFalseRep</code> cannot fit as an unsigned binary integer in the specified length.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:simpleType</code></p>

**Table 38 Properties Specific to Boolean with Binary Representation****13.11 Properties specific to Calendar with Text or Binary Representation**

The properties describe how a calendar is to be interpreted including a unparsing pattern property plus properties that qualify the pattern.

These properties can be used when a calendar has `dfdl:representation` 'text' or `dfdl:representation` 'binary' and a packed decimal representation.

Property Name	Description
<code>calendarPattern</code>	<p>String</p> <p>Defines the ICU pattern that describes the format of the calendar. The pattern defines where the year, month, day, hour, minute, second, fractional second and time zone components appear. See <code>calendarPattern</code> property section below.</p>

	<p>When the dfdl:representation is 'binary' and the representation is a packed decimal then the pattern can contain only characters and symbols that always result in the presentation of digits.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>								
calendarPatternKind	<p>Enum</p> <p>Valid values 'explicit', 'implicit'</p> <p>'explicit' means the pattern is given by dfdl:calendarPattern, 'implicit' means the pattern is derived from the XML schema date/time type.</p> <table border="1"> <thead> <tr> <th>Logical Type</th><th>Default Pattern</th></tr> </thead> <tbody> <tr> <td>xs:date</td><td>yyyy-MM-dd</td></tr> <tr> <td>xs:dateTime</td><td>yyyy-MM-dd'T'HH:mm:ss</td></tr> <tr> <td>xs:time</td><td>HH:mm:ssZ</td></tr> </tbody> </table> <p>Annotation: dfdl:element, dfdl:simpleType</p>	Logical Type	Default Pattern	xs:date	yyyy-MM-dd	xs:dateTime	yyyy-MM-dd'T'HH:mm:ss	xs:time	HH:mm:ssZ
Logical Type	Default Pattern								
xs:date	yyyy-MM-dd								
xs:dateTime	yyyy-MM-dd'T'HH:mm:ss								
xs:time	HH:mm:ssZ								
calendarCheckPolicy	<p>Enum</p> <p>Valid values are 'strict', 'lax'</p> <p>Indicates how lenient to be when parsing against the pattern.</p> <p>See Section 13.11.2 The dfdl:calendarCheckPolicy Property below for details of the specific behaviors for 'strict' and 'lax'.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>								
calendarTimeZone	<p>String</p> <p>This property provides the time zone that will be assumed if no time zone explicitly occurs in the data.</p> <p>Valid values specify a UTC time zone offset by matching the regular expression:</p> <p>(UTC) ([+ -] ([01]\d \d) ((([:][0-5]\d){1,2})?))?</p> <p>In addition, empty string can be specified to indicate "no time zone", or the IANA time zone format (also known as the Olson time zone format) may be used. (e.g, America/New_York) See <a href="#">[IANATimeZone]</a>.</p> <p>Note that this property is used when parsing only.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>								
calendarObserveDST	<p>Enum</p> <p>Valid values are 'yes', 'no'</p> <p>Whether the time zone given in dfdl:calendarTimeZone observes daylight savings time.</p> <p>Ignored if dfdl:calendarTimeZone is specified in UTC format, or if dfdl:calendarTimeZone is empty string. That is, this property is used only if the dfdl:calendarTimeZone is in IANA (also known as Olson) format <a href="#">[IANATimeZone]</a>.</p> <p>This property applies to parsing only.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>								
calendarFirstDayOfWeek	<p>Enum</p> <p>Valid values 'Monday' ... 'Sunday'</p>								

Deleted: )

	<p>The day of the week upon which a new week is considered to start.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
calendarDaysInFirstWeek	<p>Non-negative Integer</p> <p>Valid values 1 to 7</p> <p>Specify the number of days of the new year that must fall within the first week.</p> <p>The start of a year usually falls in the middle of a week. If the number of days in that week is less than the value specified here, the week is considered to be the last week of the previous year; hence week 1 starts some days into the new year. Otherwise it is considered to be the first week of the new year; hence week 1 starts some days before the new year.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
calendarCenturyStart	<p>Non-negative Integer</p> <p>Valid values 0 to 99.</p> <p>This property determines on parsing how two-digit years are interpreted. Specify the two digits that start a 100-year window that contains the current year. For example, if you specify 89, and the current year is 2006, all two-digit dates are interpreted as being in the range 1989 to 2088. A two-digit year less than 89 will be interpreted as 20nn and a two-digit year more than or equal to 89 will be treated as 19nn.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
calendarLanguage	<p>String or DFDL Expression</p> <p>The language that is used when the pattern produces a presentation in text.</p> <p>The value must match the regular expression:  <code>([A-Za-z]{1,8}([\-_][A-Za-z0-9]{1,8})*)</code></p> <p>It is a Schema Definition Error otherwise.</p> <p><u>The expression must not contain forward references to elements which have not yet been processed.</u></p> <p>All DFDL Implementations must support dfdl:calendarLanguage value "en".</p> <p>DFDL implementations may support additional values, however, the value of the dfdl:calendarLanguage property is always interpreted as a Unicode Language Identifier as defined by <a href="#">[LDML]</a>, and <a href="#">[CLDR]</a>.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 39 Properties specific to Calendar with Text or Binary Representation****13.11.1 The dfdl:calendarPattern property**

The dfdl:calendarPattern describes how to parse and unparse text and binary representations of dateTime, date and time logical types. The pattern is primarily used on unparsing to define the format but is also used to aid parsing.

The pattern is derived from the ICU SimpleDateTimeFormat class described here: [\[ICUDateTime\]](#), which uses symbols defined by [\[LDML\]](#).

An extension is the formatting symbol I which means accept a subset of ISO 8601 [\[ISO8601\]](#) compliant calendars

Symbol	Meaning	Presentation	Example	
G	era designator	Text	G	AD
y	year	Number	y, yyyy yy	1996 96
u	year (allows negative years)	Number	u	1900, 0, -500
Y	year (of the week of year)	Number	Y	1997
M	month in year	Text & Number	M, MM MMM MMMM MMMM M	09 Sept September S
d	day in month	Number	d dd	2 02
h	hour in am/pm (1~12)	Number	h hh	7 07
H	hour in day (0~23)	Number	H HH	0 00
m	minute in hour	Number	m mm	4 04
s	second in minute	Number	s ss	5 05
S	fractional second	Number	S SS SSS	2 <del>23</del> 235
E	day of week	Text	E, EE,EEE EEEE EEEEEE EEEEEE	Tues Tuesday T Tu
e	day of week (local)	Text & Number	e, ee eee eeee eeeee eeeeee	2 Tues Tuesday T Tu
D	day in year	Number	D	189
F	day of week in month	Number	F	2 (2nd Wed in July)
w	week in year	Number	w, ww	27
W	week in month	Number	W	2
a	am/pm marker	Text	A	pm

Deleted: (see note 1)

Deleted: 24

k	hour in day (0~24 )	Number	k kk	2, 24 02, 24
K	hour in am/pm (0~11)	Number	K KK	0 00
z	time zone: specific non-location	Text	z, zz, zzz zzzz	PDT Pacific Daylight Time
Z	time zone: ISO8601 basic format time zone: long localized GMT	Text	Z, ZZ, ZZZ ZZZZ	-0800, +0000 GMT-08:00, GMT+00:00
O	time zone: localized GMT	Text	O OOOO	GMT-8 GMT-08:00
v	time zone: generic non-location	Text	v vvvv	PT Pacific Time
V	time zone: short time zone ID time zone: long time zone ID time zone: exemplar city time zone: generic location.	Text	V VV VVV VVVV	uslax America/Los_Angeles Los Angeles Los Angeles Time
x	time zone: ISO8601 basic or extended format	Text	x xx xxx	-08, +0530, +0000 -0800, +0000 -08:00, +00:00
X	Time Zone: ISO8601 basic or extended format .The UTC indicator "Z" is used when local time offset is 0.	Text	X XX XXX	-08, +0530, Z -0800, Z -08:00, Z
l	ISO8601 date/time	Text	l	2006-10-07T12:06:56.568+01:00
'	escape for text	Delimiter	'	'Date='
"	single quote	Literal	"	'o'clock'

**Table 40 Symbols in the dfdl:calendarPattern Property**

The count of pattern letters determines the format as indicated in the table.

When numeric fields abut one another directly, with no intervening delimiter characters, they constitute a run of abutting numeric fields. Such runs are parsed specially as described at [ICUDateTime].

[ICUDateTime]

▼  
The maximum number of "S" symbols that may appear in the pattern is implementation-defined, but must be at least three. The stored accuracy for fractional seconds is also implementation-defined, but must be at least millisecond accuracy. When the number of "S" symbols in a pattern exceeds the supported accuracy, excess fractional seconds are truncated from the right (not

Deleted: ¶  
¶

rounded) when parsing, and zeros are added to the right when unparsing. For example, a DFDL processor allows up to six "S" symbols and has millisecond accuracy; for pattern "ss.SSSSSS", data "12.345678" would be parsed into infoset xs:time "00:00:12:345", which would be unparsed into data "12.345000".

Unlike other fields, fractional seconds, "S", are padded on the right with zero.

It is a processing error if seconds appear in that part of the SimpleContent region that represents a time zone.

The count of pattern letters determines the format as indicated in the table.

If dfdl:representation is text, any characters in the pattern that are not in the ranges of '[a'..'z]' and '[A'..'Z]' will be treated as quoted text. For instance, characters like ':', '!', ' ', '#', and '@' will appear in the formatted output even if they are not embraced within single quotes. The single quote is used to 'escape' letters. Two single quotes in a row, whether inside or outside a quoted sequence, represent a 'real' single quote.

If dfdl:representation is binary, then the pattern can contain only characters and symbols that always result in the presentation of digits.

The symbols 'z', 'zz', and 'zzz' have identical meaning, as do 'Z', 'ZZ', and 'ZZZ'.

The 'l' symbol must not be used with any other symbol with the exception of 'escape for text'. It represents calendar formats that match those defined in the restricted profile of the ISO 8601 standard proposed by the W3C at <http://www.w3.org/TR/NOTE-datetime>. The formats are referred to as 'granularities'.

- xs:dateTime. When parsing, the data must match one of the granularities. When unparsing, the fullest granularity is used.
- xs:date. When parsing, the data must match one of the date-only granularities. When unparsing, the fullest date-only granularity is used.
- xs:time. When parsing, the data must match only the time components of one of the granularities that contains time components. When unparsing, the time components of the fullest granularity are used. The literal 'T' character is not expected in the data when parsing and is not output when unparsing.
- The number of fractional second digits supported is is the same as for the "S" fractional seconds specifier described above.
- The omission of time zone from the input data when the type is xs:dateTime or xs:time is not a processing error. If that occurs then the time zone is obtained from the calendarTimeZone property.
- When unparsing and the time zone is UTC, the time zone is output as '+00:00'.

When parsing, for any pattern that omits components the values for the omitted components are supplied from the Unix epoch 1970-01-01T00:00:00.000.<sup>26</sup>

When unparsing, and the pattern contains a formatting symbol that requires a component of the date/time and the infoset value does not contain that component, it is a processing error.

When parsing a calendar element with a packed decimal representation then the nibbles from the data are converted to text digits without any trimming of leading or trailing zeros, and the result is then matched against the pattern according to the usual rules.

When unparsing, if a time zone symbol is not available for a particular time zone, a fallback may be used as defined in [\[ICUDateTime\]](#).

### 13.11.2 The dfdl:calendarCheckPolicy Property

The differences in behavior between 'strict' and 'lax' for this property can be subtle. Both are quite lenient in enforcement of many variations in format, with the 'lax' value adding additional

**Deleted:** Any number of fractional seconds "S" may be specified in the pattern and accepted by implementations, but an implementation is free to represent a limited number of fractional seconds internally. Excess fractional seconds are truncated, not rounded up. At least millisecond accuracy must be implemented. ...

**Deleted:** implementation-defined but must be at least millisecond accuracy...

<sup>26</sup> Note that DFDL does not support a pure month or day or year, as it does not support the XSD simple types xs:gMonth, xs:gDay, and xs:gYear.



tolerance of more format variations to those already allowed by the 'strict' value.

1. Lenient parsing behaviour when in 'strict' mode:
  - a. Case insensitive matching for text fields
  - b. MMM, MMMM, and MMMMM all accept either short or long form of Month
  - c. E, EE, EEE, EEEE, EEEEE, and EEEEEEE all accept either abbreviated, full, narrow and short forms of Day of Week
  - d. Accepts truncated leftmost numeric field (eg, pattern "HHmmss" allows "123456" (12:34:56) and "23456" (2:34:56) but not "3456")
2. Additional lenient parsing behaviour when in 'lax' mode:
  - a. Values outside valid ranges are normalized (eg, "March 32 1996" is treated as "April 1 1996")
  - b. Ignoring a trailing dot after a non-numeric field
  - c. Leading and trailing whitespace in the data but not in the pattern is accepted
  - d. Whitespace in the pattern can be missing in the data
  - e. Partial matching on literal strings. E.g., data "20130621d" allowed for pattern "yyyyMMdd'date' "

### 13.12 Properties Specific to Calendar with Text Representation

Property Name	Description
textCalendarJustification	<p>Enum</p> <p>Valid values 'left', 'right', 'center'</p> <p>Controls how the data is padded or trimmed on parsing and unparsing.</p> <p>Behavior as for dfdl:textStringJustification.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
textCalendarPadCharacter	<p>DFDL String Literal</p> <p>The value that is used when padding or trimming calendar elements. The value can be a single character or a single byte. If a character, then it can be specified using a literal character or using DFDL entities.</p> <p>If a byte, then it must be specified using a single byte value entity</p> <p>If a pad character is specified when dfdl:lengthUnits is 'bytes' then the pad character must be a single-byte character.</p> <p>If a pad byte is specified when dfdl:lengthUnits is 'characters' then</p> <ul style="list-style-type: none"> <li>the encoding must be a fixed-width encoding</li> <li>padding and trimming must be applied using a sequence of N pad bytes, where N is the width of a character in the fixed-width encoding.</li> </ul> <p>The string literal value is restricted in the same way as described in "Pad Character Restrictions" in the description of the dfdl:textStringPadCharacter property.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>

**Table 41 Properties Specific to Calendar with Text Representation**

### 13.13 Properties Specific to Calendar with Binary Representation

Property Name	Description
binaryCalendarRep	<p>Enum</p> <p>Valid values are 'packed', 'bcd', 'ibm4690Packed', 'binarySeconds', 'binaryMilliseconds'</p> <p>For all values, the dfdl:byteOrder property is used to determine the numeric significance of the bytes making up the representation.</p> <p>'packed' means represented as an IBM 390 packed decimal. Each byte contains two decimal digits, except for the rightmost byte, which contains a sign to the right of a decimal digit. The digits are interpreted according to the dfdl:calendarPattern property. Property dfdl:binaryPackedSignCodes is applicable.</p> <p>'bcd' means represented as a binary coded decimal with two digits per byte. The digits are interpreted according to the dfdl:calendarPattern property</p> <p>'ibm4690Packed' means represented as a variant of packed format as described in property dfdl:binaryNumberRep. The digits are interpreted according to the dfdl:calendarPattern property.</p> <p>For all packed decimals, property dfdl:binaryNumberCheckPolicy is applicable.</p> <p>For all these packed decimals, dfdl:calendarPattern can contain only characters and symbols that always result in the presentation of digits. It is a Schema Definition Error otherwise. This implies that property dfdl:calendarPatternKind must be 'explicit' because the default patterns for 'implicit' contain non-numeric characters. It is a Schema Definition Error otherwise.</p> <p>See Section 13.7 <a href="#">Properties Specific to Number with Binary Representation</a>.</p> <p>Note also that a virtual decimal point for the boundary between seconds and fractional seconds is implied from the pattern at the boundary of 's' and 'S', i.e., where the substring 'sS' appears in the pattern.</p> <p>'binarySeconds' means represented as binary xs:int, that is, as a 4 byte signed integer that is the number of seconds from the epoch (positive or negative). It is a Schema Definition Error if there is a specified length not equivalent to 4 bytes.</p> <p>'binaryMilliseconds' means represented as binary xs:long, that is, as an 8 byte signed integer that is the number of milliseconds from the epoch (positive or negative). It is a Schema Definition Error if there is a specified length not equivalent to 8 bytes.</p> <p>binarySeconds and binaryMilliseconds may only be used when the type is xs:dateTime. (It is a Schema Definition Error otherwise.)</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
binaryCalendarEpoch	<p>DateTime</p> <p>Used when dfdl:binaryCalendarRep is 'binarySeconds' or 'binaryMilliseconds'</p>

**Deleted:** Properties Specific to Numbers with Binary Representation...

	<p>The epoch from which to calculate dates and times.</p> <p>If the time zone is omitted from the value, then UTC is used.</p> <p>Annotation: dfdl:element, dfdl:simpleType</p>
--	---

**Table 42 Properties Specific to Calendar with Binary Representation**

Examples of packed decimal format calendars for December 14, 1923 and dfdl:calendarPattern of 'MMddyy' would be:

packed: (hexadecimal) 01 21 42 3C,

bcd: (hexadecimal) 12 14 23

ibm4690Packed: (hexadecimal) 12 14 23

The 'C' nibble at the end of the 'packed' representation is a sign nibble, and the leading 0 nibble is just to align to a byte boundary..

### 13.14 Properties Specific to Opaque Types (xs:hexBinary)

There are no properties specific to opaque types

### 13.15 Nil Value Processing

Sometimes it is desirable to represent an unused element, place-holder for unknown information, or inapplicable information *explicitly* with an element, rather than by the lack of an element.

For example, it may be desirable to represent a sparsely populated array of data using a distinguished nil element to fill the locations where data is absent, thereby preserving the position for the elements that are present.

As another example, it may be desirable to represent an unused simple element by a value which is not conformant to the logical type of the element.

Such cases can be represented using the DFDL nil mechanism which is based on the XML Schema nil mechanism. DFDL provides what are commonly called "in-band" nil values by way of dfdl:nilKind 'logicalValue', and also provides for two kinds of literal indicators of nil through dfdl:nilKind 'literalValue' and dfdl:nilKind 'literalCharacter'. Nil processing is used when the XSDL 'nillable' property of an element is true.

DFDL allows elements of complex type to be nillable. However, to avoid the concept of a complex element having a value, which does not exist in DFDL, the only permissible nil value is the empty string, represented by the DFDL %ES; entity and using dfdl:nilKind 'literalValue'.

On parsing, an element occurrence is nil if the element has XSDL nillable 'true' and the data is a nil representation as defined in section 9.2.1. Specifically:

1. When dfdl:nilKind is 'literalValue', the **NilLiteralValue** region of the data stream matches any of the dfdl:nilValue values.
2. When dfdl:nilKind is 'literalCharacter', all characters in the **NilLiteralCharacters** region of the data stream match the dfdl:nilValue character.
3. When dfdl:nilKind is 'logicalValue', the data contains a normal representation, and the **NilLogicalValue** region of the data stream, converted to the element's logical type, matches any of the dfdl:nilValue values.

For dfdl:nilKind 'literalValue' or 'literalCharacter':

- Determination of whether the data is a nil representation for a literal nil happens first before any consideration of whether the representation is the empty, normal, or absent representations.
- Property dfdl:nilValueDelimiterPolicy controls whether matching one of the nil values also involves matching the initiator or terminator specified by the element. This gives control over whether a nil indicator may or may not also require the delimiters that a normal data element requires.

On unparsing, an element is nil if XSDL nillable is 'true' AND the element information item in the augmented infoset has the **[nilled]** member as true, in which case what is output to the data stream is one of the following:

1. When dfdl:nilKind is 'logicalValue' then the first value of dfdl:nilValue converted to the physical representation is output as the **NilLogicalValue** region.
2. When dfdl:nilKind is 'literalValue' then the first value of dfdl:nilValue is output as the **NilLiteralValue** region.
3. When dfdl:nilKind is 'literalCharacter' then the character from dfdl:nilValue, repeated to the needed length, is output as the **NilLiteralCharacters** region.

For dfdl:nilKind 'literalValue' or 'literalCharacter' then dfdl:nilValueDelimiterPolicy determines whether any initiator or terminator also appear surrounding the literal nil in the output data.

### 13.16 Properties for Nillable Elements

These properties are used when the XSDL 'nillable' property of an element is 'true', and they control when and how the representation data are interpreted as having the logical meaning 'nil'.

