Data Format Description Language (DFDL) v1.0 Errata (Internal Committee Working Document)

Status of This Document

This working draft document provides errata information to the OGF community on the Data Format Description Language (DFDL) 1.0 specification (GFD-P-R.174). Distribution is unlimited.

Copyright Notice

Copyright © Open Grid Forum, (2011, 2012). All Rights Reserved.

Abstract

This document lists the non-editorial errata identified in the DFDL 1.0 specification.

Version

v01<u>1</u> updated 2012-<u>11</u>-<u>16</u>.

Table of Contents

<u>1</u>	Introduction	<u>3</u>
2.	Minor Technical Fixes	4
3.	Major Technical Fixes	17
4.	Appendix A. Grammar	28
5.	Authors and Contributors	30
6.	Intellectual Property Statement	<mark>31</mark>
7.	Disclaimer	32
8.	Full Copyright Notice	33
9.	References	34

1. Introduction

This document lists the non-editorial errata identified in the DFDL 1.0 specification [DFDL]. This document is a working draft and is updated each time an errata is identified.

The OGF GFD process [GFD] recognises three different kinds of error that may be found in OGF specifications:

Editorial fixes. Updates to a document which are not widely announced or publicized. This category might include headers/footers, spelling, formatting, or simple wording changes for clarity.

Minor technical fixes. Updates to a document which are not simply editorial. For example, an update to an XML schema or addition to a protocol, to bring the document into agreement with current practice.

Major technical fixes. Such fixes will often require additional technical review and result in an updated or replaced document.

The following sections list the errata that fall into the last two categories.

2. Minor Technical Fixes

The following minor technical fixes have been identified.

- **2.1.** Section 7.2.2. The ref property needs to state that circular paths are a schema definition error.
- 2.2. Section 13.13. Clarify what packed and BCD calendars mean.
 - There is no need to use a separate VDP property. The only place where a decimal point can occur is for fractional seconds. This is detectable from the pattern at the boundary of 's' and 'S', ie sS.
 - Property calendarPatternKind = 'explicit' must be used with binary calendar representations, as the defaults for 'implicit' use non-numeric characters. Schema definition error otherwise.
 - Property binaryCalendarRep should restate the rule from property calendarPattern.
 - Examples to be provided.
- **2.3.** Section 13.11.1. Does not fully state the time zone symbol behaviour, as spec quotes from http://icu-project.org/apiref/icu4c/classSimpleDateFormat.html instead of from http://userguide.icu-project.org/formatparse/datetime. It should say:

Z	Time Zone: specific non-location (z, zz, or zzz)	PDT
ZZZZ	Time Zone: specific non-location	Pacific Daylight Time
Z	Time Zone: RFC 822 (Z, ZZ, ZZZ)	-0800
ZZZZ	Time Zone: localized GMT	GMT-08:00
٧	Time Zone: generic non-location	PT
VVVV	Time Zone: generic non-location	Pacific Time
V	Time Zone: generic non-location	PT
VVVV	Time Zone: generic location	United States (Los Angeles)

Note that both Z and ZZZZ can be followed by 'U' which is a DFDL extension.

- **2.4.** Sections 22.1.1 & 22.2.1. Binary representations can have property lengthKind set to 'delimited'.
- **2.5.** Sections 22.1.2 & 22.2.2. Complex elements can have property lengthKind set to 'endOfParent'.
- **2.6.** Sections various. Spec often uses the term 'content region' but it should be more specific in terms of the grammar, and use 'SimpleContent region' or 'ComplexContent region', or both.
- **2.7.** Section 17. Text says that inputValueCalc and outputValueCalc applies to simple types, which is not correct.
- **2.8.** Section 13.16. In the description of property nilValue, state that nilLiteralCharacter test takes place on the untrimmed representation value.
- **2.9.** Section 12.3.3. Clarify that when property lengthUnits is 'bytes', using property lengthKind 'implicit' for a string interprets the min/maxLength facets as byte values and not characters when parsing or unparsing.
- **2.10.** Section 6.3. Bullet for logical value. State that the string must obey the lexical representation of the type.
- **2.11**. Section 6.3. Clarify that literal white space is only ever used as list token separator, and that entities must be used if literal white space is needed as part of the property value.

- **2.12**. Section <u>6</u>, <u>7.7</u>. Clarify that if A.xsd includes B.xsd then A can refer to a variable defined in B and reference is via QName in the usual way. <u>This is best expressed by simply saying that DFDL QNames behave like XSDL QNames in section 6. The existing text in section <u>7.7</u> can be removed.</u>
- **2.13.** Sections 9.2. Correct the grammar to reflect that a prefix length type can itself have a prefix length. This is sufficient to allow the grammar to describe the needed "one more level" of prefix (as required for modeling an ASN.1 format) without allowing recursion.

The updated grammar is in Appendix A of this document.

2.14. Section 12.3.4. Clarify that when a prefix length type itself has a prefix length, the simple types can not be the same.

Explicitly list the property restrictions that must apply to a prefix length type to comply with modeling just SimpleContent region. It is a schema definition error if the type specifies lengthKind 'delimited' or 'endOfParent' or 'explicit' where length is an expression, or a value for initiator or terminator other than empty string, or alignment other than '1', or leadingSkip or trailingSkip other than '0'.