Property Name	Description
nilKind	<p>Enum</p> <p>Valid values 'literalValue', 'logicalValue', 'literalCharacter'.</p> <p>Used when XSDL nillable is 'true'.</p> <p>Specifies how dfdl:nilValue is interpreted to represent the nil value in the data stream.</p> <p>If 'literalCharacter' then dfdl:nilValue specifies a single character or a single byte that, when repeated to the length of the element, is the nil value. 'literalCharacter' may only be specified for fixed-length elements, otherwise it is a Schema Definition Error..</p> <p>If 'literalValue' then dfdl:nilValue specifies a list of DFDL literal strings that are the possible representations for nil.</p> <p>If 'logicalValue' then dfdl:nilValue specifies a list of logical values that are the possible logical values for nil.</p> <p>Complex elements can be nillable, but dfdl:nilKind can only be 'literalValue' and dfdl:nilValue must be "%ES;". It is a Schema Definition Error otherwise.</p> <p>Annotation: dfdl:element</p>
nilValue	<p>List of DFDL String Literals, List of Logical Values, DFDL String Literal</p> <p>Specifies the text strings that are the possible literal or logical nil values of the element.</p> <p>If dfdl:nilKind is 'literalValue' then dfdl:nilValue specifies a whitespace separated list of DFDL literal strings that are the possible representations for nil. On parsing the element value is nil if the trimmed data matches one of the string literals in the list. On unparsing if the element value is nil the first string literal in the list is output.</p> <p>If dfdl:nilKind is 'logicalValue' then dfdl:nilValue specifies a whitespace separated list of logical values that are the possible logical values for nil. On parsing the element value is nil if the data, converted to its logical type, matches any of the logical values in the list. On unparsing if the element value is nil, the first</p>

Deleted: <#>¶

Deleted: , that is dfdl:lengthKind 'implicit' and 'explicit' when dfdl:length is not a DFDL expression

Deleted: (simpleType)

	<p>value from the list is converted to its physical representation and output.</p> <p>If <code>dfdl:nilKind</code> is 'literalCharacter' then <code>dfdl:nilValue</code> specifies a single character or byte that, when repeated to the length of the element, is the nil representation. If a character, then it can be specified using a literal character or using DFDL entities. If a character is specified when <code>dfdl:lengthUnits</code> is 'bytes' then the <code>dfdl:nilValue</code> must be a single-byte character. To specify a byte, it must be specified using a single <code>"%#r;"</code> entity. If a byte is specified when <code>dfdl:lengthUnits</code> is 'characters' then the <code>dfdl:encoding</code> must be a fixed-width encoding.</p> <p>On parsing, the element value is nil if all characters in the untrimmed data content match the <code>dfdl:nilValue</code> character. On unparsing, if the element value is nil the <code>dfdl:nilValue</code> character is output to the needed length.</p> <p>There are restrictions on the string literal syntax of <code>dfdl:nilValue</code>.</p> <p>When <code>dfdl:nilKind</code> is literalValue and text representation:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed</li> <li>• The DFDL byte value entity ( <code>%#rXX;</code> ) is allowed</li> <li>• DFDL Character classes NL, WSP, WSP+, WSP*, and ES are allowed.</li> </ul> <p>When <code>dfdl:nilKind</code> is literal value and binary representation:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed</li> <li>• The DFDL byte value entity ( <code>%#rXX;</code> ) is allowed</li> <li>• DFDL Character class ES is allowed.</li> <li>• Other DFDL Character classes NL, WSP, WSP+, and WSP*, are not allowed.</li> </ul> <p>When <code>dfdl:nilKind</code> is literalCharacter and text representation:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed</li> <li>• The DFDL byte value entity ( <code>%#rXX;</code> ) is allowed.</li> <li>• DFDL Character classes NL, WSP, WSP+, WSP*, and ES are not allowed.</li> </ul> <p>When <code>dfdl:nilKind</code> is literalCharacter and binary representation:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed</li> <li>• The DFDL byte value entity ( <code>%#rXX;</code> ) is allowed</li> <li>• DFDL Character classes NL, WSP, WSP+, WSP*, and ES are not allowed.</li> </ul> <p><code>dfdl:nilValue</code> is sensitive to <code>dfdl:ignoreCase</code> when <code>dfdl:nilKind</code> is 'literalValue' or 'logicalValue', but not when <code>dfdl:nilKind</code> is 'literalCharacter'</p> <p>Complex elements can be nillable, but <code>dfdl:nilKind</code> can only be 'literalValue' and <code>dfdl:nilValue</code> must be <code>"%ES;"</code>. It is a Schema Definition Error otherwise.</p> <p>Annotation: <code>dfdl:element</code></p>
<code>nilValueDelimiterPolicy</code>	<p>Enum</p> <p>Valid values are 'none', 'initiator', 'terminator' or 'both'.</p>

Deleted: (simpleType)

	<p>Indicates that when the value nil is represented, an initiator (if one is defined), a terminator (if one is defined), both an initiator and a terminator (if defined) or neither must be present.</p> <p>Ignored if both dfdl:initiator and dfdl:terminator are "" (empty string).</p> <p>Ignored if dfdl:nilKind is set to 'logicalValue' In this case the DFDL processor treats a nil representation like any other representation of the element in that it expects delimiters when parsing, outputs them when unparsing.</p> <p>'initiator' indicates that, on parsing, the dfdl:initiator followed by a dfdl:nilValue indicates that a nil representation is present. It also indicates that on unparsing when the logical value is nil that the dfdl:initiator will be output followed by the first dfdl:nilValue.</p> <p>'terminator' indicates that, on parsing, a dfdl:nilValue followed by the dfdl:terminator indicates that a nil representation is present. It also indicates that on unparsing when the logical value is nil the first dfdl:nilValue followed by the dfdl:terminator will be output.</p> <p>'both' indicates that, on parsing, both the dfdl:initiator and dfdl:terminator must be present with a dfdl:nilValue to indicate that a nil representation is present. On unparsing the dfdl:initiator followed by the first dfdl:nilValue, followed by the dfdl:terminator will be output.</p> <p>'none' indicates that a dfdl:nilValue without any dfdl:initiator or dfdl:terminator indicates that a nil representation is present. On unparsing the first dfdl:nilValue is output without any dfdl:initiator or dfdl:terminator.</p> <p><u>The value of dfdl:nilValueDelimiterPolicy should only be checked if there is a dfdl:initiator or dfdl:terminator in scope. If so, and dfdl:nilValueDelimiterPolicy is not set, it is a Schema Definition Error. If dfdl:initiator is not "" and dfdl:terminator is "" and dfdl:nilValueDelimiterPolicy is 'terminator' it is a Schema Definition Error. If dfdl:terminator is not "" and dfdl:initiator is "" and dfdl:nilValueDelimiterPolicy is 'initiator' it is a Schema Definition Error. It is not an SDE if dfdl:nilValueDelimiterPolicy is 'both' and one or both of dfdl:initiator and dfdl:terminator is "".</u></p> <p><u>This is to accommodate the common use of setting 'both' as a schema-wide setting.</u></p> <p>It is a Schema Definition Error if dfdl:nilValueDelimiterPolicy is set to 'none' or 'terminator' when the parent xs:sequence has dfdl:initiatedContent 'yes'.</p> <p>Annotation: dfdl:element</p>
useNilForDefault	<p>Enum</p> <p>Valid values are 'yes', 'no'</p> <p>When the conditions for applying a simple element default are satisfied, this property controls whether to set the Infoset item <b>[nilled]</b> boolean member, or to use the XSDL default or fixed properties to obtain a data value.</p> <p>This property has precedence over the XSDL default and XSDL fixed properties. It is only used, and must be defined, if the XSDL nillable property is 'true'.</p> <p>Defaulting occurs as described in section 9.4 Element Defaults with nil as the default value. The dfdl:nilValue property must</p>

Deleted: (simpleType)

	specify at least one nil value otherwise it is a Schema Definition Error. The dfdl:nilKind property may be any of its values. Annotation: dfdl:element (simpleType)
--	--

**Table 43 Properties for Nillable Elements**

The DFDL element defaults processing uses XSDL default, XSDL fixed or dfdl:useNilForDefault to provide a default value. See section 9.4 Element Defaults for a full description.

14. Sequence Groups

The following properties are specific to sequences.

Property Name	Description
sequenceKind	Enum Valid values are 'ordered', 'unordered' When 'ordered', this property means that the contained items of the sequence will be encountered in the same order that they appear in the schema, which is called schema-definition-order. When 'unordered', this property means that the items of the sequence will be encountered in any order. Repeating occurrences of the same element do not need to be contiguous. The children of an unordered sequence MUST be xs:element otherwise it is a Schema Definition Error. Annotation: dfdl:sequence, dfdl:group (sequence)
initiatedContent	Enum Valid values are 'yes', 'no' When 'yes' indicates that all the children of the sequence are initiated. It is a Schema Definition Error if any children have their dfdl:initiator property set to the empty string. If the child is optional then it is <u>known to exist</u> when its initiator has been found. Any subsequent error parsing the child will not cause the parser to backtrack to try other alternatives. When 'no', the children of the sequence may have their dfdl:initiator property set to the empty string. Annotation: dfdl:sequence, dfdl:choice, dfdl:group

Deleted:

Deleted:

Deleted: deemed to have been found

Table 44 Properties for Sequence Groups

A sequence can have a dfdl:initiator and/or a dfdl:terminator as described earlier.

14.1 Empty Sequences

A sequence having no children is syntactically legal in DFDL. In the data stream, such a sequence can have non-zero length **LeftFraming** and **RightFraming** regions, but the SequenceContent region in between must be empty. It is a processing error if the SequenceContent region of an empty sequence has non-zero length when parsing.

XML schema does not define an empty sequence that is the content model of a complex type definition as effective content so any DFDL annotations on such a construct would be ignored. It is a Schema Definition Error if the empty sequence is the content model of a complex type, or if a complex type has nothing in its content model at all.

A hidden group reference is indicated in DFDL using an empty sequence such as <xs:sequence dfdl:hiddenGroupRef="QName"/>. To XML Schema this is an empty sequence group; hence it is an SDE if this appears as the model group of a complex type. Otherwise this is not considered an empty sequence, but a group reference.



## 14.2 Sequence Groups with Separators

Additional properties apply to sequence groups that use text delimiters to separate one occurrence of a member of the group from the next. Such a delimiter is called a separator. DFDL provides several properties that control the parsing and writing of separators, and satisfy the requirement to model sequences where:

1. A separator has alternative potential representations in the data.
2. A separator is placed before, after, or between occurrences in the data.
3. Separators are used to indicate the position of occurrences in the data

These requirements are addressed by the properties `dfdl:separator`, `dfdl:separatorPosition` and `dfdl:separatorSuppressionPolicy`, as described below.

These properties combine to define the syntax for a sequence group with `dfdl:sequenceKind` 'ordered'. Not all combinations of the properties will give rise to consistent syntax, so some combinations are disallowed and will give rise to a Schema Definition Error.

In some sequences, the presence of separators alone is enough to establish occurrences within the sequence. Such a sequence is called a *positional* sequence.

**Positional sequence** - Each occurrence in the sequence can be identified by its position in the data. Typically, the components of such a sequence do not have an initiator. In some such sequences, the separators for optional zero-length occurrences may or must be omitted when at the end of the group. In DFDL, a sequence is considered positional if it contains only required elements and/or optional and array elements that have `dfdl:occursCountKind` 'implicit', 'fixed' or 'expression', and it has `dfdl:separatorSuppressionPolicy` 'never', 'trailingEmptyStrict' or 'trailingEmpty'.

**Non-positional sequence** - Occurrences in the sequence cannot be identified by their position in the data alone. Often the components of such a sequence have an initiator. Such sequences sometimes allow the separator to be omitted for optional zero-length occurrences anywhere in the sequence. Speculative parsing might need to be employed by the parser to identify each occurrence. In DFDL, a sequence is non-positional if it contains any optional or array elements that have `dfdl:occursCountKind` 'parsed' or 'stopValue', and/or it has `dfdl:separatorSuppressionPolicy` 'anyEmpty'.

Property Name	Description
separator	<p>List of DFDL String Literals or DFDL Expression</p> <p>Specifies a whitespace separated list of alternative DFDL String Literals that are the possible separators for the sequence. <u>Separators occur in the data either before, between or after all occurrences of the elements or groups that are the children of the sequence, in accordance with <code>dfdl:separatorPosition</code> and <code>dfdl:separatorSuppressionPolicy</code>. Elements with <code>dfdl:inputValueCalc</code> have no representation in the data stream, and so never have an associated separator.</u></p> <p>This property can be computed by way of an expression which returns a string of whitespace separated values. The expression must not contain forward references to elements which have not yet been processed. It is a Schema Definition Error if the expression returns an empty string.</p> <p>This property can be used to determine the length of an element as described in Section <a href="#">12.3.2</a> <code>dfdl:lengthKind</code> 'delimited'.</p> <p>Each string literal in the list, whether apparent in the schema, or returned as the value of an expression, is</p>

**Deleted:** A positional sequence can be modelled by setting `dfdl:separatorSuppressionPolicy` to 'never', 'trailingEmptyStrict' or 'trailingEmpty'.¶

**Deleted:** Typically

**Deleted:** is

**Deleted:** A non-positional sequence can be modelled by setting `dfdl:separatorSuppressionPolicy` to 'anyEmpty'.

**Deleted:** ¶

**Deleted:** Separators occur in the data either before, between or after all occurrences of the elements or groups that are the children of the sequence.

	<p>restricted to allow only certain kinds of DFDL String Literal syntax:</p> <ul style="list-style-type: none"> <li>• DFDL character entities are allowed.</li> <li>• DFDL Byte Value entities ( <code>%#rXX;</code> ) are allowed.</li> <li>• DFDL Character Class ES is not allowed.</li> <li>• DFDL Character Classes NL, WSP, WSP+, and WSP* are allowed.</li> <li>• The WSP* entity cannot appear on its own as one of the string literals in the list when determining the length of a component by scanning for delimiters.</li> </ul> <p>If the above rules are not followed it is a Schema Definition Error.</p> <p>The <b>Separator</b>, <b>PrefixSeparator</b> and <b>PostfixSeparator</b> regions contain one of the strings specified by the <code>dfdl:separator</code> property. When this property has "" (empty string) as its value then the separator region is of length zero.</p> <p>When parsing, the list of values is processed in a greedy manner, meaning it takes all the separators, that is, each of the string literals in the whitespace separated list, and matches them each against the data. The separator with the longest match is the one that is selected as having been 'found'. Once a matching separator is found, no other matches will be subsequently attempted (ie, there is no backtracking).</p> <p>On unparsing the first separator in the list is used as the separator.</p> <p>If a child element uses an escape scheme, then the escape scheme also applies to any separator; hence, if the separator appears within the element value, it will be escaped.</p> <p>If <code>dfdl:ignoreCase</code> is 'yes' then the case of the string is ignored by the parser.</p> <p>Annotation: <code>dfdl:sequence</code>, <code>dfdl:group</code> (sequence)</p>
separatorPosition	<p>Enum</p> <p>Valid values 'infix', 'prefix', 'postfix'</p> <p>'infix' means the separator occurs between the elements in the <b>Separator</b> grammar region.</p> <p>'prefix' means the separator occurs before each element in the <b>Separator</b> grammar region and the <b>PrefixSeparator</b> grammar region.</p> <p>'postfix' means the separator occurs after each element in the <b>Separator</b> grammar region and the <b>PostfixSeparator</b> grammar region.</p> <p>Annotation: <code>dfdl:sequence</code>, <code>dfdl:group</code> (sequence).</p>
separatorSuppressionPolicy	<p>Enum</p> <p>Valid values 'never', 'anyEmpty', 'trailingEmpty', 'trailingEmptyStrict'</p>

	<p>Only applicable if dfdl:separator is not "" (empty string) and dfdl:sequenceKind is 'ordered'.</p> <p>Controls the circumstances when separators are expected in the data when parsing, or generated when unparsing, if an optional element occurrence or a group has a zero-length representation.</p> <p>See section 14.2.1 Separators and Suppression.</p> <p>When dfdl:sequenceKind is 'unordered' then 'anyEmpty' is implied.</p> <p>Annotation: dfdl:sequence, dfdl:group (sequence)</p>
--	---

**Table 45 Properties for Sequence Groups with Separators****14.2.1 Separators and Suppression**

When parsing a sequence group that specifies a separator, the number of occurrences and separators that are expected in the data stream for a child (element or group) depends on several factors:

- Whether element occurrences are optional or required
- Whether the occurrences (element or group) have a zero-length representation
- Whether occurrences (element or group) are trailing
- Whether the sequence is positional
- The dfdl:occursCountKind of the element

Where to expect a separator for optional content of zero-length is controlled by property dfdl:separatorSuppressionPolicy.

separatorSuppressionPolicy	Implications
never	Positional sequence where all occurrences <b>MUST</b> be found in the data, along with their associated separator.
trailingEmptyStrict	Positional sequence where <i>trailing occurrences</i> that have zero length representation <b>MUST</b> be omitted from the data, along with their associated separator.
trailingEmpty	Positional sequence where <i>trailing occurrences</i> that have zero length representation <b>MAY</b> be omitted from the data, along with their associated separator.
anyEmpty	Non-positional sequence where any occurrences that have zero length representation <b>MAY</b> be omitted from the data, along with their associated separator. It must be possible for speculative parsing to identify which elements are present.

**Table 46 Sequence groups and separator suppression**

The following are definitions for terminology used in this section:

**Potentially trailing element** – An array or optional element describes an occurrence that is said to be *potentially trailing* if the element is capable of having a zero length representation and is followed in its enclosing group definition by only these kinds of schema components:

1. calculated elements (those having dfdl:inputValueCalc)
2. additional potentially trailing elements
3. potentially trailing groups

Intuitively, the array or optional element occurrence could be last.

**Potentially trailing group** – A group is said to be *potentially trailing* if the group has no framing and contains only potentially trailing element declarations/references, or recursively similar

**Deleted:** When parsing a sequence group that specifies a separator, the number of occurrences and separators that are expected in the data stream for a child element depends on several factors:

**Deleted:** ¶  
 Whether occurrences are optional or required¶  
 Whether the occurrences have a zero-length representation¶  
 Whether occurrences are trailing¶  
 Whether the sequence is positional¶  
 The dfdl:occursCountKind of the element¶

sequence or choice groups, and is followed in its enclosing group definition by only additional potentially trailing elements or potentially trailing groups.

**Trailing or Actually Trailing** – An element occurrence or group occurrence in the data is said to be *actually trailing* if it is potentially trailing and has zero-length representation and is not followed in the data by any other non-zero length element occurrence or group occurrence limited by the end of the enclosing sequence group.

In the sections that follow, it is important to remember that the `dfdl:separatorSuppressionPolicy` property is carried on the sequence, while the XSDL `minOccurs`, XSDL `maxOccurs` and `dfdl:occursCountKind` properties are carried on an *element* in that sequence.

#### 14.2.2 Parsing Sequence Groups with Separators

Parsing child elements is described first. Parsing for child groups is described in Section 14.2.2.3.

When an element is required and is not an array then one occurrence is always expected along with its separator. The `dfdl:separatorSuppressionPolicy` of the sequence has no effect (nothing is eligible for suppression). Otherwise the behaviour is dependent on `dfdl:occursCountKind`.

When `dfdl:occursCountKind` is 'fixed' then XSDL `minOccurs` must equal `maxOccurs` and that many occurrences are always expected along with their separators. The `dfdl:separatorSuppressionPolicy` of the sequence has no effect (nothing is eligible for suppression).

When `dfdl:occursCountKind` is 'expression' the number of occurrences is given by `dfdl:occursCount` and exactly that many occurrences are always expected along with their separators. The `dfdl:separatorSuppressionPolicy` of the sequence has no effect (nothing is eligible for suppression).

When `dfdl:occursCountKind` is 'parsed' any number of occurrences and their separators are expected. The `dfdl:separatorSuppressionPolicy` of the sequence must be 'anyEmpty' and it is a Schema Definition Error otherwise.

When `dfdl:occursCountKind` is 'stopValue', any number of occurrences and their separators are expected followed by the stop value and its separator. The `dfdl:separatorSuppressionPolicy` of the sequence has no effect.

When `dfdl:occursCountKind` is 'implicit', between XSDL `minOccurs` and `maxOccurs` (inclusive) occurrences and their separators are expected, according to the `dfdl:separatorSuppressionPolicy` of the sequence.

The behaviour for 'implicit' is more fully expressed in matrix form. The cells in the matrix give the number of occurrences of element values that are expected in the data stream when parsing, for the different values of `dfdl:separatorSuppressionPolicy`. The number of occurrences also depends whether XSDL `maxOccurs` is unbounded or not, and the position of the element in the sequence. The number of separators can be inferred from this, taking into account `dfdl:separatorPosition`.

dfdl: separatorSuppressionPolicy	dfdl:occursCountKind 'implicit'					
	Potentially Trailing				Not Potentially Trailing	
	maxOccurs unbounded		maxOccurs bounded		maxOccurs unbounded	maxOccurs bounded
	Element not declared last	Element declared last	Element declared last or occurrence followed by end-of-group	Element not declared last and occurrence not followed by end-of-group		

**Deleted:** The `dfdl:separatorSuppressionPolicy` is not applicable and the implied behaviour is 'never'. ¶

**Deleted:** The `dfdl:separatorSuppressionPolicy` is not applicable and the implied behaviour is 'never'. ¶

**Deleted:** The `dfdl:separatorSuppressionPolicy` is not applicable and the implied behaviour is 'never'.

**Deleted:** The `dfdl:separatorSuppressionPolicy` is not applicable and the implied behaviour is 'anyEmpty'.

**Deleted:** is not applicable and the implied behaviour is 'anyEmpty'....

**Deleted:** The `dfdl:separatorSuppressionPolicy` is applicable and determines when separators are expected for optional zero-length occurrences.

never	Schema definition error					
trailingEmptyStrict		RepDef(min) [ ~ Rep(M < INF) ~ RepNonZero(1) ]	RepDef(min) [ ~ Rep(M < max - min) ~ RepNonZero(1) ]	RepDef(min) ~ Rep(max - min)	Schema definition error	RepDef(min) ~ Rep(max - min)
trailingEmpty						
anyEmpty		RepDef(min) ~ Rep(M < INF)	RepDef(min) ~ Rep(M <= max - min)		RepDef(min) ~ Rep(M < INF)	RepDef(min) ~ Rep(M <= max - min)

**Table 47 Separator Suppression for dfdl:occursCountKind 'implicit' when Parsing**

The notation in each cell uses the "~" symbol to mean "followed by" in the data stream. Square brackets surround things that are optional, as in they may or may not appear in the data stream.

The descriptions found in the cells of the matrix do not provide a parsing algorithm, but rather state declaratively a pattern that the data must match in order to be correctly parsed.

**RepDef(min)** is short for "representation" and "defaultable". It means XSDL minOccurs occurrences of nil, empty or normal representation<sup>27</sup>. These are required occurrences so default rules apply for empty representations. XSDL minOccurs may be 0, in which case there are no required occurrences.

**Rep(M)** means M occurrences of nil, empty, normal or absent representation. These are optional occurrences so default rules do not apply for empty representations.

**RepNonZero(1)** means an occurrence of a nil, empty or normal representation where such a representation does not have zero-length<sup>28</sup>. This is an optional occurrence so default rules do not apply.

A notation like **Rep(M <= max - min)** means that there are M occurrences, where M is some value between the values of the XSDL minOccurs and maxOccurs properties. When an unbounded number of occurrences is possible this is shown explicitly by **Rep(M < INF)**, INF meaning infinity or unbounded.

#### 14.2.2.1 Errors When the Sequence is Positional

In the matrix above we see that there are some cells where the combination of properties doesn't make sense, and a Schema Definition Error is raised. These occur when an element has dfdl:occursCountKind 'implicit' and XSDL maxOccurs 'unbounded', and dfdl:separatorSuppressionPolicy implies that the sequence is positional, specifically:

- If a sequence has dfdl:separatorSuppressionPolicy 'never';
- If a sequence has dfdl:separatorSuppressionPolicy 'trailingEmptyStrict' or 'trailingEmpty' and the element is not the last declaration in the sequence. (This avoids ambiguity about which element is being suppressed.)

#### 14.2.2.2 Example Parsing Scenarios

Consider the cell of the matrix above for the element in this DFDL schema fragment:

```
<xs:sequence dfdl:separator='|' dfdl:separatorPosition='infix'
  dfdl:separatorSuppressionPolicy='trailingEmptyStrict'>
  <xs:element name='a' type='xs:int' default='0'
    maxOccurs='5' minOccurs='0'
    dfdl:representation='text' dfdl:textNumberPattern='#0'
    dfdl:occursCountKind='implicit'
    dfdl:initiator='[' dfdl:terminator=']'
    dfdl:emptyValueDelimiterPolicy='both' />
```

<sup>27</sup> Absent representation implies processing error for 'implicit' when less than or equal to minOccurs.

<sup>28</sup> Absent representation always implies zero-length. Nil, empty, and normal representations can also be zero-length with the right combinations of properties. See Section 9.2.5 Zero-length Representation.



- 'trailingEmpty' – if the child group is potentially trailing, has zero-length and it is actually trailing, its separator MAY appear or not. Additional separators are not expected.
- 'trailingEmptyStrict' – if the child group is potentially trailing, has zero-length and it is actually trailing, its separator MAY NOT appear.
- 'anyEmpty' – if the child group has zero-length its separator MAY NOT appear.

#### 14.2.3 Unparsing Sequence Groups with Separators

Unparsing child elements is described first. Unparsing for child groups is described in Section 14.2.3.2.

When an element is required and is not an array then one occurrence is always output along with its separator. The dfdl:separatorSuppressionPolicy of the sequence has no effect (nothing is eligible for suppression).

Otherwise the behaviour is dependent on dfdl:occursCountKind.

When dfdl:occursCountKind is 'fixed' or 'expression' the occurrences in the augmented Infoset are always output along with their separators. The dfdl:separatorSuppressionPolicy of the sequence has no effect (nothing is eligible for suppression).

When dfdl:occursCountKind is 'parsed' non zero-length occurrences in the augmented Infoset are output along with their separators. The dfdl:separatorSuppressionPolicy of the sequence must be 'anyEmpty' and it is a Schema Definition Error otherwise.

When dfdl:occursCountKind is 'stopValue' the occurrences in the augmented Infoset are output along with their separators followed by the stop value and its separator, according to the dfdl:separatorSuppressionPolicy of the sequence.

When dfdl:occursCountKind is 'implicit' the occurrences in the augmented Infoset are output along with their separators, according to the dfdl:separatorSuppressionPolicy of the sequence.

The behaviour for 'implicit' is more fully expressed in matrix form. The cells in the matrix give the number of occurrences of element values that are output to the data stream when unparsing, for the different values of dfdl:separatorSuppressionPolicy. The number of occurrences also depends whether XSDL maxOccurs is unbounded or not, and the position of the element in the sequence. The number of separators output can be inferred from this, taking into account dfdl:separatorPosition.

dfdl: separatorSuppressionPolicy	dfdl:occursCountKind 'implicit'					
	Potentially Trailing				Not Potentially Trailing	
	maxOccurs unbounded		maxOccurs bounded		maxOccurs unbounded	maxOccurs bounded
	Element not declared last	Element declared last	Element declared last or occurrence followed by end-of-group	Element not declared last and occurrence not followed by end-of-group		
never	Schema definition error		Unparse N occurrences ~ unparse (maxOccurs -- N) trailing zero-length occurrences		Schema definition error	Unparse N occurrences ~ unparse (maxOccurs -- N) trailing zero-length occurrences
trailingEmptyStrict	Unparse N occurrences (suppressing trailing zero-length occurrences)					
trailingEmpty						
anyEmpty	Unparse N occurrences (suppressing any optional zero-length occurrences)					

**Table 48 SeparatorSuppressions for dfdl:occursCountKind 'implicit'**

The notation in each cell uses the "~" symbol to mean "followed by" in the output data stream.

**Deleted:** The dfdl:separatorSuppressionPolicy is not applicable and the implied behaviour is 'never'.

**Deleted:** The dfdl:separatorSuppressionPolicy is not applicable and the implied behaviour is 'never'.

**Deleted:** The dfdl:separatorSuppressionPolicy is not applicable and the implied behaviour is 'anyEmpty'.

**Deleted:** non zero-length

**Deleted:** . The dfdl:separatorSuppressionPolicy is not applicable and the implied behaviour is 'anyEmpty'.

**Deleted:** . The dfdl:separatorSuppressionPolicy is applicable and helps determine whether optional zero length occurrences and their separators are output.

**N** stands for the number of elements in the augmented Infoset, which includes any defaults.

**unparse N occurrences** means output N unparsed infose items and associated separators.

**unparse(M) trailing zero length occurrences** means output M adjacent separators (according to dfdl:separatorPosition) as if separating M element occurrences.

**(suppressing trailing zero-length reps)** implies the unparser must look ahead into the infose and determine when the representations will be zero-length, and then identify those in trailing position. No separators are output corresponding to the trailing zero-length representations.

#### 14.2.3.1 Example Unparsing Scenarios

Consider the cell of the matrix above for the element in this DFDL schema fragment:

```
<xs:sequence dfdl:separator='|' dfdl:separatorPosition='infix'
  dfdl:separatorSuppressionPolicy='trailingEmpty'>
  <xs:element name='a' type='xs:int'
    maxOccurs='5' minOccurs='0'
    nillable='true'
    dfdl:representation='text' dfdl:textNumberPattern='#0'
    dfdl:occursCountKind='implicit'
    dfdl:initiator='[' dfdl:terminator=']'
    dfdl:emptyValueDelimiterPolicy='none'
    dfdl:nilKind='literalValue' dfdl:nilValue='%ES;'
    dfdl:nilValueDelimiterPolicy='none' />
</xs:sequence>
```

This example is similar to the one used above in the discussion of parsing with separator suppression. However, the element has no default value, the dfdl:emptyValueDelimiterPolicy has been removed, and the element is nillable. Element 'a' is clearly potentially trailing as it is declared last. The corresponding cell in the matrix above contains this description:

*unparse N occurrences (suppressing trailing zero length reps)*

Assume we are unparsing an infose containing five values: 1, 0, nil<sup>29</sup>, 4, nil. We will unparse the five occurrences; however, the last value is nil, which has a representation of '%ES;' meaning empty-string, and dfdl:nilValueDelimiterPolicy is 'none' meaning no initiator or terminator is to appear in the data. Since we are suppressing trailing zero-length reps the unparse results in this output:

```
[1] | [0] | | [4]
```

We see here an example where if we reparsed the data we would not get back the same infose that we started from because the trailing empty value which is the representation of the nil value, is not represented in the output, and so will not cause an Infoset item with **[nilled]** true to be created in the infose when this data is parsed. To preserve the nil we would have to change the dfdl:nilValueDelimiterPolicy to 'both', and in that case the output would be:

```
[1] | [0] | | [4] | []
```

The nils now have explicit representation in the data, and will be recreated in the Infoset when parsing.

#### 14.2.3.2 Unparsing Child Groups within Separated Sequences

When a child of a sequence is a group then a separator is output depending on dfdl:separatorSuppressionPolicy and other factors:

- 'never' - the child group's associated separator is output
- 'trailingEmpty' or 'trailingEmptyStrict' – if the child group is potentially trailing, has zero-length and it is actually trailing, its separator is not output.
- 'anyEmpty' – if the child group has zero-length its separator is not output.

<sup>29</sup> An Infoset item value of nil means the Infoset item **[nilled]** member is true, and the **[dataValue]** member has no value. See Section 4.1.2 Element Information Items.



### 14.3 Unordered Sequence Groups

The occurrences of members of a sequence group with `dfdl:sequenceKind` of 'unordered' (hereafter referred to as an 'unordered sequence') may appear in the data in any order. Occurrences of the same member do not have to be contiguous. In the infoset, sequence groups are always in schema order, so a DFDL processor must sort the members of an unordered sequence into schema order when parsing. When unparsing, the infoset must already be in schema order, and the members of the sequence will be output in schema order.

#### 14.3.1 Restrictions for Unordered Sequences

It is a Schema Definition Error if any member of the unordered sequence is not an element declaration or an element reference.

It is a Schema Definition Error if a member of an unordered sequence is an optional element or an array element and its `dfdl:occursCountKind` property is not 'parsed'

It is a Schema Definition Error if two or more members of the unordered sequence have the same name and the same namespace (see post-processing transformation below)

It is a Schema Definition Error if an unordered sequence has no members.

#### 14.3.2 Parsing an Unordered Sequence

When parsing, the semantics of an unordered sequence are expressed by way of:

1. a source-to-source transformation of the sequence group definition, and
2. a post-processing transformation of the infoset .

An implementation may use any technique consistent with this semantic.

##### 14.3.2.1 Source-to-source Transformation

The source-to-source transformation turns the declaration of an unordered sequence into an ordered sequence group that contains a repeating choice. To ensure that the resulting schema is a valid DFDL schema, the choice group is wrapped in an array element.

The unordered sequence is transformed as follows:

- the `dfdl:sequenceKind` property of the unordered sequence is changed to "ordered"
- the content of the unordered sequence is replaced by a complex element ( the 'choice element' ) with the following properties:
  - XSDL `minOccurs`="0"
  - XSDL `maxOccurs`="unbounded"
  - `dfdl:lengthKind` "implicit"
  - `dfdl:occursCountKind` "parsed"
- the content of the choice element's complex type is a choice group with the following properties:
- `dfdl:choiceLengthKind` "implicit"
- The members of the unordered sequence become the members of the choice group, with their declaration order preserved.
- The XSDL `minOccurs` and `maxOccurs` properties on each member of the choice group are both set to 1.

Using the following example as an illustration:

```
<xs:sequence dfdl:sequenceKind="unordered" dfdl:separator=",">
  <xs:element name="a" type="xs:string" dfdl:initiator="A:" />
  <xs:element name="b" type="xs:int" minOccurs="0" dfdl:initiator="B:" />
  <xs:element name="c" type="xs:string" minOccurs="0" maxOccurs="10"
    dfdl:initiator="C:" />
</xs:sequence>
```

The above unordered sequence group is conceptually rewritten into the following ordered sequence group:

```

<xs:sequence dfdl:sequenceKind="ordered" dfdl:separator=",">
  <xs:element name="choiceElement" minOccurs="0" maxOccurs="unbounded"
    dfdl:occursCountKind="parsed">
    <xs:complexType>
      <xs:choice dfdl:choiceLengthKind="implicit">
        <xs:element name="a" type="xs:string" dfdl:initiator="A:" />
        <xs:element name="b" type="xs:int" dfdl:initiator="B:" />
        <xs:element name="c" type="xs:string" dfdl:initiator="C:" />
      </xs:choice>
    </xs:complexType>
  </xs:element>
</xs:sequence>

```

Processing then constructs a temporary info set for this ordered sequence group by parsing the data.

If a member element is found to have the empty representation then the parsing of that element must use the original value of XSDL minOccurs. In this example, element "b" has minOccurs "0" and if it is found with the empty representation then it must not be defaulted.

#### 14.3.2.2 Post-processing Transformation

Post-processing consists of the following steps:

1. Sort the temporary info set to produce the real info set
2. Check scalar elements and validate

##### Sort the Temporary Info set

The temporary info set is transformed into the info set conforming to the original unordered sequence. All members of the temporary info set having the same name and namespace as the first child of the unordered sequence are placed first, in the order in which they were parsed. This algorithm repeats for the second child of the unordered sequence and so on until all members of the temporary info set have been sorted into the schema declaration order of the original unordered sequence.

For the example above, the temporary info set is transformed into the info set corresponding to:

```

<xs:sequence>
  <xs:element name="a" type="xs:string" />
  <xs:element name="b" type="xs:int" minOccurs="0" />
  <xs:element name="c" type="xs:string" minOccurs="0" maxOccurs="10" />
</xs:sequence>

```

##### Check Scalar Elements and Validate

For each element in the unordered sequence having XSDL minOccurs "1" and maxOccurs "1", the number of occurrences is checked. Each such element must occur exactly once in the info set, else it is a processing error.

If validation is enabled, the DFDL processor validates the number of occurrences of each member of the unordered sequence against XSDL minOccurs and maxOccurs.

These checks are the same as those performed for an ordered sequence group. However, in an unordered sequence the checking of XSDL minOccurs and maxOccurs must be performed after the entire group has been parsed.

#### 14.3.3 Unparsing an Unordered Sequence

When unparsing, the behavior is exactly as if dfdl:sequenceKind is 'ordered'. The info set must be presented to the unparser in schema declaration order, and the members of the unordered sequence group are output in schema declaration order.

#### 14.4 Floating Elements

Elements within an ordered sequence can be designated as floating which means that they can appear in any position within the sequence.<sup>30</sup>

Property Name	Description
floating	<p>Enum</p> <p>Valid values are 'yes', 'no'</p> <p>Whether the occurrences of an element in an ordered sequence can appear out-of-order in the representation.</p> <p>When parsing, and dfdl:floating is 'yes', occurrences of the element may be encountered in the representation in many positions within its containing sequence. If present they are placed into the infoset in schema declaration order. If the element repeats, occurrences do not need to be contiguous in the representation.</p> <p>When parsing, and dfdl:floating is 'no', occurrences of the element must be in schema declaration order, and, if present, they are placed into the infoset in schema declaration order. It is a processing error if instances of the element are not encountered in schema declaration order.</p> <p>When unparsing, occurrences of the element are expected in the infoset in schema declaration order, and are output in the representation in schema declaration order. It is a processing error if occurrences of the element are not encountered in schema declaration order,</p> <p>It is a Schema Definition Error if an unordered sequence or a choice contains any element with dfdl:floating 'yes'.</p> <p>It is a Schema Definition Error if an ordered sequence contains any element with dfdl:floating 'yes' and also contains non-element component (such as a choice or sequence model group).</p> <p>It is a Schema Definition Error if an element with dfdl:floating 'yes' is an optional element or an array element and its dfdl:occursCountKind property is not 'parsed'</p> <p>It is a Schema Definition Error if two or more elements with dfdl:floating 'yes' in the same group have the same name and the same namespace.</p> <p>Annotation: dfdl:element</p>

**Table 49 Properties for Floating Elements**

An ordered sequence with floating components is similar to an unordered sequence except only the floating elements may be out of order.

Within an ordered sequence with floating components a non-floating array element must have its occurrences appearing contiguously, so any floating elements cannot appear in between occurrences of the array element.(In other words, property dfdl:floating 'yes' only makes a statement about the floating element, not about any other elements in the sequence.)

An ordered sequence of n element children with dfdl:floating 'yes' is equivalent to an unordered sequence with the same n element children with dfdl:floating 'no'.

A complex element with dfdl:floating 'yes' can have as its content model a sequence with elements that also have dfdl:floating 'yes'.

Every element in a sequence containing one or more floating elements is a point of uncertainty, similar to the way every element in an unordered sequence is a point of uncertainty.

<sup>30</sup>The NTE segment in the X12 EDI standard is an example of a floating element.

In resolving this point of uncertainty, a parser MUST look for the element defined at that position in the schema first and only if unsuccessful with parsing that element, the parser should subsequently attempt to parse the floating elements in the order they are defined in the schema. As soon as any such parse is successful this resolves the point of uncertainty.

#### 14.5 Hidden Groups

Some fields in the physical stream provide information about other fields in the stream and are not really part of the data. For example, a field could give the number of repeats in a following array. These fields may not be of interest to an application after the data has been parsed, and so may be removed from the Infoset on parsing by containing the element declarations for them within a hidden group. A hidden group allows elements to be defined that will not be added to the Infoset on parsing and will not be expected in the Infoset on unparsing.

```
<xs:element name="root">
  <xs:complexType>
    <xs:sequence>

      <xs:sequence>
        <xs:annotation><xs:appinfo source="http://www.ogf.org/dfdl/">
          <dfdl:sequence hiddenGroupRef="tns:hiddenRepeatCount">
            </xs:appinfo></xs:annotation>
          </xs:sequence>

          <xs:element name="arrayElement" type="xs:int"
            minOccurs="0" maxOccurs="unbounded"
            dfdl:occursCountKind="expression"
            dfdl:occursCount= "{../repeatCount}"
            dfdl:representation="binary" dfdl:lengthKind="implicit" />
        </xs:sequence>
      </xs:complexType>
    </xs:element>

    <xs:group name="hiddenRepeatCount" >
      <xs:sequence>
        <xs:element name="repeatCount" type="xs:int"
          dfdl:outputValueCalc="{count(../arrayElement)}"
          dfdl:representation="binary" dfdl:lengthKind="implicit" />
      </xs:sequence>
    </xs:group>
  </xs:sequence>
</xs:element>
```

An element contained within the extent of a hidden group is commonly called a hidden element. Hidden elements are referenced via path expressions using the same DFDL expression that would be used if they were not hidden.

Hidden elements can (typically will) contain the regular DFDL annotations to define their physical properties and on unparsing to set their value. They are processed using the same behavior as non-hidden elements.

When the `dfdl:hiddenGroupRef` property is specified on an `xs:sequence` schema component, the appearance of any other DFDL properties on that component is a Schema Definition Error. It is also a Schema Definition Error if the sequence is not empty.

It is a Schema Definition Error if the sequence is the only thing in the content model of a complex type definition.

It is a Schema Definition Error if `dfdl:hiddenGroupRef` appears on a `xs:group` reference, that is, unlike most format properties that apply to sequences, `dfdl:hiddenGroupRef` cannot be combined from a `xs:group` reference.

A hidden group may appear within another hidden group.

Property Name	Description
<code>hiddenGroupRef</code>	QName

	<p>Reference to a global model group definition. Elements within this model group will not be added to the Infoset, and are called hidden elements.</p> <p>The model group within the model group definition may be a xs:sequence or xs:choice</p> <p>It is a Schema Definition Error if the value is the empty string.</p> <p>It is not possible to place this property in scope on a dfdl:format annotation.</p> <p>Annotation: dfdl:sequence</p>
--	---

**Table 50 Properties for Hidden Groups**

When unparsing a hidden group, the behaviour is the same as when elements are missing from the infoset; that is, the default-values algorithm applies. The only difference is that if a required element does not have a default value or a dfdl:outputValueCalc then it is a Schema Definition Error instead of a processing error. Note that this can be checked statically.

When unparsing a hidden group, it is a processing error if an element information item is provided in the infoset for a hidden element.

## 15. Choice Groups

The following properties are specific to xs:choice.

Property Name	Description
choiceLengthKind	Enum Valid values are 'implicit', 'explicit' 'implicit' means the branches of the choice are not filled, so the ChoiceContent region is variable length depending on which branch appears. 'explicit' means that the branches of the choice are always filled to the fixed-length specified by dfdl:choiceLength, so the ChoiceContent region is fixed-length regardless of which branch appears. Annotation: dfdl:choice, dfdl:group (choice)
choiceLength	Integer Only used when dfdl:choiceLengthKind is 'explicit'. Specifies the length of the choice in bytes, so the ChoiceContent region is fixed-length regardless of which branch appears. A <b>ChoiceUnused</b> region is therefore possible which when unparsing is filled with dfdl:fillByte. Annotation: dfdl:choice, dfdl:group (choice)
initiatedContent	Enum Valid values are 'yes', 'no' When 'yes' indicates that all the branches of the choice are initiated. It is a Schema Definition Error if any children have their dfdl:initiator property set to the empty string. The branch is <b>known to exist</b> when its initiator has been found. Any subsequent error parsing the branch will not cause the parser to backtrack. When 'no', the branches of the choice may have their dfdl:initiator property set to the empty string. Annotation: dfdl:sequence, dfdl:choice, dfdl:group
choiceDispatchKey	DFDL Expression The expression must evaluate to an xs:string. <b>It is a Schema Definition Error if the expression returns an empty string.</b> <b>It is a Schema Definition Error if the expression contains forward references to elements which have not yet been processed.</b> This property is used only when parsing. The resultant string must match <b>one of</b> the dfdl:choiceBranchKey property values of one of the branches of the choice. This match is <i>case sensitive</i> . If so, it discriminates to that branch. The parser then goes straight to that branch, ignoring consideration of any other choice branches. No backtracking of this decision occurs if there is a subsequent processing error. It is a processing error if the value of the expression does not match <b>any</b> of the dfdl:choiceBranchKey property values for <b>any of</b> the branches. It is a Schema Definition Error if any choice branch does not specify a dfdl:choiceBranchKey in a choice that carries choiceDispatchKey. It is not possible to place this property in scope on a dfdl:format annotation. Annotation: dfdl:choice

**Deleted:** deemed to have been found

**Deleted:** which must not be the empty string.

**Deleted:** ¶

**Deleted:** one

**Deleted:** ies for

choiceBranchKey	<p><a href="#">List of DFDL String Literals</a></p> <p>This literal provides an alternate way to discriminate a choice to a branch. When the <code>dfdl:choiceDispatchKey</code> expression evaluates to a string matching <u>one of</u> this property's values, the choice is discriminated to this branch. The match is case <u>sensitive</u>.</p> <p>It is a Schema Definition Error if individual <code>dfdl:choiceBranchKey</code> values are not unique across all branches of a choice that carries <code>dfdl:choiceDispatchKey</code>.</p> <p>It is a Schema Definition Error if <code>dfdl:choiceBranchKey</code> is specified on a global element, or on a sequence or choice that is the child of a global group definition.</p> <p>It is a Schema Definition Error if any choice branch does not specify a <code>dfdl:choiceBranchKey</code> in a choice that carries <code>choiceDispatchKey</code>.</p> <p>Byte value entities are not allowed.</p> <p>Character classes are not allowed.</p> <p>This property is only used when parsing.</p> <p>It is not possible to place this property in scope on a <code>dfdl:format</code> annotation.</p> <p>Annotation: <code>dfdl:element</code>, <code>dfdl:sequence</code>, <code>dfdl:choice</code>, <code>dfdl:group</code></p>
-----------------	--

Deleted: insensitive

**Table 51 Properties for Choice Groups**

A choice can have a `dfdl: initiator` and/or a `dfdl:terminator` as described earlier.

We will use this terminology:

<i>Branch</i>	A <i>branch</i> is one of the available alternatives within a choice. A branch can be an element of simple type or complex type, or it can be an embedded sequence, choice or group reference.
<i>Root of the Branch</i>	Each <i>branch</i> conceptually has a single schema component at its root which is an element, sequence, choice or group reference. This component is known as the <i>Root of the Branch</i> .

**Table 52 Choice group terminology**

The Root of the Branch MUST NOT be optional. That is XSDL `minOccurs` MUST BE greater than 0.

A choice that declares no branches in the DFDL schema is a Schema Definition Error.

When processing a choice group the parser validates any contained path expressions. If a path expression contained inside a choice branch refers to any other branch of the choice, then it is a Schema Definition Error. Note that this rule handles nested choices also. A path that navigates outward from an inner choice to another alternative of an outer choice is violating this rule with respect to the outer choice.

### 15.1 Resolving Choices

A choice corresponds to concepts called variant records, multi-format records, discriminated unions, or tagged unions in various programming languages. In some contexts choices are referred to generally as 'unions'. However, this should not be confused with XSDL unions.

When processing a choice, there are two ways to resolve the intended branch. In one, speculative parsing is used. In the other, a constant-time direct dispatch to a branch is performed.

#### 15.1.1 Resolving Choices via Speculation

Speculative resolution works as follows:

- 1) Attempt to parse the first branch of the choice.

- 2) If this fails with a processing error
  - a) If a dfdl:discriminator evaluated to true earlier on this branch then the parser is 'bound' to this branch and parsing of the entire choice construct fails with a processing error.
  - b) If the branch has a dfdl:initiator and the choice has dfdl:initiatedContent 'yes' then the parser is 'bound' to this branch and parsing of the entire choice construct fails with a processing error.
  - c) Otherwise we repeat from step 1 for the next branch of the choice.
- 3) It is a processing error if the branches of the choice are exhausted.
- 4) If a branch is successfully parsed without error, then that branch's infoset becomes the infoset for the parse of the choice construct.
- 5) If the branch is an element declaration having dfdl:occursCountKind='expression' or dfdl:occursCountKind='parsed', then zero instances are possible. If the branch parses successfully without a discriminator but produces no element occurrences, then the branch is considered missing, and the parser looks for the next branch

It is not possible for variable settings to be communicated from the speculative attempt to parse a branch to any other parsing situation. The speculative effort is completely isolated. Whether it succeeds or fails, neither the parse position in the source data, nor anything in the variable memory, nor the Infoset is affected.

Nested choices can require unbounded look ahead into the data.

#### 15.1.2 Resolving Choices via Direct Dispatch

Direct dispatch provides a constant-time dispatch to a choice branch independent of how many choice branches there are.

Direct dispatch is indicated by the dfdl:choiceDispatchKey property. This expression is evaluated to compute the string matching (case sensitive) one of the dfdl:choiceBranchKey property values of one of the choice branches.

When a match is found, it is as if a dfdl:discriminator had evaluated to true on that branch. It is selected as resolution of the choice, and there is no backtracking to try other alternative selections if a processing error occurs.

The dfdl:choiceBranchKey property can be placed on element references, local element declarations, local sequences, local choices, or group references. All values of dfdl:choiceBranchKey properties must be unique across all branches of a choice that carries a dfdl:choiceDispatchKey property and it is a Schema Definition Error otherwise.

Note that it is a Schema Definition Error if both dfdl:initiatedContent and dfdl:choiceDispatchKey are provided on the same choice. However, it is not an error if a discriminator exists on a choice branch along with a dfdl:choiceBranchKey.

#### 15.1.3 Unparsing Choices

On unparsing there is the question of how one identifies the appropriate schema choice branch corresponding to the data in the infoset. This is complicated by the fact that the children may not be elements. They may themselves be sequences or choices. The selection of the choice branch is as follows: The element in the infoset is used to search the choice branches in the schema, in schema definition order, but without looking inside any complex elements. If the element occurs in a branch, then that branch is selected and if subsequently a processing error occurs, this selection is not revisited (that is, there is no backtracking).

To avoid any unintended behavior, all the children of a choice can be modeled as elements.

Deleted:

Deleted: It



## 16. Properties for Array Elements and Optional Elements

These properties are for array elements (XSDL maxOccurs >1 or unbounded) or optional elements (XSDL minOccurs 0 and maxOccurs 1). The properties handle a logical one-dimensional array of any simple or complex type.