- **2.15.** Section 13.6. When property textNumberRep is 'zoned', the property description should state that base is assumed to be 10.
- **2.16.** Section 13.2.1. Clarify string literal content of properties escapeCharacter, escapeEscapeCharacter, extraEscapedCharacters, escapeBlockStart and escapeBlockEnd.
 - o DFDL character entities are allowed
 - o The raw byte entity (%#r) is not allowed
 - o DFDL Character classes (NL, WSP, WSP+, WSP*, ES) are not allowed
- **2.17.** Sections 13.4, 13.6, 13.9, 13.12. Clarify string literal content of properties textStringPadCharacter, textBooleanPadCharacter, textCalendarPadCharacter and textNumberPadCharacter.
 - DFDL character entities are allowed
 - The raw byte entity (%#r) is allowed subject to the restrictions already documented for these properties
 - o DFDL Character classes (NL, WSP, WSP+, WSP*, ES) are not allowed
- **2.18.** Section 13.6. Clarify string literal content of properties textStandardDecimalSeparator, textStandardGroupingSeparator, textStandardExponentCharacter, textStandardInfinityRep, textStandardNaNRep.
 - DFDL character entities are allowed
 - The raw byte entity (%#r) is not allowed
 - DFDL Character classes (NL, WSP, WSP+, WSP*, ES) are not allowed
- **2.19.** Section 13.9. Clarify string literal content of properties textBooleanTrueRep and textBooleanFalseRep.
 - o DFDL character entities are allowed
 - The raw byte entity (%#r) is not allowed
 - o DFDL Character classes (NL, WSP, WSP+, WSP*, ES) are not allowed
- **2.20.** Section 13.16. Remove restriction that property nilValue only applies when representation is text. It is not clear where this originated.

Clarify string literal content of nilValue when nilKind is 'literalValue': When representation is text:

- o DFDL character entities are allowed
- o The raw byte entity (%#r) is allowed
- o DFDL Character classes (NL, WSP, WSP+, WSP*, ES) are allowed.

When representation is binary:

- DFDL character entities are allowed
- The raw byte entity (%#r) is allowed
- DFDL Character class ES is allowed.
- Other DFDL Character classes (NL, WSP, WSP+, WSP*) are not allowed.

Clarify string literal content of nilValue when nilKind is 'literalCharacter':

When representation is text:

- o DFDL character entities are allowed
- The raw byte entity (%#r) is allowed subject to the restrictions already documented for this property
- o DFDL Character classes (NL, WSP, WSP+, WSP*, ES) are not allowed.

When representation is binary:

- DFDL character entities are allowed
- o The raw byte entity (%#r) is allowed.
- o DFDL Character classes (NL, WSP, WSP+, WSP*, ES) are not allowed.
- **2.21.** Section 13.6. Change meaning of textNumberCheckPolicy enum 'lax' to align with ICU_: "If 'lax' and dfdl:textNumberRep is 'standard' then grouping separators are ignored, leading and trailing whitespace is ignored, leading zeros are ignored, quoted characters may be omitted."

Additionally state the base behaviour of textNumberCheckPolicy enum 'strict' to align with ICU. "If 'strict' and dfdl:textNumberRep 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."

- **2.22.** Section 13.6. Disallow the use of empty string for property textStandardDecimalSeparator, and state property must be set if the pattern contains a '.' or '@' symbol (schema definition error otherwise).
- **2.23.** Section 13.6. Allow decimal separator to be a List of DFDL String Literals or a DFDL expression. This allows modelling of the EDIFACT standard where a user can choose a dynamic decimal separator in the ISA header but '.' is always allowed.
- **2.24.** Section 13.6. Disallow the use of empty string for property textStandardGroupingSeparator, and state property must be set if the pattern contains a ',' (schema definition error otherwise).
- **2.25.** Section 13.6. When property textNumberPadCharacter is '0' which it commonly is, a value of say '00000' will get trimmed to the empty string, whereas the intent is to trim to '0'. Add a new rule that says the last remaining digit is never trimmed for text numbers regardless of its value.
- **2.26.** Section 13.6. Allow the use of empty string for property textStandardExponentCharacter to model text numbers of the form nnn+mmm. Property must be set even if the pattern does not contain an 'E' symbol, to match ICU behaviour (schema definition error otherwise).
- **2.27.** Section 13.6. textStandardDecimalSeparator must be ignored when the logical type is not decimal/float/double.
- **2.28.** Section 13.6.1.1. Add support for ICU significant digits symbol '@'. Note that this is not needed as a change in 13.6.1.2.
- **2.29.** Section 13.6.1.1. Formatting. Uses terms 'minimum/maximum integer/fraction digits' but does not define them. The term 'maximum integer digits' is defined as 309 to match the ICU default, the other terms are defined by the pattern content.
- **2.30.** Section 13.9. State that textBooleanTrueRep and textBooleanFalseRep properties are used after trimming when parsing, and before padding when unparsing. If lengthKind is

'explicit' or 'implicit' and either textPadKind or textTrimKind is 'none' then the properties must have the same length else it is a schema definition error.

- **2.31.** Section 16.2. State it is a processing error if the stop value is missing from the data when parsing.
- **2.32.** Section 12.1.1. Clarify the note after Table 14 mean. "Specifying the implicit alignment in bits does not imply that dfdl:lengthUnits 'bits' can be specified for all simple types". It is really saying that alignmentUnits and lengthUnits are independent and have their own rules for when they are applicable.
- 2.33. Section 12.3.. One line descriptions of 'delimited' and 'endOfParent' are not worded correctly in the property description of lengthKind, and should be improved.

This errata is currently under review in order to provide more detail.