Property Name	Description
occursCountKind	<p>Enum</p> <p>Specifies how the actual number of occurrences is to be established.</p> <p>Valid values 'fixed', 'expression', 'parsed', 'implicit', 'stopValue'.            'fixed' means use the XSDL maxOccurs property.            'expression' means use the dfdl:occursCount property.            'parsed' means that the number of occurrences is determined solely by speculative parsing.            'implicit' means that the number of occurrences is determined by speculative parsing in conjunction with the XSDL minOccurs and maxOccurs properties.            'stopValue' means look for a mandatory logical stop value which signifies the end of the occurrences.            These values are described in detail in section 16.1.            Annotation: dfdl:element</p>
occursCount	<p>DFDL Expression</p> <p>Specifies the number of occurrences of the element.</p> <p>Required only when dfdl:occursCountKind is 'expression'.</p> <p>This property is computed by way of an expression which returns a non-negative integer. The expression must not contain forward references to elements which have not yet been processed.            Annotation: dfdl:element,</p>
occursStopValue	<p>List of DFDL Logical Values</p> <p>A whitespace separated list of logical values that specify the alternative logical stop values for the element.</p> <p>Required only when dfdl:occursCountKind is 'stopValue'.</p> <p>When parsing then if an occurrence of the element has a logical value that matches one of the values in this list then the parser must not expect any more occurrences of the element.</p> <p>On unparsing the first value will be inserted as an additional final occurrence in the array after all of the occurrences in the infoset have been output.</p> <p><u>The dfdl:occursStopValue property must not be empty string.</u></p> <p>Annotation: dfdl:element</p>

**Table 53 Properties for Array Elements and Optional Elements**

When XSDL minOccurs 1 and maxOccurs 1, the above properties are not used, and the behavior is as if dfdl:occursCountKind was 'fixed' as described in section 16.1.1.

## 16.1 The `dfdl:occursCountKind` property

### 16.1.1 `dfdl:occursCountKind` 'fixed'

The enum 'fixed' should be used when the number of occurrences is always the same. The number is provided by the XSDL `maxOccurs` property.

When parsing, `maxOccurs` occurrences are expected in the data. It is a processing error if less than `maxOccurs` occurrences are found or defaulted. The parser stops looking for occurrences when `maxOccurs` have been found or defaulted. When `maxOccurs` is 0, no occurrences are looked for in the data.

When unparsing, `maxOccurs` occurrences are expected in the infoset. It is a processing error if less than `maxOccurs` occurrences are found or defaulted. The processor stops looking for more occurrences in the infoset after `maxOccurs` occurrences have been found. When `maxOccurs` is 0, no occurrences are looked for in the infoset or written.

Deleted: , or if more than `maxOccurs` occurrences are found.

It is a Schema Definition Error if `minOccurs` is not equal to `maxOccurs`.

### 16.1.2 `dfdl:occursCountKind` 'implicit'

The enum 'implicit' should be used when the number of occurrences is to be established using speculative parsing, and there are lower and upper bounds to control the speculation. The bounds are provided by the XSDL `minOccurs` and XSDL `maxOccurs` properties.

When parsing, up to `maxOccurs` occurrences are expected in the data. It is a processing error if less than `minOccurs` occurrences are found or defaulted. The parser stops looking for occurrences when either `minOccurs` have been found or defaulted and speculative parsing does not find another occurrence, or `maxOccurs` have been found or defaulted. When `maxOccurs` is 0, no occurrences are looked for in the data.

When unparsing, up to `maxOccurs` occurrences are expected in the infoset. It is a processing error if less than `minOccurs` occurrences are found or defaulted. The processor stops looking for more occurrences in the infoset after `maxOccurs` occurrences have been found. When `maxOccurs` is 0, no occurrences are looked for in the infoset or written.

Deleted: , or if more than `maxOccurs` occurrences are found.

### 16.1.3 `dfdl:occursCountKind` 'parsed'

The enum 'parsed' should be used when the number of occurrences is to be established solely using speculative parsing.

When parsing, any number of occurrences is expected in the data. The parser stops looking for occurrences when speculative parsing does not find another occurrence. If validation is enabled, it is a validation error if less than XSDL `minOccurs` occurrences are found or defaulted, or greater than XSDL `maxOccurs` occurrences are found.

When unparsing, any number of occurrences is expected in the infoset. If validation is enabled, it is a validation error if less than `minOccurs` occurrences are found or defaulted, or if more than `maxOccurs` occurrences are found.

### 16.1.4 `dfdl:occursCountKind` 'expression'

The enum 'expression' should be used when the number of occurrences is calculated by evaluating a DFDL expression.

When parsing, the `dfdl:occursCount` expression is evaluated and provides the number of occurrences expected in the data. It is a processing error if less than `dfdl:occursCount` occurrences are found or defaulted. The parser stops looking for occurrences when `dfdl:occursCount` occurrences have been found. If validation is enabled, it is a validation error if less than XSDL `minOccurs` occurrences are found or defaulted, or more than XSDL `maxOccurs` occurrences are found. When `dfdl:occursCount` is 0, no occurrences are looked for in the data.

When unparsing, any number of occurrences are expected in the infoset. If validation is enabled, it is a validation error if less than `minOccurs` occurrences are found or defaulted, or if more than `maxOccurs` occurrences are found. The `dfdl:occurs` expression is not evaluated, The 'count' is the number of occurrences in the augmented infoset.

It is a Schema Definition Error if `dfdl:occursCount` is not provided or in scope.

#### 16.1.5 **dfdl:occursCountKind 'stopValue'**

The enum 'stopValue' should be used when the the end of the array is signaled by an occurrence having a logical value that is equal to one of the specified 'stop values'.

When parsing, any number of occurrences is expected in the data, followed by an occurrence which is a stop value as specified by `dfdl:occursStopValue`. It is a processing error if a stop value occurrence is not found in the data (including when there are zero other occurrences). The parser stops looking for occurrences once a stop value has been found. If validation is enabled, it is a validation error if less than XSDL `minOccurs` occurrences are found or defaulted, or more than XSDL `maxOccurs` occurrences are found, not including the stop value.

When unparsing, the behavior is the same as for 'parsed', with the addition that a stop value occurrence is output after the last infoSet occurrence. If `dfdl:occursStopValue` provides multiple stop values then the first is used.

The stop value itself is **not** added to the infoSet when parsing. It is a processing error if a stop value is found in the infoSet when unparsing. (This insures that the array can be reparsed, as the stop value will be placed automatically and only at the end.)

It is a Schema Definition Error if `dfdl:occursStopValue` is not provided or in scope.

It is a Schema Definition Error if the type of the element is complex.

It is a Schema Definition Error if any of the stop values provided by `dfdl:occursStopValue` do not conform to the simple type of the element.

#### 16.2 **Default Values for Arrays**

When parsing, required occurrences that have empty representation may trigger the application of a default value, as described in Section 9.4.2 Element Defaults When Parsing.

When unparsing, required occurrences that are missing from the infoSet may trigger the application of a default value, as described in Section 9.4.3 Element Defaults When Unparsing.

The application of default values is **not** dependent on `dfdl:occursCountKind`, only on whether the occurrence is required or optional, whether there is a default value specified, and whether the data contains the empty representation (parsing) or is missing (unparsing). Section 9.4 Element Defaults contains the details.

#### 16.3 **Arrays with DFDL Expressions**

If the value of a DFDL property of an array element (other than `dfdl:occursCount`) is given by a DFDL Expression, then the expression must be re-evaluated for each occurrence of the element in case the value changes.

#### 16.4 **Points of Uncertainty**

Arrays can have points of uncertainty depending on the value of `dfdl:occursCountKind`. See Section 9.3.3 Points of Uncertainty for details.

#### 16.5 **Arrays and Sequences**

In some situations arrays of elements and sequence groups of elements seem to be similar; however, there is no notion of the array itself independent of its contained elements. Arrays are distinctly different from sequence groups in this way.

A sequence can have its own initiator, and a complex element having that sequence as its content can also have its own initiator, so you could express two different initiators.

Unlike a sequence group, an array does not have its own initiator, terminator, or alignment. Those properties apply to each element occurrence of the array. To give an alignment, initiator, separator or terminator to an entire array you must enclose the element declaration for the array in a sequence group and specify the alignment, separator, initiator and terminator on the sequence group.

#### 16.6 **Forward Progress Requirement**

To prevent an infinite loop, the parsing of an array that is potentially unbounded must terminate when the following are true:

- The occurrence is a point of uncertainty;
- The position in the data does not move during the parsing of the occurrence (including any associated Separator, PrefixSeparator or PostfixSeparator region);
- The occurrence is known-to-exist with empty representation or nil representation.

An array is potentially unbounded if any of the following are true:

- dfdl:occursCountKind is 'stopValue'
- dfdl:occursCountKind is 'parsed'
- dfdl:occursCountKind is 'implicit' and XSDL maxOccurs is unbounded

When dfdl:occursCountKind is 'stopValue' this results in a processing error because the stop value will never be encountered.

Further, to prevent unnecessary consumption of resources for large bounded values of XSDL maxOccurs, the parsing of an array must terminate when the following are true:

- dfdl:occursCountKind is 'implicit';
- The occurrence is a point of uncertainty;
- The position in the data does not move during the parsing of the occurrence (including any associated Separator, PrefixSeparator or PostfixSeparator region);
- The occurrence is known-to-exist with empty representation.

### 16.7 Parsing Occurrences with Non-Normal Representation

When parsing a single array, it is possible to extract occurrences that have different representations (nil, empty, normal, absent) although with some values of dfdl:lengthKind certain combinations of representations are not possible.

Occurrences with nil representation are added to the infoset with **[nilled]** member true.

Occurrences with empty representation may or may not be added to the infoset, as described in Section 9.4. If a required occurrence is not added to the infoset, it may be a processing error, dependent on dfdl:occursCountKind as described in section 16.1.

Occurrences with absent representation are not added to the infoset. For a required occurrence it may be a processing error, dependent on dfdl:occursCountKind as described in section 16.1.

### 16.8 Sparse Arrays

Consider parsing an array where optional occurrences with empty representation are present in the data, but there are also later optional occurrences present with normal representation. Such an array is called a 'sparse array'.

If the indices of the occurrences are significant and need to be preserved, then the array may be modelled using an element with XSDL nillable 'true', dfdl:nilKind 'literalValue', dfdl:nilValue '%ES;' and dfdl:nilValueDelimiterPolicy the same as dfdl:emptyValueDelimiterPolicy. The occurrences with empty representation now become occurrences with nil representation, and will produce nil values in the infoset, so the absolute positions of all occurrences are preserved.

If the indices of the occurrences are not significant, then the array should be modelled using an element with XSDL nillable 'false'. Optional occurrences with empty representation will not create items in the infoset, so the absolute position of any optional occurrences with normal representation is not preserved. Optional occurrences with empty representation are therefore skipped.

## 17. Calculated Value Properties

This section describes properties which allow the creation of calculated elements. When parsing, the value of a calculated element is derived using a DFDL Expression, and not by processing bytes from the data stream. When unparsing, the value of a calculated element is derived using a DFDL Expression, and is not obtained from the infoset in the usual way.

Calculated elements allow a technique that is commonly called layering. In this technique, some elements are said to be in the physical layer, and some in the logical layer. When parsing, the logical layer values are computed from physical layer values. When unparsing the opposite occurs, that is the physical layer values are computed from the logical layer values.

Calculated elements are commonly used with hidden elements so as to hide the physical layer elements so that they do not become part of the infoset.

When a DFDL Schema is used to both parse and unparse data, then a calculated element on parsing will normally have one or more calculated elements on unparsing.

These properties apply to elements of simple type.

Property Name	Description
inputValueCalc	<p>DFDL Expression</p> <p>An expression that calculates the value of the element when parsing.</p> <p>It is a Schema Definition Error if the result type of the expression does not conform to the base type of the element.</p> <p>The element value created using <code>dfdl:inputValueCalc</code> is validated like any other element value (when validation is enabled).</p> <p>An element that specifies a <code>dfdl:inputValueCalc</code> expression has no representation of its own in the data stream. All other DFDL representation properties are ignored.</p> <p>When an element which carries this property appears in a sequence that has a separator, no separator is associated with the element. When parsing, no separator is expected in the input data. When unparsing, no separator is written to the output data.</p> <p>The element must not be optional nor an array nor be global.</p> <p>The DFDL Expression must not refer to this element nor cause a circular reference to this element. The expression must not contain forward references to elements which have not yet been processed.</p> <p>It is a Schema Definition Error if this property is specified on an element which has an XSDL fixed or default property.</p> <p>It is a Schema Definition Error if <code>dfdl:inputValueCalc</code> and <code>dfdl:outputValueCalc</code> are specified on the same element.</p> <p>It is not possible to place this property in scope on a <code>dfdl:format</code> annotation.</p> <p>If this property appears on an element declaration or element reference schema component, the appearance of any other DFDL properties on that component is a Schema Definition Error.</p> <p>If this property appears on an element reference, then DFDL properties expressed on the referenced global element declaration or its type are ignored.</p> <p>If this property appears on an element declaration, then DFDL properties expressed on its type are ignored.</p> <p>Annotation: <code>dfdl:element</code></p>
outputValueCalc	<p>DFDL Expression</p> <p>An expression that calculates the value of the current element when unparsing.</p>

**Deleted:** This property is not allowed to appear on a local element or element reference that is the root of a choice branch.¶

	<p>The element must not be optional nor an array nor be global.</p> <p>It is a Schema Definition Error if the result type of the expression does not conform to the base type of the element.</p> <p>The value created using dfdl:outputValueCalc is validated like any other element value (when validation is enabled).</p> <p>The value for the element, if any, in the infoset is ignored.</p> <p>The DFDL expression must not refer to this element nor cause a circular reference to this element. The expression may contain forward references to elements which have not yet been processed.</p> <p>It is a Schema Definition Error if dfdl:outputValueCalc is specified on an element which has an XSDL fixed or default property.</p> <p>It is a Schema Definition Error if dfdl:inputValueCalc and dfdl:outputValueCalc are specified on the same element.</p> <p>It is not possible to place this property in scope on a dfdl:format annotation.</p> <p>Annotation: dfdl:element</p>
--	---

**Table 54 Calculated Value Properties****17.1 Example: 2d Nested Array**

Consider this simple example. The data stream contains two elements giving the number of rows and number of columns of an array of numbers. The representation of the array is stored after these two elements.

```

<xs:complexType name="array">
  <xs:sequence dfdl:initiator="" >

    <xs:sequence dfdl:hiddenGroupRef="tns:hiddenArrayCounts"/>

    <xs:element name="rows" maxOccurs="unbounded"
      dfdl:occursCountKind="expression"
      dfdl:occursCount="{ ../nrows }">
      <xs:complexType>
        <xs:sequence>
          <xs:element name="cols" type="xs:float" maxOccurs="unbounded"
            dfdl:occursCountKind="expression"
            dfdl:occursCount=" { ../nrows } " />
          </xs:sequence>
        </xs:complexType>
      </xs:element>
    </xs:sequence>
  </xs:complexType>

  <xs:group name="hiddenArrayCounts" >
    <xs:sequence>
      <xs:element name="nrows" type="xs:unsignedInt"
        dfdl:representation="binary"
        dfdl:lengthKind="implicit"
        dfdl:outputValueCalc="{ count(../rows) }"/>
      <xs:element name="ncols" type="xs:unsignedInt"
        dfdl:representation="binary"
        dfdl:lengthKind="implicit"
        dfdl:outputValueCalc=
          "{ if ( count(../rows) ge 1 )
            then
              count(../rows[1]/cols)
            else
              0
            }"/>
    </xs:sequence>
  </xs:group>

```

```
</xs:sequence>
</xs:group>
```

In the example above we see that there are two hidden elements named 'nrows' and 'ncols'. These hidden elements' values are computed when unparsing from the number of occurrences in the 'rows' and 'cols' repeating elements. The 'rows' and 'cols' repeating elements number of occurrences are computed when parsing from the hidden elements 'nrows' and 'ncols'.

## 17.2 Example: Three-Byte Date

Logically, the data is a date.

```
<xs:element name="d" type="date"/>
```

Physically, it is stored as 3 single byte integers.

The format of this data is expressed as this schema:

```
<xs:sequence dfdl:representation="binary">
  <xs:element name="mm" type="byte" />
  <xs:element name="dd" type="byte" />
  <xs:element name="yy" type="byte"/>
</xs:sequence>
```

This physical representation can be hidden so that it does not become part of the infoset:

```
<xs:sequence>
  <xs:sequence dfdl:hiddenGroupRef="tns:hiddenpDate"/>

  <xs:element name="d" type="date">
    ...
  </xs:element>
</xs:sequence>

<xs:group name="hiddenpDate" >
  <xs:sequence>
    <xs:element name="pdate">
      <xs:complexType>
        <xs:sequence dfdl:representation="binary">
          <xs:element name="mm" type="byte" />
          <xs:element name="dd" type="byte" />
          <xs:element name="yy" type="byte"/>
        </xs:sequence>
      </xs:complexType>
    </xs:element>
  </xs:sequence>
</xs:group>
```

A calculation can be used to compute the logical date element 'd' from the physical 'pdate' when parsing:

```
<xs:sequence>
  ... hidden pdate here ...

  <xs:element name="d" type="date">
    <xs:annotation><xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:element>
        <dfdl:property name="inputValueCalc">
          {
            fn:date(fn:concat(if(..pdate/yy gt 50 )then "19" else "20",
              if ( ../pdate/yy gt 9 )
                then xs:string(..pdate/yy)
                else fn:concat("0",
                  xs:string(..pdate/yy)),
              "-",
                  xs:string(..pdate/mm),
              "-",
                  xs:string(..pdate/dd)))
```

```

    }
  </dfdl:property>
</dfdl:element>
</xs:appinfo></xs:annotation>
</xs:element>
...
</xs:sequence>

```

The expression above assembles a string resembling, for example, "2005-12-17" or "1957-3-9" which is the string representation of a date that is acceptable to the `fn:date` constructor function. The hidden element `'pdate'` is referenced by relative paths. The expression `'../pdate/yy'` accesses an element of type `'int'`, and the `xs:string` constructor function turns it into an integer.

Finally, we must handle the unparse case where the physical layer is computed from the logical layer:

```

<xs:sequence dfdl:representation="binary">
  <xs:element name="mm" type="byte">
    dfdl:outputValueCalc="{ fn:month-from-date(..d) }" />
  <xs:element name="dd" type="byte">
    dfdl:outputValueCalc="{ fn:day-from-date(..d) }" />
  <xs:element name="yy" type="byte">
    dfdl:outputValueCalc="{ fn:year-from-date(..d) idivmod 100 }"/>
</xs:sequence>

```

The entire example in one place:

```

<xs:sequence>
  <xs:sequence dfdl:hiddenGroupRef="tns:hiddenpDate"/>

  <xs:element name="d" type="date">
    <xs:annotation><xs:appinfo source="http://www.ogf.org/dfdl/">
      <dfdl:element>
        <dfdl:property name="inputValueCalc">
          {
            fn:date(fn:concat(if(..pdate/yy gt 50) then "19" else "20",
              if ( ../pdate/yy gt 9 )
                then xs:string(..pdate/yy)
                else fn:concat("0",
                  xs:string(..pdate/yy)),
              "-",
                  xs:string(..pdate/mm),
              "-",
                  xs:string(..pdate/dd)))
          }
        </dfdl:property>
      </dfdl:element>
    </xs:appinfo></xs:annotation>
  </xs:element>
  ...
</xs:sequence>

<xs:group name="hiddenpDate" >
  <xs:sequence>
    <xs:element name="pdate">
      <xs:complexType>
        <xs:sequence dfdl:representation="binary">
          <xs:element name="mm" type="byte">
            dfdl:outputValueCalc="{ fn:month-from-date(..d) }" />
          <xs:element name="dd" type="byte">
            dfdl:outputValueCalc="{ fn:day-from-date(..d) }" />
          <xs:element name="yy" type="byte">
            dfdl:outputValueCalc="{ fn:year-from-date(..d) idivmod 100 }"
          />

```



```
</xs:sequence>  
  </xs:complexType>  
  </xs:element>  
</xs:sequence>  
</xs:group>
```

The above sequence contains logically only a single date element.

### 18. External Control of the DFDL Processor

In addition to providing the DFDL schema and data to be parsed or serialized, DFDL Schemas can also be parameterized by external variables.

DFDL processors can provide implementation-defined means to specify:

1. The data to be processed: a data stream when parsing or an infoset when unparsing.
2. The DFDL schema to be used
3. The *distinguished global element declaration* to be used (specifying both name of element and namespace of that name)
4. Values for external variables

Notice also that a DFDL Schema, like any XML schema, can have multiple top-level element declarations; hence, the distinguished global element declaration is necessary to indicate which of these top-level element declarations is to be the starting point for processing data. The distinguished global element declaration may be omitted if the DFDL schema contains only one top-level element declaration.

The mechanism by which a DFDL processor is controlled is not specified by this standard. For example, command line DFDL processors may use command line options, but DFDL processors embedded in other kinds of software systems may need other mechanisms.

**19. Built-in Specifications**

For convenience, a standard set of named DFDL format definitions may be provided with DFDL processors. These built-in format definitions may be imported by DFDL schema authors.

**20. Conformance**

DFDL conformance can be claimed for schema documents and for processors

A schema document conforms to this specification if it conforms to the subset of XML Schema 1.0 defined in section 5.1 DFDL Subset of XML Schema and consists of components which individually and collectively satisfy all the relevant constraints specified in this document.

Conformance may be claimed separately for a DFDL parser, a DFDL unparser or a DFDL processor that parses and unparses.

1. A DFDL processor claiming conformance **MUST** identify the level of conformance and version specification claimed.
2. A minimal conforming DFDL processor conforms to this specification when it implements all the non-optional features defined in this document.
3. An extended conforming DFDL processor conforms to the specification when it implements all the non-optional features and some of the optional features defined in this document.
4. A fully conforming DFDL processor conforms to the specification when it implements all the features defined in this document.

See Section 21 Optional DFDL Features for the list of optional feature

It is the intention of the DFDL Work Group to provide a conformance test suit to help verify conformance with this specification.

## 21. Optional DFDL Features

The following table lists the features of the DFDL language that are considered optional for DFDL processor implementations. This list admits very small subsets of the full DFDL specification. For example, a binary-only subset without any expressions or variables is specifically allowed.

Feature	Detection
Validation	External switch
Named Formats	dfdl:defineFormat or dfdl:ref
Choices	xs:choice in xsd
Arrays where size not known in advance	dfdl:occursCountKind 'implicit', 'parsed', 'stopValue'
Expressions	Use of a DFDL expression in any property value
End of parent	dfdl:lengthKind "endOfParent"
Simple type restrictions	xs:simpleType in xsd
Text representation for types other than String	dfdl:representation "text" for Number, Calendar or Boolean types
Delimiters	dfdl:separator <> "" or dfdl:initiator <> "" or dfdl:terminator <> "" or dfdl:lengthKind "delimited"
Nils	XSDL nillable 'true' in xsd
Defaults	XSDL default or fixed in xsd
Lengths in Bits	dfdl:alignmentUnits 'bits' or dfdl:lengthUnits 'bits'
Delimited lengths and representation binary element	dfdl:representation 'binary' (or implied binary) and dfdl:lengthKind 'delimited'
Regular expressions	dfdl:lengthKind 'pattern', dfdl:assert with dfdl:testkind 'pattern', dfdl:discriminator with dfdl:testkind 'pattern'
Zoned numbers	dfdl:textNumberRep 'zoned'
IBM 390 packed numbers	dfdl:binaryNumberRep 'packed'
IBM 390 packed calendars	dfdl:binaryCalendarRep 'packed'
IBM 390 floats	dfdl:binaryFloatRep 'ibm390Hex'
Unordered sequences	dfdl:sequenceKind 'unordered'
Floating elements	dfdl:floating 'yes'
dfdl functions in expression language	dfdl:functions in expression
Hidden groups	dfdl:hiddenGroupRef <> "
Calculated values	dfdl:inputValueCalc <> " or dfdl:outputValueCalc <> "
Escape schemes	dfdl:defineEscapeScheme in xsd
Extended encodings	Any dfdl:encoding value beyond the core list
<u>UTF-16 Variable Width Characters</u>	<u>dfdl:utf16Width="variable"</u>
Asserts	dfdl:assert in xsd
Discriminators	dfdl:discriminator in xsd
Prefixed lengths	dfdl:lengthKind 'prefixed'
Variables	dfdl:defineVariable,

Deleted: Bi-Directional text.

	<code>dfdl:newVariableInstances</code> , <code>dfdl:setVariable</code> Variables in DFDL expression language Note that variables as a feature is dependent on the Expressions feature.
BCD calendars	<code>dfdl:binaryCalendarRep "bcd"</code>
BCD numbers	<code>dfdl:binaryNumberRep "bcd"</code>
Multiple schemas	<code>xs:include</code> or <code>xs:import</code> in <code>xsd</code>
IBM 4690 packed numbers	<code>dfdl:binaryNumberRep "ibm4690Packed"</code>
IBM 4690 packed calendars	<code>dfdl:binaryCalendarRep "ibm4690Packed"</code>
DFDL Byte Value Entities	Use of <code>%#r</code> syntax in a DFDL String Literal other than the <code>dfdl:fillByte</code> property
DFDL Standard Character Set Encodings	<code>dfdl:encoding</code> name begins with "X-DFDL-".
Bit Order - Least Significant Bit First	<code>dfdl:bitOrder</code> with value 'leastSignificantBitFirst'

**Table 55 Optional DFDL features**

In order to provide portability of a DFDL schema, a minimal or extended conforming processor must issue warnings about any DFDL properties it does not implement. This warning can simply state that the property was not recognized.

(This allows the implementation to simply have no knowledge of properties it does not need for the subset of features it implements.)

For example if the hidden groups feature is not implemented, then the implementation will most likely not recognize the `dfdl:hiddenGroupRef` property at all. Such an implementation must issue a warning that the `dfdl:hiddenGroupRef` property was not recognized.

It is a Schema Definition Error if a DFDL schema uses an optional feature that is not supported by a minimal or extended conforming processor.

## 22. Property Precedence

### 22.1 Parsing

The following list gives the order in which DFDL properties are examined when the DFDL parser is positioned at a particular component in the DFDL schema, and about to parse the bitstream modeled by that component.

#### 22.1.1 **dfdl:element (simple)** and **dfdl:simpleType**

- *Parsing: calculated value (does not apply to dfdl:simpleType or to global elements)*
  - dfdl:inputValueCalc
- *Parsing: common*
  - dfdl:bitOrder
  - dfdl:encoding
    - 'UTF-16' 'UTF-16BE' 'UTF-16LE'
      - dfdl:utf16Width
  - dfdl:encodingErrorPolicy
  - dfdl:ignoreCase
- *Parsing: nillable*
  - XSDL nillable (does not apply to dfdl:simpleType)
    - dfdl:nilKind
      - "literalValue", "logicalValue", "literalCharacter"
        - dfdl:nilValue
- *Parsing: occurrences (does not apply to dfdl:simpleType)*
  - dfdl:floating
  - (maxOccurs > 1 or unbounded) or (minOccurs = 0 and maxOccurs = 1)
    - dfdl:occursCountKind
      - "expression"
        - dfdl:occursCount
      - "fixed", "implicit"
        - minOccurs
        - maxOccurs
      - "parsed"
      - "stopValue"
        - dfdl:occursStopValue
- *Parsing: identification, framing & extraction*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits
  - dfdl:alignment
    - dfdl:alignmentUnits
  - dfdl:initiator
    - dfdl:nilValueDelimiterPolicy (does not apply to dfdl:simpleType)
    - dfdl:emptyValueDelimiterPolicy
  - dfdl:representation "text" or xs:simpleType is 'string'
    - dfdl:lengthKind
      - "implicit"

- XSD maxLength or dfdl:textBooleanTrueRep/dfdl:textBooleanFalseRep
      - dfdl:lengthUnits
    - "explicit"
      - dfdl:length
      - dfdl:lengthUnits
    - "prefixed"
      - dfdl:prefixLengthType
      - dfdl:prefixIncludesPrefixLength
      - dfdl:lengthUnits
    - "pattern"
      - dfdl:lengthPattern
    - "delimited", "endOfParent"
      - None
  - dfdl:textTrimKind
    - dfdl:textStringPadCharacter, dfdl:textNumberPadCharacter, dfdl:textBooleanPadCharacter or dfdl:textCalendarPadCharacter
    - dfdl:textStringJustification, dfdl:textNumberJustification, dfdl:textBooleanJustification or dfdl:textCalendarJustification
  - dfdl:escapeSchemeRef
- dfdl:representation "binary" or xs:simpleType is 'hexBinary'
  - dfdl:lengthKind
    - "implicit"
      - XSD maxLength or xs:simpleType
      - dfdl:lengthUnits
    - "explicit"
      - dfdl:length
      - dfdl:lengthUnits
    - "prefixed"
      - dfdl:prefixLengthType
      - dfdl:prefixIncludesPrefixLength
      - dfdl:lengthUnits
    - "delimited", "endOfParent"
      - None
- dfdl:terminator
  - dfdl:nilValueDelimiterPolicy (does not apply to dfdl:simpleType)
  - dfdl:emptyValueDelimiterPolicy
  - dfdl:documentFinalTerminatorCanBeMissing
- dfdl:trailingSkip
  - dfdl:alignmentUnits
- Parsing: conversion
  - XSD type property
    - "Number"
      - dfdl:decimalSigned
      - dfdl:representation

**Deleted:** <#>dfdl:textBidi¶  
 <#>dfdl:textBidiOrdering¶  
 <#>dfdl:textBidiOrientation ¶



- *"text"*
  - dfdl:textNumberRep
    - *"standard"*
      - dfdl:textNumberPattern
      - dfdl:textStandardDecimalSeparator
      - dfdl:textStandardGroupingSeparator
      - dfdl:textStandardExponentRep
      - dfdl:textNumberCheckPolicy
      - dfdl:textStandardInfinityRep
      - dfdl:textStandardNaNRep
      - dfdl:textNumberRounding
        - *"explicit"*
        - dfdl:textNumberRoundingMode
        - dfdl:textNumberRoundingIncrement
      - dfdl:textStandardZeroRep
      - dfdl:textStandardBase
    - *"zoned"*
      - dfdl:textNumberPattern
      - dfdl:textNumberCheckPolicy
      - dfdl:textNumberRounding
        - *"explicit"*
        - dfdl:textNumberRoundingMode
        - dfdl:textNumberRoundingIncrement
      - dfdl:textZonedSignStyle
- *"binary"*
  - dfdl:byteOrder
  - *xs:decimal and restrictions*
    - dfdl:binaryNumberRep
      - *"packed"*
        - dfdl:binaryPackedSignCodes
        - dfdl:binaryDecimalVirtualPoint
        - dfdl:binaryNumberCheckPolicy
      - *"bcd", "ibm4690Packed"*
        - dfdl:binaryDecimalVirtualPoint
        - dfdl:binaryNumberCheckPolicy

**Commented [MJB3]:** These two bullets should be indented one more level, but MS-Word will not do so.

**Commented [MJB4]:** These two bullets should be indented one more level but MS-Word will not do so.

- *"binary"*
      - dfdl:binaryDecimalVirtualPoint
  - *xs:float, xs:double*
    - dfdl:binaryFloatRep
- *"String"*
- *"Calendar"*
  - dfdl:representation
    - *"text"*
      - dfdl:calendarPatternKind
        - *"explicit"*
          - dfdl:calendarPattern
      - dfdl:calendarCheckPolicy
      - dfdl:calendarTimeZone
      - dfdl:calendarObserveDST
      - dfdl:calendarFirstDayOfWeek
      - dfdl:calendarDaysInFirstWeek
      - dfdl:calendarCenturyStart
      - dfdl:calendarLanguage
    - *"binary"*
      - dfdl:byteOrder
      - dfdl:binaryCalendarRep
        - *"packed"*
          - dfdl:packedDecimalSignCodes
          - dfdl:binaryNumberCheckPolicy
          - dfdl:calendarPatternKind
            - *"explicit"*
            - dfdl:calendarPattern
          - dfdl:calendarCheckPolicy
          - dfdl:calendarTimeZone
          - dfdl:calendarObserveDST
          - dfdl:calendarFirstDayOfWeek
          - dfdl:calendarDaysInFirstWeek
          - dfdl:calendarCenturyStart
        - *"bcd", "ibm4690Packed"*
          - dfdl:binaryNumberCheckPolicy
          - dfdl:calendarPatternKind
            - *"explicit"*
            - dfdl:calendarPattern
          - dfdl:calendarCheckPolicy
          - dfdl:calendarTimeZone
          - dfdl:calendarObserveDST
          - dfdl:calendarFirstDayOfWeek
          - dfdl:calendarDaysInFirstWeek

- dfdl:calendarCenturyStart
    - "binarySeconds", "binaryMilliseconds"
      - dfdl:binaryCalendarEpoch
  - "Opaque"
  - "Boolean"
    - dfdl:representation
      - "text"
        - dfdl:textBooleanTrueRep
        - dfdl:textBooleanFalseRep
      - "binary"
        - dfdl:byteOrder
        - dfdl:binaryBooleanTrueRep
        - dfdl:binaryBooleanFalseRep
  - dfdl:useNilForDefault (does not apply to dfdl:simpleType)
    - "true"
      - None
    - "false"
      - XSDL default or:fixed

#### 22.1.2 dfdl:element (complex)

- *Parsing: common*
  - dfdl:bitOrder
  - dfdl:encoding
    - 'UTF-16' 'UTF-16BE' 'UTF-16LE'
      - dfdl:utf16Width
  - dfdl:encodingErrorPolicy
  - dfdl:ignoreCase
- *Parsing: nillable*
  - XSDL nillable
    - dfdl:nilKind
      - "literalValue"
        - dfdl:nilValue (must be "%ES;")
- *Parsing: occurrences*
  - dfdl:floating
  - (maxOccurs > 1 or unbounded) or (minOccurs = 0 and maxOccurs = 1)
    - dfdl:occursCountKind
      - "expression"
        - dfdl:occursCount
      - "fixed", "implicit"
        - minOccurs
        - maxOccurs
      - "parsed"
- *Parsing: identification, framing & extraction*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits

- dfdl:alignment
  - not *"implicit"*
    - dfdl:alignmentUnits
- dfdl:initiator
  - dfdl:nilValueDelimiterPolicy
  - dfdl:emptyValueDelimiterPolicy
- dfdl:lengthKind
  - *"explicit"*
    - dfdl:length
    - dfdl:lengthUnits
  - *"prefixed"*
    - dfdl:prefixLengthType
    - dfdl:prefixIncludesPrefixLength
    - dfdl:lengthUnits
  - *"pattern"*
    - dfdl:lengthPattern
  - *"implicit", "delimited", "endOfParent"*
    - None
- dfdl:terminator
  - dfdl:nilValueDelimiterPolicy
  - dfdl:emptyValueDelimiterPolicy
  - dfdl:documentFinalTerminatorCanBeMissing
- dfdl:trailingSkip
  - dfdl:alignmentUnits

### 22.1.3 dfdl:sequence and dfdl:group (when reference is to a sequence)

- *Parsing: hidden (xs:sequence only)*
  - dfdl:hiddenGroupRef
- *Parsing: common*
  - dfdl:bitOrder
  - dfdl:encoding
    - 'UTF-16' 'UTF-16BE' 'UTF-16LE'
      - dfdl:utf16Width
  - dfdl:encodingErrorPolicy
  - dfdl:ignoreCase
- *Parsing: identification, framing & extraction*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits
  - dfdl:alignment
    - not *"implicit"*
      - dfdl:alignmentUnits
  - dfdl:initiator
  - dfdl:sequenceKind
  - dfdl:initiatedContent
  - dfdl:separator

- dfdl:separatorPosition
- dfdl:separatorSuppressionPolicy
- dfdl:terminator
  - dfdl:documentFinalTerminatorCanBeMissing
- dfdl:trailingSkip
  - dfdl:alignmentUnits

#### 22.1.4 dfdl:choice and dfdl:group (when reference is to a choice)

- *Parsing: common*
  - dfdl:bitOrder
  - dfdl:encoding
    - 'UTF-16' 'UTF-16BE' 'UTF-16LE'
      - dfdl:utf16Width
  - dfdl:encodingErrorPolicy
  - dfdl:ignoreCase
- *Parsing: identification, framing & extraction*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits
  - dfdl:alignment
    - *not "implicit"*
      - dfdl:alignmentUnits
  - dfdl:initiator
  - dfdl:choiceLengthKind
    - *"explicit"*
      - dfdl:choiceLength
  - dfdl:initiatedContent
  - dfdl:choiceDispatchKey
  - dfdl:choiceBranchKey (on elements)
  - dfdl:terminator
    - dfdl:documentFinalTerminatorCanBeMissing
  - dfdl:trailingSkip
    - dfdl:alignmentUnits

## 22.2 Unparsing

The following list gives the order in which DFDL properties are examined when the DFDL unparsing is positioned at a particular component in the DFDL Infoset, and about to unparsing and thereby create the bitstream which is the representation of that component.

### 22.2.1 dfdl:element (simple) and dfdl:simpleType

- *Unparsing: calculated value (does not apply to dfdl:simpleType or to global elements)*
  - dfdl:inputValueCalc (if set then element is ignored)
  - dfdl:outputValueCalc
- *Unparsing: common*
  - dfdl:outputNewLine
  - dfdl:bitOrder
  - dfdl:encoding
    - 'UTF-16' 'UTF-16BE' 'UTF-16LE'

- dfdl:utf16Width
  - dfdl:encodingErrorPolicy
  - dfdl:fillByte
- *Unparsing: occurrences (does not apply to dfdl:simpleType)*
  - (maxOccurs > 1 or unbounded) or (minOccurs = 0 and maxOccurs = 1)
    - dfdl:occursCountKind
      - "expression"
        - dfdl:occursCount
      - "fixed", "implicit"
        - minOccurs
        - maxOccurs
      - "parsed"
      - "stopValue"
        - dfdl:occursStopValue
- *Unparsing: conversion*
  - dfdl:useNilForDefault (does not apply to dfdl:simpleType)
    - "true"
      - None
    - "false"
      - XSDL default or fixed
  - XSDL nillable (does not apply to dfdl:simpleType)
    - dfdl:nilKind
      - "literalValue", "logicalValue", "literalCharacter"
        - dfdl:nilValue
  - XSD type property
    - "Number"
      - dfdl:decimalSigned
      - dfdl:representation
        - "text"
          - dfdl:textNumberRep
            - "standard"
              - dfdl:textNumberPattern
              - dfdl:textStandardBase
              - dfdl:textStandardDecimalSeparator
              - dfdl:textStandardGroupingSeparator
              - dfdl:textStandardExponentRep
              - dfdl:textNumberCheckPolicy
              - dfdl:textStandardInfinityRep
              - dfdl:textStandardNaNRep
              - dfdl:textNumberRounding
              - "explicit"
                - dfdl:textNumberRoundingMode

- dfdl:textNumberRoundingIncrement
    - dfdl:textStandardZeroRep
  - "zoned"
    - dfdl:textNumberPattern
    - dfdl:textNumberCheckPolicy
    - dfdl:textNumberRounding
    - "explicit"
      - dfdl:textNumberRoundingMode
      - dfdl:textNumberRoundingIncrement
    - dfdl:textZonedSignStyle
- "binary"
  - dfdl:byteOrder
  - xs:decimal and restrictions
    - dfdl:binaryNumberRep
      - "packed"
        - dfdl:binaryPackedSignCodes
        - dfdl:binaryDecimalVirtualPoint
      - "bcd", "ibm4690Packed"
        - dfdl:binaryDecimalVirtualPoint
      - "binary"
        - dfdl:binaryDecimalVirtualPoint
    - xs:float, xs:double
      - dfdl:binaryFloatRep
  - "String"
  - "Calendar"
    - dfdl:representation
      - "text"
        - dfdl:calendarPatternKind
          - "explicit"
            - dfdl:calendarPattern
        - dfdl:calendarCheckPolicy
        - dfdl:calendarTimeZone
        - dfdl:calendarObserveDST
        - dfdl:calendarFirstDayOfWeek
        - dfdl:calendarDaysInFirstWeek
        - dfdl:calendarLanguage
      - "binary"
        - dfdl:byteOrder

**Deleted:** <#>dfdl:textBidi¶  
 <#>dfdl:textBidiOrdering¶  
 <#>dfdl:textBidiOrientation¶  
 <#>dfdl:textBidiNumeralShapes ¶

**Deleted:** <#>dfdl:textBidi¶  
 <#>dfdl:textBidiOrdering¶  
 <#>dfdl:textBidiOrientation¶  
 <#>dfdl:textBidiSymmetric¶  
 <#>dfdl:textBidiShaped¶

**Deleted:** <#>dfdl:textBidi¶  
 <#>dfdl:textBidiOrdering¶  
 <#>dfdl:textBidiOrientation¶  
 <#>dfdl:textBidiSymmetric¶  
 <#>dfdl:textBidiShaped¶

- dfdl:binaryCalendarRep
  - "packed"
    - dfdl:packedDecimalSignCodes
    - dfdl:decimalVirtualPoint
    - dfdl:calendarPatternKind
      - "explicit"
        - dfdl:calendarPattern
    - dfdl:calendarCheckPolicy
    - dfdl:calendarTimeZone
    - dfdl:calendarObserveDST
    - dfdl:calendarFirstDayOfWeek
    - dfdl:calendarDaysInFirstWeek
    - dfdl:calendarCenturyStart
  - "bcd", "ibm4690Packed"
    - dfdl:decimalVirtualPoint
    - dfdl:calendarPatternKind
      - "explicit"
        - dfdl:calendarPattern
    - dfdl:calendarCheckPolicy
    - dfdl:calendarTimeZone
    - dfdl:calendarObserveDST
    - dfdl:calendarFirstDayOfWeek
    - dfdl:calendarDaysInFirstWeek
    - dfdl:calendarCenturyStart
  - "binarySeconds", "binaryMilliseconds"
    - dfdl:binaryCalendarEpoch
- "Opaque"
- "Boolean"
  - dfdl:representation
    - "text"
      - dfdl:textBooleanTrueRep
      - dfdl:textBooleanFalseRep
    - "binary"
      - dfdl:byteOrder
      - dfdl:binaryBooleanTrueRep
      - dfdl:binaryBooleanFalseRep
- *Unparsing: insertion & framing*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits
  - dfdl:alignment
    - not "implicit"
      - dfdl:alignmentUnits
  - dfdl:representation "text" or xs:simpleType 'string'
    - dfdl:escapeSchemeRef

**Deleted:** <#>dfdl:textBidi¶  
 <#>dfdl:textBidiOrdering¶  
 <#>dfdl:textBidiOrientation¶  
 <#>dfdl:textBidiSymmetric¶  
 <#>dfdl:textBidiTextShaped ¶



- dfdl:lengthKind
  - *"implicit"*
    - XSD maxLength or dfdl:textBooleanTrueRep/dfdl:textBooleanFalseRep
    - dfdl:lengthUnits
    - dfdl:textPadKind
      - dfdl:textStringPadCharacter, dfdl:textNumberPadCharacter, dfdl:textBooleanPadCharacter or dfdl:textCalendarPadCharacter
      - dfdl:textStringJustification, dfdl:textNumberJustification, dfdl:textBooleanJustification or dfdl:textCalendarJustification
    - dfdl:truncateSpecifiedLengthString
  - *"explicit"*
    - *not expression*
      - dfdl:length
      - dfdl:truncateSpecifiedLengthString
    - *expression*
      - XSD minLength or dfdl:textOutputMinLength
    - dfdl:lengthUnits
    - dfdl:textPadKind
      - dfdl:textStringPadCharacter, dfdl:textNumberPadCharacter, dfdl:textBooleanPadCharacter or dfdl:textCalendarPadCharacter
      - dfdl:textStringJustification, dfdl:textNumberJustification, dfdl:textBooleanJustification or dfdl:textCalendarJustification
  - *"prefixed"*
    - dfdl:prefixLengthType
    - dfdl:prefixIncludesPrefixLength
    - dfdl:lengthUnits
    - dfdl:textPadKind
      - dfdl:textStringPadCharacter, dfdl:textNumberPadCharacter, dfdl:textBooleanPadCharacter or dfdl:textCalendarPadCharacter
      - dfdl:textStringJustification, dfdl:textNumberJustification, dfdl:textBooleanJustification or dfdl:textCalendarJustification
      - XSD minLength or dfdl:textOutputMinLength
  - *"pattern", "delimited", "endOfParent"*
    - dfdl:textPadKind

- dfdl:textStringPadCharacter,  
dfdl:textNumberPadCharacter,  
dfdl:textBooleanPadCharacter or  
dfdl:textCalendarPadCharacter
  - dfdl:textStringJustification,  
dfdl:textNumberJustification,  
dfdl:textBooleanJustification or  
dfdl:textCalendarJustification
  - XSD minLength or dfdl:textOutputMinLength
- dfdl:representation *"binary"* or *xs:simpleType 'hexBinary'*
  - dfdl:lengthKind
    - *"implicit"*
      - XSD maxLength or xs:simpleType
      - dfdl:lengthUnits
    - *"explicit"*
      - dfdl:length
      - dfdl:lengthUnits
    - *"prefixed"*
      - dfdl:prefixLengthType
      - dfdl:prefixIncludesPrefixLength
      - dfdl:lengthUnits
    - *"delimited", "endOfParent"*
      - *None*
- dfdl:initiator
  - dfdl:nilValueDelimiterPolicy (does not apply to dfdl:simpleType)
  - dfdl:emptyValueDelimiterPolicy
- dfdl:terminator
  - dfdl:nilValueDelimiterPolicy (does not apply to dfdl:simpleType)
  - dfdl:emptyValueDelimiterPolicy
- dfdl:trailingSkip
  - dfdl:alignmentUnits

#### 22.2.2 dfdl:element (complex)

- *Unparsing: common*
  - dfdl:outputNewLine
  - dfdl:bitOrder
  - dfdl:encoding
    - 'UTF-16' 'UTF-16BE' 'UTF-16LE'
      - dfdl:utf16Width
  - dfdl:encodingErrorPolicy
  - dfdl:fillByte

- *Unparsing: nillable*
  - XSDL nillable (*does not apply to dfdl:simpleType*)
    - dfdl:nilKind
      - "literalValue"
        - dfdl:nilValue (must be "%ES;")
- *Unparsing: occurrences*
  - (maxOccurs > 1 or unbounded) or (minOccurs = 0 and maxOccurs = 1)
    - dfdl:occursCountKind
      - "expression"
        - dfdl:occursCount
      - "fixed", "implicit"
        - minOccurs
        - maxOccurs
      - "parsed"
- *Unparsing: insertion & framing*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits
  - dfdl:alignment
    - *not* "implicit"
      - dfdl:alignmentUnits
  - dfdl:initiator
    - dfdl:nilValueDelimiterPolicy
    - dfdl:emptyValueDelimiterPolicy
  - dfdl:lengthKind
    - "explicit"
      - dfdl:length
      - dfdl:lengthUnits
    - "prefixed"
      - dfdl:prefixLengthType
      - dfdl:prefixIncludesPrefixLength
      - dfdl:lengthUnits
    - "implicit", "pattern", "delimited", "endOfParent"
      - None
  - dfdl:terminator
    - dfdl:nilValueDelimiterPolicy
    - dfdl:emptyValueDelimiterPolicy
  - dfdl:trailingSkip
    - dfdl:alignmentUnits

#### 22.2.3 dfdl:sequence and dfdl:group (when reference is a sequence)

- *Unparsing: hidden (xs:sequence only)*
  - dfdl:hiddenGroupRef
- *Unparsing: common*
  - dfdl:outputNewLine

- dfdl:bitOrder
- dfdl:encoding
  - 'UTF-16' 'UTF-16BE' 'UTF-16LE'
    - dfdl:utf16Width
- dfdl:encodingErrorPolicy
- dfdl:fillByte
- *Unparsing: insertion & framing*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits
  - dfdl:alignment
    - *not "implicit"*
      - dfdl:alignmentUnits
  - dfdl:initiator
  - dfdl:separator
    - dfdl:separatorPosition
    - dfdl:separatorSuppressionPolicy
  - dfdl:terminator
  - dfdl:trailingSkip
    - dfdl:alignmentUnits

#### 22.2.4 dfdl:choice and dfdl:group (when reference is a choice)

- *Unparsing: common*
  - dfdl:outputNewLine
  - dfdl:bitOrder
  - dfdl:encoding
    - 'UTF-16' 'UTF-16BE' 'UTF-16LE'
      - dfdl:utf16Width
  - dfdl:encodingErrorPolicy
  - dfdl:fillByte
- *Unparsing: insertion & framing*
  - dfdl:leadingSkip
    - dfdl:alignmentUnits
  - dfdl:alignment
    - *not "implicit"*
      - dfdl:alignmentUnits
  - dfdl:initiator
  - dfdl:choiceLengthKind
    - *"explicit"*
      - dfdl:choiceLength
  - dfdl:terminator
  - dfdl:trailingSkip
    - dfdl:alignmentUnits

### 23. Expression language

The DFDL expression language allows the processing of values conforming to the data model defined in the DFDL Infoset. It allows properties in the DFDL schema to be dependent on the value of an occurrence of an element or the value of a DFDL variable. For example the length of the content of an element can be made dependent on the value of another element in the document.