- **2.34.** Section 12.3.2. Rule 3 for resolving ambiguity between delimiters, which says "When the separator and terminator on a group have the same value, the separator has precedence", needs clarifying to say "When the separator and terminator on a group have the same value, then at a point where either separator or terminator could be found, the separator is tried first."
- **2.35.** Section 17. InputValueCalc. Description talks about returning an empty string being ok if minLength permits this. Replace sentence with a fuller clarification that inputValueCalc value is validated like a parsed value, so schema definition error if value does not conform to base type, and validation error if validation enabled and value conforms to base type but not actual type.
- **2.36.** Section 16. Spec allows occursCount to be a non-negative integer. This is superfluous as the property is only used when occursCountKind is expression. Change so that occursCount is only allowed to be a DFDL Expression.
- **2.37.** Section 23.3. Clarify that DFDL expression syntax "{}" is invalid, as it results in an empty XPath 2.0 expression, which is not legal. In particular, clarify that setting a property to {} does not give the same result as setting a property to the empty string.
- **2.38**. Section 11. Specification does not definitively list which binary reps are subject to byteOrder. Clarify that byteOrder applies to all Numbers and Calendars with representation binary. Specifically that is binary integers, packed decimals, BCD, binary floats, binary seconds and binary milliseconds.
- **2.39.** Section 13.11.1. When parsing an xs:date or xs:datetime, if a calendarPattern doesn't specify some parts (other than time zone), say, calendarPattern="MM", then the Unix epoch 1970-01-01T00:00:00:00.000 is used to provide the missing parts.

Noted that if a pure month or day or year is needed, then this would be achieved by a future DFDL extension to expand the supported simple types to include xs:gMonth, xs:gDay, xs:gYear types.

- **2.40**. Section 13.6.1.1. The paragraphs that describe the V symbol (virtual point) and P symbol (scaling factor) talk about 'number region'. This relates to the BNF and so it should say 'vpinteger region' instead, which is where in the pattern the V and P symbols reside. Table 20 should also be updated so that it matches the BNF.
- **2.41**. Section 13.6.1.2. This section does not explicitly say that its content is effectively a delta on section 13.6.1.1. The section should be rewritten to make it clear which behaviour is the same as 13.6.1.1 and which is different.
- 2.42. Section 13.6. Clarify string literal content of property textStandardZeroRep.

- DFDL character entities are allowed
- The raw byte entity (%#r) is not allowed
- DFDL Character classes (NL, ES) are not allowed
- DFDL Character classes (WSP, WSP+, WSP*) are allowed
- **2.43**. Section 13.11. Property calendarTimeZone is defined as an Enum of type string, but in reality is better defined as a String constrained by a regular expression: $(UTC)([+-]([01]\d|\d)((([:][0-5]\d)\{1,2\})?))?)$ See also errata 2.50 and 2.65.
- **2.44.** Section 13.11. Property calendarLanguage is defined as an Enum of type string, but in reality is better defined as a String constrained by a regular expression: $([A-Za-z]\{1,8\})([-][A-Za-z0-9]\{1,8\})^*)$
- **2.45**. Sections 13.17 and 22. State that property useNilForDefault is only examined when xs:nillable is "true", and must be set when xs:nillable is "true".
- **2.46.** Section 13.6. Allow multiple characters for property textStandardExponentCharacter to handle representations like 1.23x10^4 as ICU allows that. Note that property name will therefore change to textStandardExponentRep.
- **2.47**. Section 13.6. Change name of property textStandardNanRep to textStandardNaNRep to reflect common usage of NaN and avoid typographical errors in models.
- **2.48**. Section 14.1. Spec states that an empty sequence that is the content of a complex type is a schema definition error. Many schema processors are not able to distinguish this condition from a complex type with no content at all (it is not required to do so by the XML Schema specification). As a complex type with no content is not useful in DFDL, change the spec to state that both conditions are schema definition errors.
- **2.49.** Section 23.3. Clarify that when a property can be either a DFDL String Literal or a DFDL Expression, then if the value is a DFDL String Literal and the first character is '{' then it MUST be escaped as '{{'}. For such a property, a value '{xxx' will be treated as an (invalid) expression, and not as a string literal.
- 2.50. Section 13.11. The calendarTimeZone property is used to supply a time zone when there is none in the data (and by implication none in the pattern). However this means DFDL is not compatible with XML Schema 1.0 where "no time zone" is an allowable state for a calendar infoset value. Further, XML Schema 1.0 validation validates a calendar value against facets according to http://www.w3.org/TR/2004/REC-xmlschema-2-20041028/#dateTime-order which has rules that cater for "no time zone". It is desirable therefore for DFDL to permit a calendar value to have "no time zone". Accordingly, the calendarTimeZone property will allow
- **2.51.** Section 13.11. The calendarTimeZone property will apply when parsing only. This avoids any problem with values changing after validation has taken place.

a value of empty string to indicate "no time zone".

- **2.52.** Section 13.11.1. DFDL only allows 'y' as year symbol in dfdl:calendarPattern. This only allows positive values, any negative value is ignored unless the 'G' (era) symbol is also specified, and there is no year 0 (which means that negative astronomical dates are one year out). DFDL will also support the ICU 'u' extended year symbol which allows year 0 (means 1BC), corresponds to astronomical years and matches XML Schema 1.0 (post errata for year 0).
- **2.53**. Section 13.13. Property binaryCalendarEpoch is used when the binaryCalendarRep is either binarySeconds or binaryMilliseconds, and is of type xs:dateTime. It is allowable to omit the time zone component from the binaryCalendarEpoch property value, and if this occurs UTC is used as the time zone.

2.54. Section 13.11. ICU has some lax behaviour when parsing numbers using the supplied textNumberPattern. The pattern is primarily used to control formatting. The specification should make clear what this lax behaviour is.

This errata is currently under review in order to provide more detail.

2.55. Section 9.2. Grammar terminal FinalUnusedRegion is intended to handle unmodeled bytes in the data that arise due to 'specified length' settings of lengthKind on a complex xs:element and choiceLengthKind on a xs:choice. It is not doing so correctly when a terminator is present on a xs:sequence or xs:choice. To fix this, the terminal is removed from the grammar and replaced by two new terminals ElementUnused and ChoiceUnused.

The <u>updated</u> grammar <u>is in Appendix A of this document.</u>

- **2.56.** Section 13.13. State that when property binaryCalendarRep is set to 'binaryMilliseconds' or 'binarySeconds', it is a schema definition error if the type is xs:time or xs:date. This is because when unparsing it is not possible to obtain a milliseconds or seconds value from just an xs:time or xs:date and the epoch.
- **2.57**. Section 13.13. State that when property binaryCalendarRep is set to 'binarySeconds' or 'binaryMilliseconds', the value in the data is treated as signed. This lets DFDL support POSIX/Unix times, which are allowed to be negative.
- **2.58**. Section 13.16. State that property nilValueDelimiterPolicy is ignored when property nilKind is set to 'logicalValue', and the behaviour of the DFDL processor is to expect delimiters when parsing and to output delimiters when unparsing, if delimiters are specified for the element. This is to simplify implementations.
- **2.59**. Section 7.1.3.3. State that short form property syntax is not allowed on the xs:schema object as an equivalent to the element form property syntax of the default dfdl:format (or any other global DFDL) annotation.
- **2.60**. Section 13.3. Remove redundancy from the names of some of the bidirectional text properties, specifically textBidiTextOrdering becomes textBidiOrdering, and textBidiTextShaped becomes textBidiShaped.
- **2.61.** Section 9.2. The child content of xs:sequence and xs:choice are almost the same, so the grammar can be refactored to remove duplication.

The updated grammar is in Appendix A of this document.

- **2.62.** Section 11. Clarify that when the encoding property is specified as 'UTF-8' then that is a strict definition of UTF-8 and does not include variants such as CESU-8. This is in keeping with ICU's interpretation of UTF-8.
- **2.63**. Sections 12.2, 14.2. Properties that use the empty string as a special value to switch off use of the property, and that allow the value to be a DFDL expression, should not be able to set empty string by evaluating the expression. It must be possible to evaluate statically whether the property is used or not. This affects properties initiator, terminator, separator. It is a schema definition error if an expression returns the empty string.
- **2.64**. Section 6.3.1. Update to say that empty string is not allowed as a string literal value, unless explicitly stated otherwise in the description of a property, in which case any semantic must be a special behaviour of the property and not the literal empty string (for which DFDL provides entity %ES;).
- **2.65**. Section 13.11. Property calendarTimezone is changed to accept either a UTC offset or an Olson format time zone. Property calendarObserveDST is changed so that it is only used

when calendarTimeZone is Olson format. If it is a UTC offset then calendarObserveDST is ignored.

- **2.66.** Section 13.11. When unparsing and property calendar Pattern contains a formatting symbol for time zone (zzz, Z, VVVV etc) and the infoset value does not contain a time zone, it is a processing error. This matches the behaviour on parsing when the data does not contain a time zone but the pattern does.
- **2.67**. Section 4.1.2. The second paragraph of the description of [dataValue] should be replaced with:

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

2.68. Section 13.11.1. Calendar formatting symbols 'l' and 'T' are intended as a short-hand way of accepting a subset of ISO 8601 variants. However the description of the behaviour is not correct and is unnecessarily complex. The 'T' symbol is dropped altogether, and the 'l' symbol behaviour is defined as the following:

"The 'I' 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. 'IU' is permitted for xs:date but
 the 'U' is ignored as there is no time zone in the date-only granularities.
- 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 implementation dependent but must be at least one.
- For a granularity that omits components, when parsing the values for the omitted components are supplied from the Unix epoch 1970-01-01T00:00:00.000."
- **2.69**. Section 23.3. The last paragraph is inconsistent with the rest of the section. It should say:

"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."

The sentence "If the expression returns an empty sequence it will be treated as returning nil" is removed.

2.70. Section 12.2, 14.2. Clarify the parser matching algorithm used for properties initiator, terminator and separator.

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 white space separated list, and matches them each against the data. In each case the longest possible match is found. The initiator with the longest match is the one that is selected as having been 'found', with length-ties being resolved so that the matching initiator is selected that is first in the order written in the schema. Once a matching initiator is found, no other matches will be subsequently attempted (ie, there is no backtracking).

- **2.71**. Section 13.16. Clarify that nilValue is sensitive to ignoreCase when nilKind is 'literalValue' or 'logicalValue', to be consistent with properties such as textBooleanTrueRep, but not when nilKind is 'literalCharacter', to be consistent with properties such as textBooleanPadCharacter.
- **2.72**. Section 12.3.6. Additional constraints and clarifications apply to the use of lengthKind 'endOfParent' beyond those already documented:

The parent element lengthKind must not be 'implicit' or 'delimited'.

When looking for end of parent, the parser is not sensitive to any in-scope terminating delimiters.

If the element is in a sequence then:

- o the sequence must be the content of a complex type
- the separatorPosition of the sequence must not be 'postFix'
- the sequenceKind of the sequence must be 'ordered'
- o no terminator on the sequence
- no trailingSkip on the sequence
- o no floating elements in the sequence

If the element is in a choice where choiceLengthKind is 'implicit' then

- o the choice must be the content of a complex type
- no terminator on the choice
- o no trailingSkip on the choice

A simple element must have either type xs:string or representation 'text' or type xs:hexBinary or (representation 'binary' and binaryNumber/CalendarRep 'packed', 'bcd, 'ibm4690Packed').

As noted in errata 2.5, a complex element can have 'endOfParent'. If so then its last child element can be any lengthKind including 'endOfParent'.