The main uses of the expression language are as follows:

1. When a DFDL property needs to be set dynamically at parse time from the value of one or more elements of the data. Properties such as initiator, terminator, length, occursCount and separator accept an expression.
2. In a dfdl:assert annotation
3. In a dfdl:discriminator annotation to resolve uncertainty when parsing
4. In a dfdl:inputValueCalc property to derive the value of an element in the logical model that doesn't exist in the physical data.
5. In a dfdl:outputValueCalc property to compute the value of an element on unparsing.
6. As the value in a dfdl:setVariable annotation or the dfdl:defaultValue in a dfdl:defineVariable or dfdl:newVariableInstance.

The DFDL expression language is a subset of XPath 2.0 [\[XPath\]](#). DFDL uses a subset of XML schema and has a simpler information model, so only a subset of XPath 2.0 expressions is meaningful in DFDL Schemas. For example there are no attributes in DFDL so the attribute axis is not needed.

[XPath 2.0 specification \[XPath2\] allows implementation-dependent evaluation of expressions thereby allowing either lazy \(sequential\) evaluation or full \(parallel\) evaluation of expressions with OR and AND clauses. This flexibility is not desirable in DFDL 1.0 implementations, so the specification is changed to prescribe lazy \(sequential\) evaluation left-to-right.](#)

In addition, DFDL expressions never return node-sequences having more than one node. DFDL expressions either return a simple value, a node sequence containing exactly one node/value, or an empty node sequence. Node sequences of length greater than one can be used within the expression, just not as the final result. Alternatively, one can state this as there are no constructs in DFDL which can accept a node sequence of more than one node; hence, DFDL expressions can never return a node sequence of more than one node as their final result.

For nilled elements, an attempt to get the value of a nilled element returns an empty node sequence.

DFDL implementations MUST comply with the error code behaviour in Appendix G of the XPath 2.0 spec and map these to the correct DFDL failure type. All but one of XPath's errors map to a Schema Definition Error. The exception is XPTY0004, which is used both for static and dynamic cases of type mismatch. A static type mismatch maps to a Schema Definition Error, whereas a dynamic type mismatch maps to a processing error. A DFDL implementation should distinguish the two kinds of XPTY0004 error if it is able to do so, but if unable it should map all XPTY0004 errors to a Schema Definition Error

Implementation Note: DFDL implementations may use off-the-shelf XPath 2.0 processors, but will need to pre-process DFDL expressions to ensure that the behaviour matches the DFDL specification:

- Ensure that what is returned as the result is not a sequence with length > 1 by appropriate use of fn:exactly-one().
- Check for the disallowed use of those XPath 2.0 functions that are not in the DFDL subset

[XPath 2.0 specification \[XPath2\] defines its functions to be in namespace http://www.w3.org/2005/xpath-functions. The DFDL specification assumes namespace prefix "fn:" is bound to this namespace.](#)

### 23.1 Expression Language Data Model

The DFDL expression language operates on the DFDL infoset with the addition of the hidden elements. That is, it operates on the *augmented* infoset.

In general, a DFDL expression can only reference an element that precedes the position in the schema where the expression is declared, and it is a schema definition otherwise, with the following exceptions:

- An assert or discriminator on a component may reference an element that is a descendent of the component.
- A `dfdl:outputValueCalc` property may reference an element that follows the position in the schema where the property is specified.

Implementations may have implementation-~~defined~~ limitations on the use of forward or backward reference, or may provide controls for bounding the reach of such references. These mechanisms are beyond the scope of this specification.

Deleted: dependent

### 23.2 Variables

A variable is a binding between a (qualified) name and a (typed) value. Variables are defined using the `dfdl:defineVariable` annotation (see 7.7); defining a variable causes an initial instance also to be created. Further instances of variables are created using the `dfdl:newVariableInstance` annotation. Instances of variables are assigned a value using the `dfdl:setVariable` annotation. Variables are referenced in expressions by preceding the QName with '\$'.

This section describes the semantics of variables. Any implementation consistent with the behavior described here is acceptable.

The memory where the information about a variable is stored during DFDL processing is called the *variable memory*. A variable is a name that is associated with a storage tuple in the variable memory.

Specifically, the variable memory contains:

- a counter used to generate locations for new tuples. Initial value is 1.
- an ordered list of locations. Each location contains a tuple of values:
  - has-been-set flag. This Boolean is originally false. `dfdl:setVariable` changes this flag to true.
  - has-been-referenced flag. This Boolean is originally false. Evaluation of an expression that uses the variable value changes the value to true.
  - has-value flag. This Boolean is originally true if the `dfdl:defineVariable` or `dfdl:newVariableInstance` annotation has a default value specified, or if a default value has been supplied externally. Otherwise it is false, but is set to true if a `dfdl:setVariable` annotation is processed.
  - typeId. This string is a type identifier taken from the type specified in the `dfdl:defineVariable` annotation.
  - value. This is a typed value, or the distinguished value "unknown". The type of the value must correspond to the typeId. The value is optionally specified in `dfdl:defineVariable` or `dfdl:newVariableInstance` annotations in which case we refer to it as the *default value* for the variable. A default value may also be provided by the DFDL processor when the variable is defined with external "true".

The variable memory is initialized when a `dfdl:defineVariable` annotation is encountered.

Each time a `dfdl:newVariableInstance` annotation is encountered, the parser captures the current value of the counter from the variable memory. It then creates a new variable memory where the location counter's value is one greater, and where the list of locations has been augmented with a new tuple at the location given by the prior value of the location counter. The tuple is initialized based on the specifics of the `dfdl:defineVariable` annotation.

### 23.2.1 Rewinding of Variable Memory State

Upon exit of the scope where the new variable instance was created, the newly created variable memory is discarded and the prior variable memory is restored.

Note that the above algorithm insures that each time a `dfdl:newVariableInstance` is encountered, a fresh location is initialized for it, and once the scope containing that variable goes out of scope, the instance tuple for the variable can no longer be reached. A different variable instance tuple may now be visible if there is one still in an enclosing scope.

### 23.2.2 Variable Memory State Transitions

The flags in the variable memory tuples are interpreted and modified as follows:

DFDL annotation	before annotation processed			after annotation processed		
	has-been-set	has-been-referenced	has-value	has-been-set	has-been-referenced	has-value
defineVariable (without default or external value)	tuple doesn't exist			false	false	false
defineVariable (with default value)	tuple doesn't exist			false	false	true
defineVariable (with external value)	tuple doesn't exist			false	false	true
newVariableInstance (without default value)	tuple doesn't exist			false	false	false
newVariableInstance (with default value)	tuple doesn't exist			false	false	true
setVariable	tuple doesn't exist			Schema Definition Error		
	false	false	false	true	false	true
	false	false	true	true	false	true (also value changed to new value)
	false	true	true	Schema Definition Error – set after reference not allowed.		
	true	any	true	Schema Definition Error - double set not allowed.		
reference variable (from DFDL expression)	tuple doesn't exist			Schema Definition Error		
	false	false	false	Schema Definition Error – undefined variable		
	any	any	true	false	true (value is returned)	true

**Table 56 Memory States for Expression Language Variables**

The above table describes a set of rules which might be abbreviated as:

- write once, read many
- no write after the value has been read

An exception to this behavior occurs whenever the DFDL processor backtracks because it is processing multiple arms of a choice or as a result of speculative parsing. In this case the variable state is also rewound.

It is a Schema Definition Error if a `dfdl:setVariable` or a variable reference occurs and there is no corresponding variable name defined by a `dfdl:defineVariable` annotation.

It is a Schema Definition Error if a `dfdl:setVariable` provides a value of incorrect type which does not correspond to the type specified by the `dfdl:defineVariable`.

It is a Schema Definition Error if a variable reference in an expression is able to return a value of incorrect type for the evaluation of that expression. That is, DFDL - including the expressions contained in it - is a statically type-checkable language. DFDL implementations may issue these Schema Definition Errors prior to processing time.

Even if the errors are detected at processing time, the errors associated with write-after-read, and double-write are Schema Definition Errors because they indicate the schema is not properly designed to use variables consistent with their single-assignment behavior.

### 23.3 General Syntax

DFDL expressions follow the XPath 2.0 syntax rules but are always enclosed in curly braces "{" and "}".

When a property accepts either a DFDL string literal or a DFDL expression, and the value is a string literal starting with a "{" character, then "{{" must be used to escape the "{" character.

The syntax "{}" is a Schema Definition Error as it results in an empty XPath 2.0 expression which is not legal. It is not the equivalent of setting the property to empty string.

Examples

```
{ /book/title }
{ $x+2 }
{ if (fn:exists(..field1)) then 1 else 0 }
```

The result of evaluating the expression must be a single atomic value of the type expected by the context, and it is a Schema Definition Error otherwise. Some XPath expressions naturally return a sequence of values, and in this case it is also Schema Definition Error if an expression returns a sequence containing more than one item.

Additionally:

- Every property that accepts an expression states exactly what the expression is expected to return. To ensure the returned value is of the correct type, an expression must use XPath constructors or the correct literal values.
- What appears lexically as the syntax of an expression follows XPath 2.0 rules. Note specifically that this is not the same as XSDL default and fixed property lexical syntax. Specifically, XSDL default and fixed properties do not accept expressions. They are always interpreted as XML Schema string literals. See [\[XSDL\]](#) for details.
- No extra auto-casting is performed over and above that provided by XPath 2.0. XPath 2.0 has rules for when it promotes types and when it allows types to be substituted. These are in Appendix B.1 of the XPath 2.0 spec.
- If the property is not expecting an expression to return a DFDL string literal, the returned value is never treated as a DFDL string literal.
- If expecting an expression to return a DFDL string literal, the returned value is always treated as a DFDL string literal.
- Within an expression, a string is never interpreted as a DFDL string literal.

### 23.4 DFDL Expression Syntax

Refer to XML Path Language (XPath) 2.0 [\[XPath\]](#) for a description of XPath expressions



DFDL Expression	::=	"{" Expr "}"
Expr	::=	ExprSingle
ExprSingle	::=	IfExpr   OrExpr
IfExpr	::=	"if" "(" Expr ")" "then" ExprSingle "else" ExprSingle
OrExpr	::=	AndExpr ( "or" AndExpr )*
AndExpr	::=	ComparisonExpr ( "and" ComparisonExpr )*
ComparisonExpr	::=	AdditiveExpr ( ValueComp ) AdditiveExpr?
AdditiveExpr	::=	MultiplicativeExpr ( "+"   "-" MultiplicativeExpr )*
MultiplicativeExpr	::=	<u>IntersectExceptExpr</u> ( "*"   "div"   "idiv"   "mod" ) <u>IntersectExceptExpr</u> *
<u>IntersectExceptExpr</u>	::=	<u>UnaryExpr ( ("intersect"   "except") UnaryExpr )*</u>
UnaryExpr	::=	("-"   "+")* ValueExpr
ValueExpr	::=	PathExpr
ValueComp	::=	"eq"   "ne"   "lt"   "le"   "gt"   "ge"
PathExpr	::=	("/" RelativePathExpr?)   RelativePathExpr   FilterExpr
RelativePathExpr	::=	StepExpr ( "/" StepExpr )*
StepExpr	::=	AxisStep
AxisStep	::=	(ReverseStep   ForwardStep) Predicate?
ForwardStep	::=	(ForwardAxis NodeTest)   AbbrevForwardStep
ForwardAxis	::=	("child" "::")   ("self" "::")
AbbrevForwardStep	::=	NodeTest   ContextItemExpr
ReverseStep	::=	(ReverseAxis NodeTest)   AbbrevReverseStep
ReverseAxis	::=	("parent" "::")
AbbrevReverseStep	::=	".."
NodeTest	::=	NameTest
NameTest	::=	QName
FilterExpr	::=	PrimaryExpr Predicate?
Predicate	::=	"[" Expr "]"
PrimaryExpr	::=	Literal   VarRef   ParenthesizedExpr   ContextItemExpr   FunctionCall
Literal	::=	NumericLiteral   StringLiteral
NumericLiteral	::=	IntegerLiteral   DecimalLiteral   DoubleLiteral
VarRef	::=	"\$" VarName
VarName	::=	QName
ParenthesizedExpr	::=	"(" Expr ")"
ContextItemExpr	::=	"."

Deleted: UnaryExpr

Deleted: UnaryExpr

FunctionCall	::=	QName "(" (ExprSingle ("," ExprSingle)*)? ")"
--------------	-----	---

**Table 57 DFDL Expression Language**

Notes:

1. Only *If* and *path* expression types are supported
2. Only the *child*, *parent*, and *self* axes are supported
3. Predicates are only used to index arrays and so must be integer expressions otherwise a Schema Definition Error occurs
4. A subset of the XPath 2.0 operators are supported
5. NameTest - These QNames are path steps that refer to elements in the DFDL infoset. If such an element is in a namespace, then the NameTest QName must have a prefix which is bound to the namespace. Specifically, any default namespace is not used to implicitly qualify these NameTest QNames. This behavior is consistent with XPath expression usage in XML Schema [W3C] such as in the path property of the xs:selector and xs:field elements within xs:key and xs:unique constraints, and in related XML standards such as XSLT. Note however, that this behavior is different from the way QNames are used in other places in XML and DFDL Schemas such as the ref property of an element reference, or the dfdl:ref property of a DFDL format annotation. There a QName with no prefix must always be referring to a global declaration or definition, and so is augmented with the default namespace when needed.

Commented [MJB5]: Test that hyperlink works to reference.

**23.5 Constructors, Functions and Operators**

In the function signatures below a '?' following an argument name, argument type or result type indicates that the argument/result can be a node or value of the expected type or it can have no value.

**23.5.1 Constructor Functions for XML Schema Built-in Types**

The arguments to the constructors are all of type xs:anyAtomicType. Since the expression language can be statically type checked, it is a Schema Definition Error if the type of the argument is not one of the DFDL-supported subtypes of xs:anyAtomicType,

However, many statically type-correct values will still not be convertible to the result type. It is a processing error if the supplied argument value is not convertible to the constructed type.

The following constructor functions for the built-in types are supported:

Function
xs:string(\$arg as xs:anyAtomicType) as xs:string
xs:boolean(\$arg as xs:anyAtomicType) as xs:boolean
xs:decimal(\$arg as xs:anyAtomicType) as xs:decimal
xs:float(\$arg as xs:anyAtomicType) as xs:float
xs:double(\$arg as xs:anyAtomicType) as xs:double
xs:dateTime(\$arg as xs:anyAtomicType) as xs:dateTime
xs:time(\$arg as xs:anyAtomicType) as xs:time
xs:date(\$arg as xs:anyAtomicType) as xs:date
xs:hexBinary(\$arg as xs:anyAtomicType) as xs:hexBinary
xs:integer(\$arg as xs:anyAtomicType) as xs:integer
xs:long(\$arg as xs:anyAtomicType) as xs:long

<code>xs:int(\$arg as xs:anyAtomicType) as xs:int</code>
<code>xs:short(\$arg as xs:anyAtomicType) as xs:short</code>
<code>xs:byte(\$arg as xs:anyAtomicType) as xs:byte</code>
<code>xs:nonNegativeInteger(\$arg as xs:anyAtomicType) as xs:nonNegativeInteger</code>
<code>xs:unsignedLong(\$arg as xs:anyAtomicType) as xs:unsignedLong</code>
<code>xs:unsignedInt(\$arg as xs:anyAtomicType) as xs:unsignedInt</code>
<code>xs:unsignedShort(\$arg as xs:anyAtomicType) as xs:unsignedShort</code>
<code>xs:unsignedByte(\$arg as xs:anyAtomicType) as xs:unsignedByte</code>

**Table 58 Basic Constructors**

A special constructor function is provided for constructing a `xs:dateTime` value from an `xs:date` value and an `xs:time` value.

Function
<code>fn:dateTime(\$arg1 as xs:date, \$arg2 as xs:time) as xs:dateTime</code>

**Table 59 Special Constructor for xs:dateTime**

## 23.5.2 Standard XPath Functions

### 23.5.2.1 Boolean functions

The following additional constructor functions are defined on the boolean type.

Function	Meaning
<code>fn:true()</code>	Constructs the <code>xs:boolean</code> value 'true'.
<code>fn:false()</code>	Constructs the <code>xs:boolean</code> value 'false'.

**Table 60 Boolean functions**

The following functions are defined on boolean values. The return type of these functions is `xs:boolean`:

Function	Meaning
fn:not(\$arg?)	<p>If \$arg is the empty sequence or a node with <b>[nilled]</b> true, fn:not returns true.</p> <p>If \$arg is a sequence containing a node with <b>[nilled]</b> false or <b>[nilled]</b> having no value (that is, a node corresponding to a non-nillable element), fn:not returns false.</p> <p>If \$arg is a value of type xs:boolean or a derived from xs:boolean, fn:not returns the boolean inverse of \$arg.</p> <p>If \$arg is a value of type xs:string or a type derived from xs:string, fn:not returns true if the operand value has zero length; otherwise it returns false.</p> <p>If \$arg is a value of any numeric type or a type derived from a numeric type, fn:not returns true if the operand value is NaN or is numerically equal to zero; otherwise it returns false.</p> <p>In all other cases, fn:not raises a processing error.</p> <p>Inverts the xs:boolean value of the argument.</p>

**Table 61 Boolean functions****23.5.2.2 Numeric Functions**

The following functions are defined on numeric types. Each function returns a value of the same type as the type of its argument. The argument must be convertible to a number type.

Function	Meaning
fn:abs(\$arg as numeric)	Returns the absolute value of the argument.
fn:ceiling(\$arg as numeric)	Returns the smallest number with no fractional part that is greater than or equal to the argument.
fn:floor(\$arg as numeric)	Returns the largest number with no fractional part that is less than or equal to the argument.
fn:round(\$arg as numeric)	Rounds to the nearest number with no fractional part. When the value is x.5, it rounds toward positive infinity.
fn:round-half-to-even(\$arg as numeric) fn:round-half-to-even(\$arg as numeric, \$precision as xs:integer)	Takes a number and a precision and returns a number rounded to the given precision. If the fractional part is exactly half, the result is the number whose least significant digit is even.

**Table 62 Numeric Functions****23.5.2.3 String Functions**

The following functions are defined on values of type xs:string and types derived from it. In the functions below which compare strings, DFDL always uses the default unicode collation algorithm (which is a comparison of codepoint values).

Function	Meaning
fn:concat( \$arg1 as xs:anyAtomicType, \$arg2 as xs:anyAtomicType, ... )	Concatenates two or more xs:anyAtomicType arguments cast to xs:string.

fn:substring(\$sourceString as xs:string, \$startingLoc as xs:double) fn:substring(\$sourceString as xs:string, \$startingLoc as xs:double, \$length as xs:double)	Returns the xs:string located at a specified place within an argument xs:string.
fn:string-length(\$arg as xs:string)	Returns the length of the argument as an xs:integer
fn:upper-case(\$arg as xs:string)	Returns the upper-cased value of the argument.
fn:lower-case(\$arg as xs:string)	Returns the lower-cased value of the argument.
fn:contains(\$arg1 as xs:string, \$arg2 as xs:string)	Returns xs:boolean indicating whether one xs:string contains another xs:string.
fn:starts-with(\$arg1 as xs:string, \$arg2 as xs:string)	Returns xs:boolean indicating whether the value of one xs:string begins with the characters of another xs:string.
fn:ends-with(\$arg1 as xs:string, \$arg2 as xs:string)	Returns xs:boolean indicating whether the value of one xs:string ends with the characters of another xs:string.
fn:substring-before(\$arg1 as xs:string, \$arg2 as xs:string)	Returns the characters of one xs:string that precede in that xs:string the characters of another xs:string.
fn:substring-after(\$arg1 as xs:string, \$arg2 as xs:string)	Returns the characters of xs:string that follow in that xs:string the characters of another xs:string.

**Table 63 String Functions****23.5.2.4 Date and Time Functions**

Function	Meaning
fn:year-from-dateTime(\$arg as xs:dateTime)	Returns the year from an xs:dateTime value as an xs:integer.
fn:month-from-dateTime(\$arg as xs:dateTime)	Returns the month from an xs:dateTime value as an xs:integer.
fn:day-from-dateTime(\$arg as xs:dateTime)	Returns the day from an xs:dateTime value as an xs:integer.
fn:hours-from-dateTime(\$arg as xs:dateTime)	Returns the hours from an xs:dateTime value as an xs:integer.
fn:minutes-from-dateTime(\$arg as xs:dateTime)	Returns the minutes from an xs:dateTime value as an xs:integer.
fn:seconds-from-dateTime(\$arg as xs:dateTime)	Returns the seconds from an xs:dateTime value as an xs:decimal.
fn:year-from-date(\$arg as xs:date)	Returns the year from an xs:date value as an xs:integer.
fn:month-from-date(\$arg as xs:date)	Returns the month from an xs:date value as an xs:integer.

fn:day-from-date(\$arg as xs:date)	Returns the day from an xs:date value as an xs:integer.
fn:hours-from-time(\$arg as xs:time)	Returns the hours from an xs:time value as an xs:integer.
fn:minutes-from-time(\$arg as xs:time)	Returns the minutes from an xs:time value as an xs:integer.
fn:seconds-from-time(\$arg as xs:time)	Returns the seconds from an xs:time value as an xs:decimal.

**Table 64 Date and Time Functions****23.5.2.5 Node Sequence Test Functions**

The following functions are defined on sequences. (Note that DFDL v1.0 does not support sequences of length > 1.)

Function	Meaning
fn:empty(\$arg?)	Indicates whether or not the provided sequence is empty.
fn:exists(\$arg?)	Indicates whether or not the provided sequence is not empty.
fn:exactly-one(\$arg?)	<u>Returns the input sequence if it contains exactly one item. Raises an error otherwise</u>
fn:count(\$arg)	Returns the number of items in the value of \$arg as an xs:integer. Returns 0 if \$arg is the empty sequence.