An element with 'endOfParent' must be the last thing in its 'box'. A 'box' is defined as a portion of the data stream that has an established length prior to the parsing of its children. Specifically a box is either:

- A complex element with lengthKind 'explicit', 'prefixed' or 'pattern' AND no (sequence right framing or sequence postfix separator or choice right framing)
- A choice with choiceLengthKind 'explicit'

When unparsing, if the parent is a complex element with lengthKind 'explicit' or a choice with choiceLengthKind 'explicit' then the element with lengthKind 'endOfParent' is padded or filled in the usual manner to the required length.

- **2.73**. Section 12.3.6. An element with lengthKind 'endOfParent' is allowed to be the root element of a parse or unparse.
- **2.74**. Sections 13.2.1, 22.2.1. During unparsing, the application of escape scheme processing should take place before the application of the emptyValueDelimiterPolicy property.
- **2.75**. Section 4.1.1. Replace the existing description of the Document Information Item's [schema] member with 'This member is reserved for future use'.

- **2.76**. Section 12.3.4. When property prefixIncludesPrefixLength is 'yes' there are some restrictions that need to be added to enable reliable lengths to be calculated:
 - If the prefix type is lengthKind 'implicit' or 'explicit' then the lengthUnits properties of both the prefix type and the element must be the same.
 - If the prefixType is lengthKind 'pattern' then the lengthUnits of the element must be 'characters'.
- **2.77.** Sections 12.3.4, 12.3.2. The sections for lengthKind 'prefixed' and 'delimited' need the equivalent of Table 16 to express their rules for binary data.
- **2.78**. Section 12.3.4. Add a note to cover the scenario where lengthUnits is 'bits' and lengthKind is 'prefixed'. When parsing, any number of bits can be precisely extracted from the data stream, but when unparsing the number of bits written will always be a multiple of 8 as the Infoset does not contain bit-level information.
- **2.79**. Section 13.6.1.1. Clarify text number pattern rules for use of V and P symbols in conjunction with # symbol.
 - A pattern with a V symbol must not have # symbols to the right of the V symbol.
 - $\circ\quad$ A pattern with P symbols at the left end must not have # symbols .
 - o A pattern with P symbols at the right end can have # symbols.
- **2.80**. Section 13.6.1.1. Clarify text number pattern rules for use of V and P symbols in conjunction with @ and E and * symbols.
 - A pattern with a V symbol must not have @ or * symbols.
 - A pattern with P symbols must not have @ or E or * symbols.

This means that a V symbol and an E symbol may occur in the same text number pattern. The BNF in Figure 5 is revised to allow this.

2.81. Section 15.2. The specification currently says "On unparsing the choice branch supplied in the infoset is output". This does not handle the case where one or more branches of a choice is a sequence or a choice (or a group ref to such). Here, the element in the Infoset is one of the children of the branch sequence but it might not be the first in the sequence, or the element in the Infoset is one of the children of the branch choice. To handle this scenario, 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 chosen. If the chosen branch causes a processing error, no other branches are chosen (that is, there is no backtracking).

To avoid any unintended behaviour, a branch sequence may be wrapped in an element.

- **2.82**. Section 12.3.5. The behaviour for unparsing when lengthKind is 'pattern' is the same as for 'delimited', ie, for a simple element use textPadKind to determine whether to pad, for a complex element the length is that of the ComplexContent region.
- 2.83. Section 23.3. Clarifications on what is returned by an expression.
 - Every property that accepts an expression must state exactly what the expression is expected to return
 - To ensure the returned value is of the correct type, use XPath constructors or the correct literal values
 - What is returned lexically by an expression follows XPath 2.0 rules, which this is not the same as xs:default and xs:fixed lexical content.
 - No extra auto-casting is performed over and above that provided by XPath 2.0. <u>XPath 2.0 has rules for when it promotes types and when it allows types to be substituted.</u>
 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 expression to return a DFDL string literal, the returned value is always treated as a DFDL string literal.
 - o Within an expression, a string is never interpreted as a DFDL string literal

- 2.84. Section 23.5.3. The dfdl:property() function is removed.
- **2.85.** Section 23.5.3. Two new functions are provided to assist in the creation of expressions that return DFDL string literals.

	Returns 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:stringLiteralFromString (\$arg)	Use this function 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.
dfdl:containsEntity (\$arg)	Returns a Boolean indicating whether the \$arg string argument contains one or more DFDL entities.

- **2.86.** Section 24. State that DFDL regular expressions do not interpret DFDL entities.
- **2.87**. Section 12.3.7. State that 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 fillByte property. (The fillByte is *not* used to trim an element of type xs:hexBinary when parsing.)
- **2.88.** Section 13.5. Add support for HP NonStop Tandem zoned decimals. In this architecture, the negative sign is incorporated in the last byte of the number in the usual manner, but the overpunching occurs on the highest bit (ie, value 8) of the nibble. Consequently, a new enum value 'asciiTandemModified' is added to property textZonedSignStyle.
- **2.89**. Section 12.1. In the description of the alignment property, remove the rule that states 'The alignment of a child component must be less than or equal to the alignment of the parent element, sequence or choice'. It is overly restrictive.
- **2.90**. Sections 12.3, 12.3.7.2. Additionally allow lengthUnits 'bits' to apply to binary signed integer types, to support the modeling of signed integer bit fields in the C language. The physical bits are interpreted as a two's complement integer.
- **2.91.** Section 12.3.4. State that the global simple type referenced by prefixLengthType only obtains values for missing properties from its own schema's default dfdl:format annotation. 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.
- **2.92.** Section 13.6. When property textNumberRep is 'zoned', the property description should state that 'zoned' is only allowed for SBCS encodings (schema definition error otherwise).
- 2.93. Sections 13.6, 13.7. State that when unparsing a number and excess precision is supplied in the Infoset and rounding is not in effect, it is a processing error. Applies to text numbers when rounding is not enabled (matches ICU behaviour), and to binary numbers (always no rounding).
- 2.94. Sections 6.3.1.3, 12.2. Correct the wording for NL mnemonic in Table 5 to make it clear that when parsing it means either %LF; or %CR; or %CR; %LF% or %NEL; or %LS; and not combinations of those. Similarly, state that outputNewLine can only be either %LF; or %CR; or %CR; %LF% or %NEL; or %LS; and not combinations of those.

- **2.95.** Section 12.1. State that that all delimiters are implicitly aligned on a byte boundary and this cannot be overridden using the alignment property. If part way through a byte then bits will be skipped. State that if representation is text or type is string, then alignment must be on a byte boundary (schema definition error otherwise).
- 2.96. Section 23.5.3. Three of the DFDL-specific functions are renamed: dfdl:position() -> dfdl:occursIndex() dfdl:count() -> dfdl:occursCount() dfdl:countWithDefault() -> dfdl:occursCountWithDefault()
- **2.97**. Section 12.3.2. Additionally allow lengthKind 'delimited' for elements of simple type xs:hexBinary.
- **2.98.** Section 13.7. State that the maximum allowed value for two's complement binary integers is implementation independent but must be at least 8 bytes.
- **2.99**. Section 13.7, 13.13. Add support for the IBM 4690 point of sale variant of a packed decimal. This has the following characteristics:
 - o Nibbles represent digits 0 9 in the usual BCD manner
 - o A positive value is simply indicated by digits
 - o A negative number is indicated by digits with the leftmost nibble being xD
 - If a positive or negative value packs to an odd number of nibbles, an extra xF nibble is added on the left

Existing properties binaryNumberRep and binaryCalendarRep each take a new enum 'ibm4690Packed'. For numbers, properties byteOrder and binaryDecimalVirtualPoint actively apply. For calendars, properties byteOrder calendarPatternKind and calendarPattern actively apply (same restrictions as for 'packed' and 'bcd'). Property 'binaryPackedSignCodes' does not apply. Property 'binaryNumberCheckPolicy' applies but currently has no effect.

Where the DFDL specification provides for general behaviours for 'packed' and 'bcd', those behaviours apply also to 'ibm4690Packed'. Specifically:

- The same lengthKind enums and rules apply.
- There is no rounding when unparsing, so a value that can't be accommodated is a processing error.
- If logical type is unsigned and a negative value is received, it is a processing error.
- If invalid bytes are parsed, it is a processing error.
- 2.100. Section 12.3.1. State that when unparsing an element with lengthKind 'explicit' and where length is an expression, then the data in the Infoset is treated as variable length and not fixed length. The behaviour is the same as lengthKind 'prefixed'.
- **2.101**. Section 23.4. The BNFL for DFDL expressions allows a variable to appear as a path segment. This is not supported by DFDL, which only allows variables to return a simple value, and XPath does not permit variables to return simple values in path segments.
- **2.102**. Section 23. State that it is a schema definition error if an array element appears as a segment in a path location and is not qualified by a predicate.
- **2.103**. Section 12.1. Clarify that when the alignment properties are applied to an array element, the properties are applied to each occurrence of the element (as implied by the grammar).
- **2.104.** Section 13.11.1. State that when parsing a calendar element with binaryCalendarRep 'packed', 'bcd' or 'ibm4690Packed' 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 calendarPattern according to the usual ICU rules.

- **2.105**. Section 9.1.1. State that the presence of a separator is not sufficient to cause the parser to assert that a component is known to exist.
- 2.106. Section 13.6. State that textStandardDecimalSeparator, textStandardGroupingSeparator, textStandardExponentRep, textStandardInfinityRep, textStandardNanRep and textStandardZeroRep must all be entirely distinct from one another, and it is a schema definition otherwise. This is in the interests of clarity, and is an extra constraint compared to ICU. If any property value is an expression, the checking of this constraint cannot take place until processing.
- 2.107. Section 11. The list of enums for the encoding property is extended to include 'US-ASCII-7-bit-packed' in order to support data formats where ASCII characters are encoded in 7 bits with no padding bit. The new enum is neither a CCSID or an IANA charset. The new enum is not in the set of encodings that a DFDL processor must accept in order to be minimally conformant.
- **2.108**. Section 3. Update the Glossary concerning annotations, as follows, and use the new or changed terms as appropriate throughout the specification:
 - Add: Annotation point A location within a DFDL schema where DFDL annotation elements are allowed to appear.
 - Add: Statement annotations The annotation elements dfdl:assert, dfdl:discriminator, dfdl:setVariable, and dfdl:newVariableInstance. Also called DFDL Statements.
 - Add: Defining annotations The annotation elements dfdl:defineFormat, dfdl:defineVariable, and dfdl:defineEscapeScheme
 - <u>Change:</u> Format annotations The annotation elements dfdl:format, dfdl:element, dfdl:simpleType, dfdl:group, dfdl:sequence, and dfdl:choice.
 - o Remove: Format Annotations
 - Change: 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.
 - Add: 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.
- **2.109**. Section 6.2. Clarify that at any single annotation point of the schema, there can be only one format annotation (as defined in 2.108).
- 2.110. Section 7.3.1, 7.4.1. When testKind is 'pattern' for an assert or discriminator:
 - The pattern is applied to the data position corresponding to the beginning of the representation. Consequently the framing (including any initiator) is visible to the pattern.
 - o It is a schema definition error if there is no value for encoding in scope.
 - It is a schema definition error if alignment is other than 1.
 - It is a schema definition error if leadingSkip is other than 0.
- **2.111**. Sections 5.2, 23.5.3. Correct the XML Schema facets and attributes that are used by the dfdl:checkConstraints() function. Specifically, the function does not use the default, minOccurs and maxOccurs attributes.
- **2.112**. Section 3. Update the Glossary concerning arrays, as follows, and use the new or changed terms as appropriate throughout the specification:
 - o Remove: Scalar Element
 - O Remove: Fixed-Occurrence Item
 - o Remove: Variable Occurrence Item
 - Optional Item