**Deleted:** True if the provided sequence contains exactly one node/value....

**Table 65 Node Sequence Test Functions****23.5.2.6 Node functions**

This section discusses functions and operators on nodes.

Function	Meaning
fn:local-name() fn:local-name(\$arg)	Returns the local name of the context node or the specified node as an xs:string.
fn:namespace-uri() fn:namespace-uri(\$arg)	Returns the namespace URI as an xs:string for the argument node or the context node if the argument is omitted. Returns empty string if the argument/context node is in no namespace.

**Table 66 Node functions****23.5.2.7 Nillable Element Functions**

This section discusses functions related to nillable elements.

Function	Meaning
fn:nilled(\$arg?)	Returns an xs:boolean true when the argument node Infoset member <b>[nilled]</b> is true and false when <b>[nilled]</b> is false. If the argument is not an element node, returns the empty sequence. If the argument is the empty sequence, returns the empty sequence. If the argument is an element node and <b>[nilled]</b> has no value returns the empty sequence.

**Table 67 Nillable Element Functions**

## 23.5.3 DFDL Functions

Function	Meaning
dfdl:contentLength(\$node, \$lengthUnits)	<p>Returns the length of the supplied node's SimpleContent region for elements of simple type, or ComplexContent region for elements of complex type. These regions are defined in Section 9.2 DFDL Data Syntax Grammar. The value is returned as an xs:unsignedLong.</p> <p>The second argument is of type xs:string and must be 'bytes', 'characters', or 'bits' (Schema Definition Error otherwise) and determines the units of length.</p>
dfdl:valueLength(\$node, \$lengthUnits)	<p>Returns the length of the supplied node's SimpleLogicalValue region for elements of simple type, or ComplexValue region for elements of complex type. These regions are defined in Section 9.2 DFDL Data Syntax Grammar. The value is returned as an xs:unsignedLong.</p> <p>For simple types, the valueLength() function returns a length which excludes any padding or filling.</p> <p>The second argument is of type xs:string and must be 'bytes', 'characters', or 'bits' (Schema Definition Error otherwise) and determines the units of length.</p>
dfdl:testBit(\$data, \$bitPos)	<p>Returns Boolean true if the bit number given by the xs:nonNegativeInteger \$bitPos is set on in the xs:unsignedByte given by \$data, otherwise returns Boolean false.</p>
dfdl:setBits(\$bit1, \$bit2, ... \$bit8)	<p>Returns an unsigned byte being the value of the bit positions provided by the Boolean arguments, where true is 1, false is 0. The number of arguments must be 8.</p>
dfdl:occursIndex()	<p>Returns the position of the current item of an array as an xs:nonNegativeInteger.</p> <p>The first element is at position 1.</p> <p>The function may be used on non-array elements so long as it appears within the dynamic scope of some array element.</p> <p>In this case it returns the index of the current item of the innermost enclosing array element.</p> <p>It is a Schema Definition Error if this function is called when there is no enclosing array element.</p>
dfdl:checkConstraints(\$node)	<p>Returns boolean true if the specified node value satisfies the XML schema facet constraints that are associated with it. Returns false if the specified node does not meet the constraints or does not exist.</p> <p>The facets that are checked are</p> <ul style="list-style-type: none"> <li>• minLength, maxLength</li> <li>• pattern</li> <li>• enumeration</li> <li>• maxInclusive, maxExclusive, minExclusive, minInclusive</li> </ul>

Deleted: SimpleValue or NilLogicalValue

Deleted: ComplexContent

	<ul style="list-style-type: none"><li>totalDigits</li><li>fractionDigits</li></ul> <p>See Section 5.2 for which facets are checked for each simple type.</p> <p>Additionally the XSD fixed property is checked.</p> <p>It is a Schema Definition Error if the argument is a complex element.</p>	
dfdl:encodeDFDLEntities(\$arg)	Returns a string containing a DFDL string literal constructed from the \$arg string argument. If \$arg contains any '%' and/or space characters, then the return value replaces each '%' with '%%' and each space with '%SP;', otherwise \$arg is returned unchanged.	
dfdl:decodeDFDLEntities (\$arg)	Returns a string constructed from the \$arg string argument. If \$arg contains syntax matching DFDL Character Entities syntax, then the corresponding characters are used in the result. Any characters in \$arg not matching the DFDL Character Entities syntax remain unchanged in the result.	It is a Schema Definition Error if \$arg contains syntax matching DFDL Byte Value Entities syntax.
dfdl:containsDFDLEntities(\$arg)	Returns a Boolean indicating whether the \$arg string argument contains one or more DFDL entities.	
dfdl:timeZoneFromDateTime(\$arg) dfdl:timeZoneFromDate(\$arg) dfdl:timeZoneFromTime (\$arg)	Returns the timezone component, if any, of \$arg as an xs:string. The \$arg is of type xs:dateTime, xs:date and xs:time respectively.	If \$arg has a timezone component, then the result is a string in the format of an ISO Time zone designator. Interpreted as an offset from UTC, its value may range from +14:00 to -14:00 hours, both inclusive. The UTC time zone is represented as "+00:00". If the \$arg has no timezone component, then "" (empty string) is returned.
<u>dfdl:checkRangeInclusive(\$node, \$val1, \$val2)</u> <u>dfdl:checkRangeExclusive(\$node, \$val1, \$val2)</u>	<u>Returns boolean true if the specified node value is in the range given by \$val1 and \$val2.</u> <u>The type of \$val1 and \$val2 must be compatible with the type of \$node, and must be a derivative of xs:decimal, xs:float or xs:double. It is a Schema Definition Error if the \$node argument is a complex element.</u>	

**Table 68 DFDL Functions**

Notes:

dfdl:valueLength(path, lengthUnits) - returns the value length which excludes any padding or filling which might be added for a *specified length*

If the element declaration in the DFDL schema corresponding to the infoSet item is not potentially represented, then the unpadded length is defined to be 0.

The value length includes the length contributions from introduced escape characters needed to escape contained delimiters (if such are defined, and will appear in the output representation).

The value length is also a function of the dfdl:encoding property. Multi-byte and variable-width character set encodings will commonly contribute more bytes to the value length than a single-byte character set would.



The value length is computed from the DFDL infoSet value, ignoring the `dfdl:length` or `dfdl:textOutputMinLength` property. Other DFDL properties which affect the length of a text or binary representation are respected, it is only an explicit length which is ignored.

For a complex type, this means a bottom up totaling of the `dfdl:contentLength()` of all the contents and framing of the complex type.

`dfdl:contentLength(path, lengthUnits)` – returns the length of the content of the infoSet data item as identified by the path argument. This includes padding or filling or truncation which might be carried out for a *specified length* item.

If the element declaration in the DFDL schema corresponding to the infoSet item is not potentially represented (e.g., has an `dfdl:inputValueCalc` property), then the length is defined to be 0.

When unparsing with `dfdl:lengthKind` "explicit", the calculation of `dfdl:contentLength()` returns the value of the `dfdl:length` property.

For both `dfdl:contentLength()` and `dfdl:valueLength()`, the content length excludes any alignment filling as well as excluding any leading or trailing skip bytes. That is, the returned length is about the length of the content, and not about the position of that content in the output data stream.

Use `dfdl:encodeDFDLEntities()` when the value of a DFDL property is obtained from the data stream using an expression, and the type of the property is DFDL String Literal or List of DFDL String Literals, and the values extracted from the data stream could contain '%' or space characters. If the data already contains DFDL entities, this function should not be used.

Use `dfdl:decodeDFDLEntities()` when you need to create a value which contains characters for which DFDL Character Entities are needed. An example is to create data containing the NUL (character code 0) codepoint. This character code is not allowed in XML documents, including DFDL Schemas; hence, it must be specified using a DFDL Character Entity. Within a DFDL Expression, use this function to obtain a string containing this character.

#### 23.5.4 DFDL Constructor Functions

There is sometimes a need to create a number type from hex binary, and a hex binary type from a number. Accordingly the following new DFDL specific functions are provided.

Function	Meaning
<code>dfdl:byte (\$arg)</code> <code>dfdl:unsignedByte (\$arg)</code> <code>dfdl:short (\$arg)</code> <code>dfdl:unsignedShort (\$arg)</code> <code>dfdl:int (\$arg)</code> <code>dfdl:unsignedInt (\$arg)</code> <code>dfdl:long (\$arg)</code> <code>dfdl:unsignedLong (\$arg)</code>	<p>These constructor functions behave identically to the XPath 2.0 constructor functions of the same names, with one exception. The argument can be a quoted string beginning with the letter 'x', in which case the remainder of the string is hexadecimal digits that represent a big-endian two's complement representation of a binary number.</p> <p>If the string begins with 'x', it is a Schema Definition Error if a character appears other 0-9, a-f, A-F.</p> <p>Each constructor function has a limit on the number of hex digits, with no more digits than 2, 4, 8, or 16 for the byte, short, int and long versions respectively. It is a Schema Definition Error if more digits are encountered than are suitable for the type being created</p>

dfdl:hexBinary (\$arg)	<p>This constructor function behaves identically to the XPath 2.0 constructor function of the same name, with one exception. The argument can also be a long, unsignedLong, or any subtype thereof, and in that case a xs:hexBinary value containing a number of hex digits is produced. The ordering and number of the digits correspond to a binary big-endian two's-complement implementation of the type of the argument. Digits 0-9, A-F are used.</p> <p>The number of digits produced depends on the type of \$arg, being 2, 4, 8 or 16. If \$arg is a literal number then the type is the smallest signed type (long, int, short, byte) that can contain the value.</p> <p>If a literal number is not able to be represented by a long, it is a Schema Definition Error.</p>
------------------------	--

**Table 69: DFDL Constructor Functions**

Examples:

- dfdl:unsignedInt("xa1b2c3d4") is the unsigned int value 2712847316.
- dfdl:int("FFFFFFFF") is the signed int value -1.
- dfdl:unsignedByte("xFF") is the unsigned byte value 255.
- dfdl:byte("xff") is the signed byte value -1.
- dfdl:byte("x7F") is the signed byte value 127.
- dfdl:byte("x80") is the signed byte value -128.
- dfdl:unsignedByte("x80") is the unsigned byte value 128.
- dfdl:byte("x0A3") is a Schema Definition Error (too many digits for type).
- dfdl:byte("xG3") is a Schema Definition Error (invalid digit).
- dfdl:hexBinary(xs:unsignedByte(208)) is the hexBinary value "D0".
- dfdl:hexBinary(208) is the hexBinary value "00D0".
- dfdl:hexBinary(-2084) is the hexBinary value "F7DC".

**23.5.5 Miscellaneous Functions**

Function	Meaning
fn:error()	Causes a processing error.
fn:error(\$id as xs:string)	This function does not return a value. A processing error ends the evaluation of the expression.
fn:error(\$id as xs:string, \$desc as xs:string, \$obj?)	The \$id argument is an error code identifier string that distinguishes this error from others. The string should have the structure of an XSD QName; the namespace URI conventionally identifies the component, subsystem, or authority responsible for defining the meaning of the error code, while the local part identifies the specific error condition. This information is incorporated into any diagnostic messages created by the DFDL implementation in response to the processing error in an implementation-dependent manner. If the \$id argument string does not have the form of an XSD QName, or the QName cannot be interpreted as a meaningful namespace prefix and

Commented [MJB6]: Verify formatting of table once change is accepted.

	<p><u>local identifier, then the processing error still occurs but the diagnostic message is created in an implementation-dependent manner.</u></p> <p><u>The \$desc is a natural-language description of the error condition. This string will appear in any diagnostic messages created by the DFDL implementation in response to the processing error.</u></p> <p><u>The \$obj? argument is an arbitrary value used to convey additional information about the error and it is used to construct the diagnostic message in an implementation-dependent manner.</u></p> <p><u>If any argument is not supplied the processing error occurs but the diagnostic message created is implementation-dependent.</u></p>
--	---

**23.6 Unparsing and Circular Expression Deadlock Errors**

It is possible for expressions and lengths of elements in a DFDL schema to interact badly, resulting in circular deadlocks where an expression is unable to evaluate because it depends in some way on the length of something that depends on the expression itself.

Expression deadlocks are always Schema Definition Errors.

One scenario where such a deadlock can arise is due to what is called the *interior-alignment problem*. If a dfdl:outputValueCalc expression depends on the valueLength of a following complex element which due to interior alignments, has a length that depends on its starting position, then we have a circular deadlock.

Commented [7]:

## 24. DFDL Regular Expressions

A DFDL regular expression may be specified for the `dfdl:lengthPattern` format property and the `dfdl:testPattern` property of the `dfdl:assert` and `dfdl:discriminator` annotations. DFDL regular expressions do not interpret DFDL entities.

A DFDL regular expression is defined by a set of valid pattern characters. For portability, a DFDL regular expression pattern is restricted to the inclusive subset of the ICU regular expression [\[ICURegex\]](#) and the Java(R) 7 regular expression [\[JavaRegex\]](#) with the Unicode flags `UNICODE_CASE` and `UNICODE_CHARACTER_CLASS` turned on. DFDL regular expressions thereby conform to Unicode Technical Standard #18, Unicode Regular Expressions, level 1 [\[UnicodeRegex\]](#).

The following regular expression constructs are not common to both ICU and Java(R) 7 and it is a Schema Definition Error if any are used in a DFDL regular expression:

Construct	Meaning	Notes
<code>\N{UNICODE CHARACTER NAME}</code>	Match the named character	ICU only
<code>\X</code>	Match a Grapheme Cluster	ICU only
<code>\Uhhhhhhh</code>	Match the character with the hex value hhhhhhhh	ICU only
<code>(?# ... )</code>	Free-format comment	ICU only
<code>(?w-w)</code>	<code>UREGEX_UWORD</code> - Controls the behaviour of <code>\b</code> in a pattern	ICU only
<code>(?d-d)</code>	<code>UNIX_LINES</code> - Enables Unix lines mode	Java 7 only
<code>(?u-u)</code>	<code>UNICODE_CASE</code> - Enables Unicode-aware case folding	Java 7 only (1)
<code>(?U-U)</code>	<code>UNICODE_CHARACTER_CLASS</code> - Enables the Unicode version of predefined character classes and POSIX character classes	Java 7 only (2)

**Table 70 Disallowed Regular Expression Constructs**

Notes:

(1) Implementations using Java 7 must set flag `UNICODE_CASE` by default to match ICU.

(2) Implementations using Java 7 must set flag `UNICODE_CHARACTER_CLASS` by default to match ICU.

Additionally, the behaviour of the word character construct `(\w)` is not consistent in ICU and Java 7. In Java 7 `\w` is `[\p{Alpha}\p{gc=Mn}\p{gc=Me}\p{gc=Mc}\p{Digit}\p{gc=Pc}]`, which is a larger set than ICU where `\w` is `[\p{Li}\p{Lu}\p{Lt}\p{Lo}\p{Nd}]`.

The use of `\w` is not recommended in DFDL regular expressions in conjunction with Unicode encodings, and an implementation must issue a warning if such usage is detected.

Character properties are detailed by the Unicode Regular Expressions [\[UnicodeRegex\]](#).

**25. Security Considerations**

All locations must be properly initialized before writing so as to prevent accidental (or purposeful) transmission of data in the unused parts of data formats. Even when a DFDL description does not specify that data should be written to a particular part of the output representation, a defined pattern should always be written.

When unparsing data it is a Schema Definition Error if the representation properties that control filling and padding are not defined by the DFDL schema. The DFDL processor must fail if they are not defined so that it is certain no region of the output data has unspecified contents.

If regions within a DFDL-described data object are encrypted, then when decrypting them proper means must be used to assure secure passage of passwords to the decrypting software. Such means are beyond the scope of the DFDL language specification.

In addition, if encryption passwords/keys are stored in DFDL schema-described data, then proper means must be used to assure that the decrypted form of these passwords is not revealed. Such means are beyond the scope of the DFDL language specification.

**26. Authors and Contributors**

Michael J. Beckerle, (corresponding author)

Tresys Technology/[Owl Cyber Defense](#)

Columbia, MD

USA

Email: [mbeckerle@tresys.com](mailto:mbeckerle@tresys.com), [mbeckerle.dfdl@gmail.com](mailto:mbeckerle.dfdl@gmail.com)

Stephen M. Hanson, (corresponding author)

IBM Software Group,

Hursley,

Winchester, UK

Email: [smh@uk.ibm.com](mailto:smh@uk.ibm.com)

Alan W. Powell,

Email: [apowell888@googlemail.com](mailto:apowell888@googlemail.com)

We greatly acknowledge the contributions made to this document by the following and all the other people who provided constructive and valuable input in the group discussions.

Tim Kimber, IBM Software Group, Hursley, UK

Suman Kalia, IBM Software Group, Markham, Ontario, Canada

Stephanie Fetzer, IBM Software Group, Charlotte, USA

Martin Westhead, Groupon, Stanford Univ, CA, USA

James Myers, Univ. of Michigan, MI, USA

Jonathan W. Cranford, MITRE Corporation, USA

Tom Sugden, UBS Investment Bank, UK

Tara Gibson, PNNL, Richland, WA, USA

Robert McGrath, NCSA, Urbana-Champaign, IL, USA

Geoff Judd, IBM Software Group, Hursley, UK

Dewey M. Sasser, MA, USA

David A. Loose, Opera Solutions, MA, USA

Eric S. Smith, eBay Enterprise, MA, USA

Kristoffer H. Rose, Two Sigma Investments, NY, USA

Simon Parker, Bloomberg Polar Lake, Dublin, IE

Peter A. Lambros, IBM Software Group, Hursley, UK

Dave Glick, SIPC, USA

Steve Marting, Progeny, USA

Alejandro Rodriguez, Amazon, WA, USA

**27. Intellectual Property Statement**

The OGF takes no position regarding the validity or scope of any intellectual property or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; neither does it represent that it has made any effort to identify any such rights. Copies of claims of rights made available for publication and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the OGF Secretariat.

The OGF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights which may cover technology that may be required to practice this recommendation. Please address the information to the OGF Executive Director.

**28. Disclaimer**

This document and the information contained herein is provided on an "As Is" basis and the OGF disclaims all warranties, express or implied, including but not limited to any warranty that the use of the information herein will not infringe any rights or any implied warranties of merchantability or fitness for a particular purpose.



**29. Full Copyright Notice**

Copyright (C) Open Grid Forum (2005-2014). Some Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included as references to the derived portions on all such copies and derivative works. The published OGF document from which such works are derived, however, may not be modified in any way, such as by removing the copyright notice or references to the OGF or other organizations, except as needed for the purpose of developing new or updated OGF documents in conformance with the procedures defined in the OGF Document Process, or as required to translate it into languages other than English. OGF, with the approval of its board, may remove this restriction for inclusion of OGF document content for the purpose of producing standards in cooperation with other international standards bodies.

The limited permissions granted above are perpetual and will not be revoked by the OGF or its successors or assignees.

ICU - Copyright (c) 1995-2014 International Business Machines Corporation and others

XPATH - Copyright © 2007 [W3C](#)® ([MIT](#), [ERCIM](#), [Keio](#)), All Rights Reserved. W3C [liability](#), [trademark](#) and [document use](#) rules apply.

### 30. References

- [ASN1] "Introduction to ASN.1",  
<http://www.itu.int/en/ITU-T/asn1/Pages/introduction.aspx>
- [ASN1CER] X.690 : Information technology - ASN.1 encoding rules: Specification of Basic Encoding Rules (BER), Canonical Encoding Rules (CER) and Distinguished Encoding Rules (DER)  
<https://www.itu.int/rec/T-REC-X.690-200811-l/en>
- [ASN1ECN] "ASN.1 Encoding Control Notation,"  
<http://www.itu.int/en/ITU-T/asn1/Pages/ecn.aspx>
- [ASN1PER] "ASN.1 Packed Encoding Rules (PER),"   
<http://www.itu.int/ITU-T/studygroups/com17/languages/X.691-0207.pdf>
- [AVRO] "Avro,"  
<http://avro.apache.org/docs/1.3.0/spec.html>
- [BFD] "Binary Format Description (BFD) Language,"  
<http://collaboratory.emsl.pnl.gov/sam/bfd/>
- [CARalia] Disc Interchange Service Company, Inc., Westford, MA, USA., "EBCDIC to ASCII Conversion of Signed Fields,"  
<http://www.3480-3590-data-conversion.com/article-signed-fields.html>
- [CCSID] "Coded Character Set Identifiers (CCSID),"   
[http://www.ibm.com/software/globalization/ccsid/ccsid\\_registered.html](http://www.ibm.com/software/globalization/ccsid/ccsid_registered.html)
- [EXI] W3C, "Efficient XML Interchange (EXI) Format 1.0 (Second Edition),"   
<http://www.w3.org/TR/exi>
- [HDF] Hierarchical Data Format (HDF)  
<http://www.hdfgroup.org/>
- [IANA] IANA, "Character Sets,"  
<http://www.iana.org/assignments/character-sets>
- [ICUDateTime] icu-project.org, "Formatting Dates and Times,"  
<http://userguide.icu-project.org/formatparse/datetime>
- [ICUDecimal] icu-project.org, "icu::DecimalFormat Class Reference,"

	<a href="http://icu.sourceforge.net/apiref/icu4c/classDecimalFormat.html#_details">http://icu.sourceforge.net/apiref/icu4c/classDecimalFormat.html#_details</a>
[ICULocale]	icu-project.org, "Locale," <a href="http://userguide.icu-project.org/locale">http://userguide.icu-project.org/locale</a>
[ICURegex]	icu-project.org, "Regular Expressions," <a href="http://userguide.icu-project.org/strings/regexp">http://userguide.icu-project.org/strings/regexp</a>
[ISO10646]	Universal Coded Character Set <a href="http://standards.iso.org/ittf/PubliclyAvailableStandards/c056921_ISO_IEC_10646_2012.zip">http://standards.iso.org/ittf/PubliclyAvailableStandards/c056921_ISO_IEC_10646_2012.zip</a>
[ISO8601]	Date and Time Format - ISO 8601 <a href="http://www.iso.org/iso/home/standards/iso8601.htm">http://www.iso.org/iso/home/standards/iso8601.htm</a>
[JavaRegex]	Oracle, "Class Pattern," <a href="http://docs.oracle.com/javase/7/docs/api/java/util/regex/Pattern.html">http://docs.oracle.com/javase/7/docs/api/java/util/regex/Pattern.html</a>
[OBSOLETE_DF DL]	Michael J Beckerle, Stephen M Hanson, Alan W Powell. <a href="#">GFD-P-R.174: Data Format Description Language (DFDL) v1.0 Specification</a> . Open Grid Forum. January 2011. <a href="http://www.oqf.org/documents/GFD.174.pdf">http://www.oqf.org/documents/GFD.174.pdf</a>
[IANATimeZone]	IANA - Internet Assigned Numbers Authority, "Time Zone Database," <a href="http://www.iana.org/time-zones">http://www.iana.org/time-zones</a>
[JSON]	Introducing JSON <a href="http://www.json.org">http://www.json.org</a>
[NETCDF]	Network Common Data Form (NetCDF) <a href="http://www.unidata.ucar.edu/software/netcdf/">http://www.unidata.ucar.edu/software/netcdf/</a>
[OMG_EAI]	OMG "CAM" TD Model: Object Management Group (OMG) "UML Profile and Interchange Models for Enterprise Application Integration (EAI) Specification" formal/04-03-26, March 2004. Section 7.3.2. Available at <a href="http://www.omg.org/cgi-bin/doc?formal/2004-03-26">http://www.omg.org/cgi-bin/doc?formal/2004-03-26</a> Available at <a href="http://www.omg.org/cgi-bin/doc?formal/2004-03-26">http://www.omg.org/cgi-bin/doc?formal/2004-03-26</a>
[RDP]	William H. Burge, "Recursive Programming Techniques," Addison-Wesley Longman, Incorporated, 1975. ISBN-10: 0201144506  , ISBN-13: 978-0201144505 0
[RFC2119]	S. Bradner, "RFC 2119: Key words for use in RFCs to Indicate Requirement Levels," IETF (Internet Engineering Task Force). <a href="#">RFC 2119: Key words for use in RFCs to Indicate Requirement Levels</a> . S. Bradner. 1997.