- Remove: Number Of Occurrences
- <u>Change: Required Element. An element declaration or reference where minOccurs is</u> greater than zero.
- Change: Optional Element. An element declaration or reference where minOccurs is equal to zero.
- Add: Fixed Array Element. An array element where minOccurs is equal to maxOccurs.
- Add: Variable Array Element. An array element where minOccurs is not equal to maxOccurs.
- Add: Occurrence. An instance of an element in the data, or an item in the DFDL Infoset.
- o Add: Count. The number of occurrences of an element. .
- o Add: Index. The position of an occurrence in a count, starting at 1.
- Add: Required Occurrence. An occurrence with an index less than or equal to minOccurs.
- Add: Optional Occurrence. An occurrence with an index greater than maxOccurs.
- **2.113.** Section 23. 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:
- 1. Wrap path locations in a call to fn:exactly-one() except when the path location occurs within certain functions which operate on arrays
- Check for the disallowed use of those XPath 2.0 functions that are not in the DFDL subset
- **2.114.** Section 23. 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.

3. Major Technical Fixes

The following major technical fixes have been identified.

3.1. Section 14.5. Changes to property hiddenGroupRef.

Change to behave like ref property. That is, it can not be placed in scope by a format annotation, and is only set at its point of use. Empty string is not an allowed value. This reflects that there is no hiddenGroupRef value that applies universally.

The spec is not clear as to whether this property is allowed on the sequences and choices that are the direct children of global groups, or on group references. Clarify that it is allowed on any xs:sequence or xs:choice but not on any xs:group, including group reference.

3.2. Section 17. Changes to properties inputValueCalc and outputValueCalc.

Change to behave like ref property. That is, they can not be placed in scope by a format annotation, and are only set at their point of use. Empty string is not an allowed value. This reflects that there is no inputValueCalc or outputValueCalc property value that applies universally

The spec is confused as to whether these properties are applicable to simple types. Remove any references to these properties in relation to simple types, as they are applicable to elements only. Any application to simple types is a future extension.

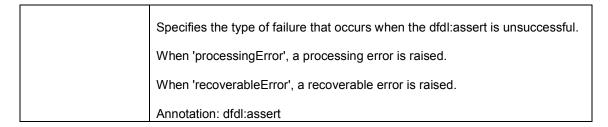
The spec is not clear as to whether these properties are allowed on global elements or element references. Clarify that they are allowed on local element and element references but not on global elements.

3.3. Section 12.3. Clarify that when property is lengthKind 'explicit', 'implicit' (simple only), 'prefixed' or 'pattern', it means that delimiter scanning is turned off and in-scope delimiters are not looked for within or between elements.

Consequently remove the last paragraph of section 5.2.2 starting "It is a processing error when a fixed-length string is found to have a number of characters not equal to the fixed number".

- **3.4.** Sections 2 and 7.3. Add a new failure type 'recoverable error' for use by the assert annotation when parsing, to permit the checking of physical constraints without terminating a parse. For example, using an assert to check a physical length constraint when property lengthKind is 'delimited'. Details:
 - After a recoverable error the parser will continue.
 - o Importantly, it does not cause backtracking to take place when speculating.
 - o It can be raised via a new enum attribute on dfdl:assert called 'failureType'.
 - o An error occurring during evaluation of a dfdl:assert remains a processing error.
 - All existing stated processing errors remain as such.
 - Discriminators remain unchanged.
 - The issuing of recoverable errors is independent of whether validation is enabled.

Property Name	Description
failureType	Enum (optional)
,,	
	Valid values are 'processingError', 'recoverableError'.
	Default value is 'processingError'.



Considered extending validation error to cover this, but the spec is quite clear that a validation error is a logical check performed on the infoset and the behaviour of the DFDL processor is unspecified.

- **3.5.** Section 13.8. The spec is not clear which variants of IEEE binary floats are supported. Clarify that support is for IEEE 754-1985, the same as XSDL 1.0. The implications of this are:
 - xs:float must have a physical length of 4 bytes for both 'ieee' and 'ibm390Hex' (schema definition error if explicit length is other than 4).
 - xs:double must have a physical length of 8 bytes for both 'ieee' and 'ibm390Hex' (schema definition error if explicit length is other than 8).
 - Add statement that there may be precision/rounding issues when converting IBM float/double to/from infoset float/double which is IEEE
 - Half-precision IEEE and quad-precision IEEE/IBM are not supported

Noted 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.

- **3.6.** Section 4. It was observed that the content of the DFDL infoset after parsing is not sufficient to build a W3C Post Schema Validation infoset (PSVI). Specifically, two things are missing:
 - o whether an element is valid
 - o for a simple element with a union type, which member the value matched.

In order to achieve this the DFDL infoset is modified as follows:

- Add a new Boolean [valid] member to element information item. A complex element information is not valid if any of its [children] are not valid. Empty if validation is not enabled.
- Add a new string [unionMemberSchema] member to simple element information item. This is 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 these properties are ignored. However, the augmented infoset which is built from the unparse operation should contain values for these properties if validation is enabled during unparsing.