[SCD]	W3C, "W3C XML Schema Definition Language (XSD): Component Designators," <a href="http://www.w3.org/TR/xmlschema-ref/">http://www.w3.org/TR/xmlschema-ref/</a>
[Thrift]	M. Slee, A. Agarwal and M. Kwiatkowski, "Thrift: Scalable Cross-Language Services Implementation," <a href="http://thrift.apache.org/static/files/thrift-20070401.pdf">http://thrift.apache.org/static/files/thrift-20070401.pdf</a>
[UML]	Unified Modeling Language <a href="http://www.uml.org/">http://www.uml.org/</a>
[Unicode]	The Unicode Consortium, "Unicode," <a href="http://www.unicode.org/">http://www.unicode.org/</a> .
[UnicodeCLDR]	Unicode, Inc., "Unicode Common Locale Data Repository," <a href="http://sites.google.com/site/cldr/">http://sites.google.com/site/cldr/</a> .
[UnicodeRegex]	Unicode, Inc., "Unicode Regular Expressions," <a href="http://www.unicode.org/reports/tr18/">http://www.unicode.org/reports/tr18/</a>
[UnicodeLDML]	Unicode, Inc., "Unicode Locale Data Markup Language (LDML)," <a href="http://www.unicode.org/reports/tr35/">http://www.unicode.org/reports/tr35/</a> .
[Walmsley]	<u>Definitive XML Schema (Walmsley, ISBN 0-13-065567-8) page 390, Section 17.8, Table 17-6 says "A child element-type name which must be prefixed if it is in a namespace".</u>
[XDR]	RFC 4506: XDR: External Data Representation Standard IETF (Internet Engineering Task Force), 2006. <a href="http://tools.ietf.org/html/rfc4506">http://tools.ietf.org/html/rfc4506</a>
[XML1.0]	W3C, "Extensible Markup Language (XML) 1.0 (Fifth Edition)," 26 November 2008. <a href="http://www.w3.org/TR/REC-xml">http://www.w3.org/TR/REC-xml</a>
[XML1.1]	W3C, "Extensible Markup Language (XML) 1.1 (Second Edition)," 16 August 2006. <a href="http://www.w3.org/TR/xml11/">http://www.w3.org/TR/xml11/</a>
[XMLInfoSet]	W3C, "XML Information Set (Second Edition)," 4 February 2004. <a href="http://www.w3.org/TR/xml-infoSet">http://www.w3.org/TR/xml-infoSet</a>
[XMLNamespace s]	W3C, "Namespaces in XML 1.0 (Third Edition)," 8 December 2009. <a href="http://www.w3.org/TR/REC-xml-names/">http://www.w3.org/TR/REC-xml-names/</a>
[XSDL]	W3C, "XML Schema," <a href="http://www.w3.org/XML/Schema">http://www.w3.org/XML/Schema</a>
[XPath]	W3C, "XML Path Language (XPath) 2.0 (Second Edition)," <a href="http://www.w3.org/TR/xpath20/">http://www.w3.org/TR/xpath20/</a>
[XSDLV1]	W3C, XML Schema Part 1: Structures <a href="http://www.w3.org/TR/xmlschema-1/">http://www.w3.org/TR/xmlschema-1/</a> XML Schema Part 2: Datatypes <a href="http://www.w3.org/TR/xmlschema-2/">http://www.w3.org/TR/xmlschema-2/</a> , 28 October 2004.

### 31. Appendix A: Escape Scheme Use Cases

#### 31.1 Escape Character Same as dfdl:escapeEscapeCharacter

dfdl:escapeKind 'escapeCharacter', dfdl:escapeCharacter '/', dfdl:escapeEscapeCharacter '/',  
dfdl:separator ';', dfdl:extraEscapedCharacters '?'

Logical Data	Physical Data / Representation
.....	.....
...../.....	.....//.....
...../.../.....	.....//...//.....
.....//.....	.....///.....
/.....	//.....
...../	.....//
/...../.....	//.....//.....
...../...../	.....//.....//
.....;.....	...../;.....
...../;.....	.....//;.....
;.....	/;.....
.....?.....	...../?.....

Table 71 Examples of Escape Character Same as dfdl:escapeEscapeCharacter

#### 31.2 Escape Character Different from dfdl:escapeEscapeCharacter

dfdl:escapeKind 'escapeCharacter', dfdl:escapeCharacter '/', dfdl:escapeEscapeCharacter  
'%%'<sup>31</sup>, dfdl:separator ';', dfdl:extraEscapedCharacters '?'

Logical Data	Physical Data / Representation
.....	.....
...../.....	.....%/.....
...../.../.....	.....%/...%/.....
.....//.....	.....%/%/.....
/.....	%/.....
...../	.....%/
/...../.....	%/.....%/.....
...../...../	.....%/.....%/
.....;.....	...../;.....
...../;.....	.....%///;.....
;.....	/;.....
.....?.....	...../?.....

<sup>31</sup> The '%' character needs escaping when used as a literal character in DFDL properties. It is escaped by doubling it, i.e., preceding it by another '%' character.

.....%.....	.....%.....
.....%/.....	.....%%/.....
...../%.....	.....%/.....

**Table 72 Examples (1) of Escape Character Different from dfdl:escapeEscapeCharacter**

Dfdl:escapeKind 'escapeCharacter', dfdl:escapeCharacter '/', dfdl:escapeEscapeCharacter '%%',  
dfdl:separator 'sep', dfdl:extraEscapedCharacters '?'

Logical Data	Physical Data / Representation
.....sep.....	...../sep.....
...../sep.....	.....%/sep.....
sep.....	/sep.....

**Table 73 Examples (2) of Escape Character Different from dfdl:escapeEscapeCharacter****31.3 Escape Block with Different Start and End Characters**

dfdl:escapeKind 'escapeBlock', dfdl:escapeBlockStart '[', dfdl:escapeBlockEnd ']',  
dfdl:escapeEscapeCharacter '%%', dfdl:separator ';', dfdl:extraEscapedCharacters '?'

Logical Data	Physical Data / Representation
.....	.....
[.....]	[[.....]]
].....	].....
.....[.....	.....[.....
.....].....	.....].....
.....]	.....]
[[.....]	[[[.....]]
.....]]	.....]]
.....[[.....	.....[[.....
.....]].....	.....]].....
[.....]	[[.....%]]
[.....]	[[.....%].....]
.....[.....]	.....[.....]
[.....[.....]	[[.....[.....%]]
[.....]	[[.....%].....%]]
[[.....]	[[[.....%]]
[.....]]	[[.....%]%%]]
[[.....]]	[[[.....%]%%]]
.....%	.....%.....
.....%%.....	.....%%.....
.....%[.....	.....%[.....

.....%].....	.....%].....
%[.....	%[.....
.....%]	.....%]
%[.....%]	%[.....%]
[.....%.....]	[[.....%.....%]]
[.....%].....]	[[.....%%].....%]]
.....;.....	[.....;.....]
.....%;.....	[.....%;.....]
[.....;.....]	[[.....;.....%]]
.....?.....	[.....?.....]

Table 74 Examples of Escape Block with Different Start and End Characters

**31.4 Escape Block with Same Start and End Characters**

dfdl:escapeKind 'escapeBlock', dfdl:escapeBlockStart ' &apos;', dfdl:escapeBlockEnd '&apos;',  
dfdl:escapeEscapeCharacter '%', dfdl:separator ';', dfdl:extraEscapedCharacters '?'

Logical Data	Physical Data / Representation
.....	.....
'.....	'% '.....'
.....'	.....'
.....'	.....'
".....	'% '% '.....'
....."	....."
....."	....."
'.....'	'% '.....% '"
'.....'	'% '.....% '.....'
.....'	.....'
'.....'	'% '.....% '.....% '"
"....."	'% '% '.....% '"
'....."	'% '.....% '% '"
"....."	'% '% '.....% '% '"
.....%	.....%
.....%%.....	.....%%.....
.....%'.....	.....%'.....
% '.....	% '.....
.....% '.....	.....% '.....
'.....% '.....	'% '.....% '% '"
% '.....% '.....	% '.....% '.....

'.....%.....'	'% '.....%.....%''
'.....% '.....'	'% '.....% % '.....%''
.....;.....	'.....;.....'
.....%;.....	'.....%;.....'
'.....;.....'	'% '.....;.....%''
.....?.....	'.....?.....'

Table 75 Examples of Escape Block with Same Start and End Characters



### 32. Appendix B: Rationale for Single-Assignment Variables

DFDL is intended to be a description language. That is, the capture of a data format should be as descriptive/declarative as possible.

An additional quite critical goal for DFDL is that it allows very high performance implementations, including use of parallel processing wherever possible.

DFDL contains an expression language with variables for use in creating parameterized DFDL schemas.

However, the way variables can be used in DFDL is quite constrained. Specifically, the variables are single-assignment.

Single-assignment variables solve a number of problems.

First, they keep the schema more declarative, because the name of a variable represents a value, not a location. Before assignment, the value is not yet known, after the assignment the value is known, but the consumer of the value need only know the name, and need not be aware of the mechanism by which it gets its value or when.

Second, single-assignment variables avoid over-constraining the implementation, thereby preserving the potential for high-performance and parallel processing.

Some digression is useful here: Any variable creates a data dependency in order of processing. The part of the schema reading/using the variable's value depends upon the data value coming from the part of the schema providing that value. This kind of data dependency is inherent and inescapable. Values must be created before they can be used.

However, if you consider a variable to be a location that can be assigned repeatedly, then things are more complex because you not only have data dependency on the value (one part of the schema writes the location, another reads that location), but you have the dependency in the other direction: you must read the location before it can be used again for the *next* value. This is usually called anti-dependency. Anti-dependency is the enemy of high-performance and parallel execution. It forces specific and artificial sequential ordering on things that is due to the way variable names are allocated to storage locations.

If variables are single-assignment only, then only data-dependencies exist. Anti-dependencies don't exist, and implementations are free to work in any way consistent with the (inescapable) data dependencies.

### 33. Appendix C: Processing of DFDL String literals

A DFDL String Literal describes characters or bytes that can appear in the data. The syntax of a DFDL String Literal is described in Section 6.3. This appendix describes how a DFDL String Literal must be processed by a DFDL processor.

#### 33.1 Interpreting a DFDL String Literal

A DFDL String literal is a sequence of one or more *string literal parts*. Each string literal part describes zero or more characters or exactly one byte. To process a DFDL string literal correctly it is necessary to parse it into its string literal parts using the grammar in Table 2 DFDL Character Entity, Character Class Entity, and Byte Value Entity Syntax of Section 6.3.1 DFDL String Literals.

#### 33.2 Recognizing a DFDL String Literal

When parsing, a DFDL processor must be able to recognize a DFDL String Literal in the data. The recognition algorithm is used whenever the DFDL processor needs to process a property of type 'DFDL String Literal'.

The recognition algorithm is as follows:

```
for each DfdlStringLiteralPart
  if this DfdlStringLiteralPart is found in the data at the current offset
    advance the offset by the width of the characters matched
  else
    return false
  end if
end if
return true
```

#### 33.3 Recognizing DFDL String Literal Part

Each type of string literal part and/or the DfdlESEntity is recognised as described in the table below:

String Literal Part	Recognition algorithm
LiteralString	The literal string is specifying a sequence of Unicode characters. The string part is recognized if the same Unicode characters are found in the data at the current offset. Implementation note: It is not necessary for either the LiteralString or the data to be encoded using a Unicode encoding. The matching algorithm must operate as if the LiteralString and the data were both converted to Unicode before comparison.
DfdlCharEntity	The entity is specifying a single Unicode character. The recognition algorithm is the same as for a single-character LiteralString.
ByteValue	The entity is specifying a single raw byte value. The string part is recognized if the literal byte in the data at the current byte offset matches the raw byte value. This entity is specifying a raw byte value and not a character. It will not always be safe to interpret the byte value as a character in the component's encoding.
DfdlNLEntity	The entity is specifying a newline character or character sequence. The string part is recognized if a newline character sequence is found in the data at the current byte offset: The newline character sequences are defined in <a href="#">Table 4 DFDL Character Class Entities</a> , in Section 6.3.1.2 DFDL Character Entities, Character Class Entities, and Byte Values in String Literals.

Deleted: Table 4 DFDL Character Class Entities

	The recognition algorithm must be greedy. If the CR character is matched then the DFDL processor must check to see whether it is followed by an LF character and if so, consume that character also.
DfdIWSPEntity	The entity is specifying a single whitespace character. The string part is recognized if one of the whitespace characters is found in the data at the current byte offset. The whitespace characters are defined in <a href="#">Table 4 DFDL Character Class Entities</a> , in Section 6.3.1.2 DFDL Character Entities, Character Class Entities, and Byte Values in String Literals.
DfdIWSPStarEntity	The entity is specifying zero or more whitespace characters. The string part is recognized when the entire sequence of whitespace characters starting at the current offset has been consumed. The recognition of this entity cannot fail because it can match zero whitespace characters. The recognition algorithm must be greedy. All possible whitespace characters must be consumed.
DfdIWSPPlusEntity	The entity is specifying one or more whitespace characters. The string part is recognized when the entire non-empty sequence of whitespace characters starting at the current offset has been consumed. The recognition algorithm must be greedy. All possible whitespace characters must be consumed.
DfdIESEntity	The entity is specifying an empty string. The string part is recognized if the data available for matching is zero-length. This entity is only allowed in contexts where the available data will be constrained by other DFDL properties.

Deleted: Table 4 DFDL Character Class Entities

**Table 76 Recognizing DFDL String Literal Part**

### 34. Appendix D: DFDL Standard Encodings<sup>32</sup>

This appendix provides a specification template for the definition of DFDL standard encodings, and defines two such encodings for use in DFDL schemas.

These encodings are optional - DFDL implementations are not required to support them to be considered to be in conformance with the DFDL specification.

However, if implemented, then a conformant DFDL implementation must follow the specification in this appendix.

#### 34.1 Purpose

Additional encoding names are needed for several reasons:

1. **Incomplete Information:** the definitions of IANA [\[IANA\]](#) and CCSID [\[CCSID\]](#) character set encodings are not sufficiently specific about the encoding, omitting important characteristics such as alignment or bit order. An example of this is the encoding X-DFDL-US-ASCII-7-BIT-PACKED, defined herein.
2. **Unnamed or Legacy Encodings:** There are many legacy encodings that have no IANA or CCSID identifiers. For example ITA2 [\[ITA2\]](#) is a 5-bit-per character encoding that is very old, yet still has modern uses.

#### 34.2 Conventions

All indexing of positions are 1-based meaning the first position is position 1.

Bytes are unsigned integers with values from 0 to 255 (decimal).

Hexadecimal numbers have the prefix "0x".

Decimal numbers have no prefix when the context is clear, and are explicitly labeled with (decimal) to resolve ambiguity.

#### 34.3 Specification Template

A DFDL standard encoding specification must include:

1. **Name:** All DFDL standard encoding names will begin with the prefix "X-DFDL-". Note that encoding names are not case-sensitive.
2. **Translation table:** An encoding/decoding algorithm that provides a complete bi-directional mapping between encoded characters and UTF-16 character codes. Many encodings share character code definitions and describe only different ways of representing those values; hence, this section may reference other IANA, CCSID or DFDL standard encoding definitions.
3. **Width:** Whether the character set is fixed or variable width. If fixed width, the size of a code unit in bits or bytes.
4. **Alignment:** Alignment and alignment units for a code unit.
5. **Byte order:** When a code unit consists of more than one byte of data, the byte order of the code units. Note that the `dfdl:byteOrder` property does not apply to text.

#### 34.4 Encoding X-DFDL-US-ASCII-7-BIT-PACKED

This encoding is used by MIL-STD-2045 [\[MILSTD2045\]](#) and a variety of related binary data standards.

##### 34.4.1 Name

X-DFDL-US-ASCII-7-BIT-PACKED

##### 34.4.2 Translation table

The characters are exactly the 128 US-ASCII characters with the same code unit values.

---

<sup>32</sup> It is anticipated this appendix will be migrated into a separate document to facilitate the addition of new encodings without revising the DFDL specification.

### 34.4.3 Width

Fixed width. Each code unit is 7-bits wide.

This encoding differs from US-ASCII in that it specifically uses only 7-bits per character within each byte<sup>33</sup>. A subsequent character or the next data field may begin in the very next bit after a 7-bit character code of this encoding.

Hence, 8 characters will fit exactly into 7 bytes of data as no bits are wasted.

### 34.4.4 Alignment

1 bit. That is, a code unit may begin on any bit boundary.

### 34.4.5 Byte Order

Not applicable - code units are always smaller than 1 byte.

### 34.4.6 Example 1

The table below shows how the 8 character string 'UNIT1234' is encoded into 7 bytes, when the bit order is least-significant-bit first.

Logical character	7-bit code unit	Bit placement 2 <sup>7</sup> 2 <sup>0</sup>	Byte value	Byte number
U	1010101	X1010101		
N	1001110	0XXXXXXX XX100111	0x55	1
I	1001001	01XXXXXX XXX10010	0x67	2
T	1010100	100XXXXX XXXX1010	0x92	3
1	0110001	0001XXXX XXXXX011	0x1A	4
2	0110010	10010XXX XXXXXX01	0x93	5
3	0110011	110011XX XXXXXXX0	0xCD	6
4	0110100	0110100X	0x68	7

**Table 77: Example of DFDL Standard Encoding X-DFDL-US-ASCII-7-BIT-PACKED**

Below is the same data shown horizontally with the bytes numbered from right to left. The bits corresponding to the characters are shown in distinct colors. The characters that each bit contributes to are written below each of the bits and in the corresponding color.

	01010101	01100111	10010010	00011010	10010011	11001101	01101000
Byte	-----1--	-----2--	-----3--	-----4--	-----5--	-----6--	-----7--
Char	UUUUUUU	IINNNNN	TTTTIII	1111TTTT	2222111	3333322	4444443

This kind of data is often better visualized by numbering the bytes in the opposite direction that is *right to left*, which leads to this presentation:

	01101000	11001101	10010011	00011010	10010010	01100111	01010101
Byte	-----7--	-----6--	-----5--	-----4--	-----3--	-----2--	-----1--

<sup>33</sup> Though US-ASCII contains only 128 characters DFDL assumes the US-ASCII encoding to use one character per 8-bit byte.

Char 44444443 33333322 22222111 1111TTTT TTTIIIII IINNNNNN NUUUUUUU

In the above, the bits corresponding to each character code unit are more easily recognized, but the characters appear right-to-left (i.e., backward for English).

34.4.7 Example 2

The bits below represent a 3-bit unsigned integer containing value 7, followed by the ASCII string 'ABC' followed by the ASCII DEL character (character code 0x7F). This illustrates a string not beginning on a byte boundary. Again the bit ordering is least-significant-bit first.

The bits are written in increasing position and place value from *right-to-left*:

1111111	1000011	1000010	1000001	111
DEL	C	B	A	7

In the above example, if we number the bits from the right, starting with position 1, the character code for 'A' occupies bits 4 to 10. If we divide the data above into bytes with vertical bars we must start on the right to get:

char:	DEL	C	B	A
bits:	01111111	1000011 1	000010 10	00001 111
byte:	7F	87	1C	0F
byte pos:	4	3	2	1

### 34.5 Encoding X-DFDL-US-ASCII-6-BIT-PACKED

This encoding is used by MIL-STD-6016D (See [\[MILSTD6016\]](#)) and a variety of related binary data standards.

#### 34.5.1 Name

X-DFDL-US-ASCII-6-BIT-PACKED

#### 34.5.2 Translation Table

The characters are a subset of the 128 US-ASCII characters. The assignment of numeric code units to the characters is different for some of the characters and the same for others.

The correspondence of ASCII characters to their code point values is given by the table below.

Character	Code Point Value (Decimal)	Notes
@	0	These characters' code points are the same as the US-ASCII code point, minus 64 (decimal)
A to Z	1 through 26	
[	27	
\	28	
]	29	
^	30	
_	31	These characters' code points are the same as in US-ASCII.
SPACE	32	
!	33	
"	34	
#	35	
\$	36	
%	37	
&	38	
' APOSTROPHE	39	
(	40	
)	41	
*	42	
+	43	
, COMMA	44	
-	45	
. PERIOD	46	
/	47	
0 to 9	48 through 57	
:	58	
;	59	

<	60	
=	61	
>	62	
?	63	

**Table 78: Translation Table for DFDL Standard Encoding X-DFDL-US-ASCII-6-BIT-PACKED****34.5.3 Width**

Fixed width. Each code unit is 6-bits wide. A subsequent character or the next data field may begin in the very next bit after a 6-bit character code of this encoding.

Hence, 4 characters will fit into 3 bytes of data as no bits are wasted.

**34.5.4 Alignment**

1 bit. That is, a code unit may begin on any bit boundary.

**34.5.5 ByteOrder**

Not applicable. Code units are always less than 1 byte in size.

**34.5.6 Example 1**

The table below shows how the 8 character string '1234' is encoded into 3 bytes. The bit order is least-significant-bit first:

Logical character	6-bit code unit	Bit placement 2 <sup>7</sup> 2 <sup>0</sup>	Byte value	Byte number
1	110001	XX110001		
2	110010	10XXXXXX XXXX1100	0xB1	1
3	110011	0011XXXX XXXXXX11	0x3C	2
4	110100	110100XX	0xD3	3

**Table 79: Example of DFDL Standard Encoding X-DFDL-US-ASCII-6-BIT-PACKED**



### 34.6 References for Appendix D<sup>34</sup>

- [CCSID] "Coded Character Set Identifiers (CCSID)," [http://www.ibm.com/software/globalization/ccsid/ccsid\\_registered.html](http://www.ibm.com/software/globalization/ccsid/ccsid_registered.html)
- [DFDL] Michael J Beckerle, Steven M Hanson, Alan W Powell. Data Format Description Language (DFDL) v1.0 Specification. Open Grid Forum. (<http://redmine.ogf.org/dmsf/dfdl-wg>) GFD-P-R.207 (2014)
- [IANA] IANA, "Character Sets," <http://www.iana.org/assignments/character-sets>
- [ITA2] Wikipedia: Baudot code, [http://en.wikipedia.org/wiki/Baudot\\_code#ITA2](http://en.wikipedia.org/wiki/Baudot_code#ITA2)
- [MILSTD2045] CONNECTIONLESS DATA TRANSFER APPLICATION LAYER STANDARD, MIL-STD-2045-47001D w/CHANGE 1, 23 June 2008 (available publicly from US Dept. of Defense at <http://assistdocs.com/>)
- [MILSTD6016] Tactical Data Link Message Standard MIL-STD-6016  
This standard requires authorization to access.  
[http://en.wikipedia.org/wiki/Link\\_16](http://en.wikipedia.org/wiki/Link_16) provides some information about the standard.  
[http://www.assistdocs.com/search/document\\_details.cfm?ident\\_number=123964](http://www.assistdocs.com/search/document_details.cfm?ident_number=123964) provides links for those with authorization to access the various versions of the standard.

---

<sup>34</sup> This list is maintained separately from the other references in the specification to facilitate future migration of this material into a separate document.