3.7. Section 4, 9, 11, 12.3.7.1.3. Forcing a DFDL author to explicitly model a Unicode byte order mark (BOM) is a significant usability issue. Most authors working with Unicode data will expect a DFDL processor to handle BOMs in the same way as other software applications. Accordingly the DFDL specification is enhanced to add automatic detection and generation of Unicode BOMs.

A new string [unicodeByteOrderMark] member is added to the DFDL infoset document information item. 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 BOM. It there is a BOM 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.

The grammar production for the overall document changes to accommodate a BOM as shown in Appendix A of this document.

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 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 must be explicitly modelled as such, otherwise if parsed as part of an xs:string it will be interpreted as a zero-width non-breaking space (ZWNBS) character.

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. If a BOM does appear at any other point, then it must be explicitly modelled, otherwise if parsed as part of an xs:string it will be interpreted as a ZWNBS character.

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 unparser 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.

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 unparser will **not** look at the document information item [unicodeByteOrderMark] member and will **not** output a BOM. The byte order to use is implicit in the encoding. If a BOM does need to be output at the start of the data stream, then it must be explicitly modelled as such.

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

The unparser never outputs a BOM at any other point in the data stream. If a BOM needs to appear, then it must be explicitly modelled as such.

3.8. Section 2.2. Clarification is needed to schema definition error reporting criteria. The intent of the spec is that a DFDL processor only needs to report schema definition errors that directly affect its processing of the data. This is because the nature of DFDL's scoping rules mean that often it is not possible to validate an object definition for correctness in isolation.

Clarify that a DFDL processor:

- That only implements a DFDL parser does not have to validate 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 validate properties or annotations required by those optional features, but MUST issue a warning that an unrecognized property or annotation has been encountered.
- Need not validate global objects as they may legitimately be incomplete, with the following exceptions which must be validated:
 - 1. Global simple types that are referenced by prefixLengthType property
 - 2. Global elements that are the document root.

Clarify what action a DFDL processor should take when it encounters an object that explicitly carries properties that are not relevant to the object as defined.

- Property not applicable to the object's DFDL annotation.
 Schema definition error. Example is lengthKind on xs:sequence.
- Property not applicable because of simple type.
 Warning (optional). Example is calendarPatternKind on xs:string.
- Property not applicable because of another DFDL property setting.
 Warning (optional). Example is binaryNumberRep when representation is text.
- Property not applicable because object is local, global or reference.
 Warning (optional). Example is occursCountKind on a global xs:element.¹
- **3.9.** Section 12.3.5.1. The spec currently allows lengthKind 'pattern' to be used when the representation of the current element, or of a child element, is binary, but imposes restrictions on the encoding that can be in force. However encoding is not necessarily examined for binary elements, so this would introduce another reason for needing encoding.

Change the spec so that lengthKind 'pattern' is only applicable

- elements of simple type with representation 'text'
- elements of complex type

dfdl-wg@ogf.org

Excludes inputValueCalc and outputValueCalc..

For an element of complex type:

- all simple child elements must have representation 'text' and have the same encoding as the parent complex element, and
- 2. all complex child elements must themselves follow 1 and 2 (recursively).

Note that the same restrictions do **not** apply to testKind="pattern" on asserts and discriminators.

Table 16 can accordingly be deleted.

3.10. Sections 5.1, 13.15. Allow complex elements to be nillable. There are advantages in permitting complex elements to be nillable as well as simple elements. For example, it provides better interoperability with XML infosets. However, to avoid the concept of a complex element having a value, which is not possible in DFDL, the only permissible nil value is the empty string, represented by the DFDL %ES; entity.

If a complex element has xs:nillable set to 'true', it is a schema definition error if nilKind is not 'literalValue' or nilValue is not the single value '%ES;'.

Allowing complex elements to be nillable also solves another problem, that of preserving the position of optional complex elements in an array that contains explicit gaps. An infoset item with the special value nil is created for each such gap.

Property nilValueDelimiterPolicy is applicable.

The grammar changes to reflect this, as shown in Appendix A of this document.

3.11. Section 16. If the occurrences of an element in the data are not fixed (that is, the element is a variable array or is optional) and the count of the number of elements is not provided in the data nor is there a stop value, then the DFDL language only provides one mechanism for deducing the number of elements when parsing, namely occursCountKind 'parsed'. This causes the parser to speculate indefinitely until no more elements can be established. However there are circumstances where the minimum and maximum number of elements is known, and these facts could be used to guide the parse.

A new occursCountKind enumeration called 'implicit' is added, which takes into account minOccurs and maxOccurs settings.

Behaviour when parsing:

```
Expect up to maxOccurs in the data

Processing error if minOccurs not found or defaulted

Stop looking if > minOccurs found and a processing error is encountered

Stop looking if and when maxOccurs found (if not unbounded)
```

Behaviour when unparsing:

```
Expect up to maxOccurs in infoset
Processing error if minOccurs not in infoset or defaulted
Processing error if more then maxOccurs in infoset
```

To be added: A full behavioural description of all occursCountKind enumerations for clarification purposes.

3.12. Section 2.4. Validation checks are constraints expressed in XSDL, and they apply to the logical content of the infoset. Currently the spec says 'an unparse validation error occurs

when the physical representation being output would generate a validation error when parsing the data representation using the same DFDL schema.' This is a convenient definition, but problematic, because the original infoset used by the unparser could have been invalid, and the act of DFDL unparsing created a data stream which when parsed created a valid infoset. This can occur because of rounding, for example.

The specification will be changed to say that validation on parsing takes place on the infoset that is created by the parse, and that 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.

The new approach is in keeping with the way that XML Schema 1.0 defines validation against its PSVI.

3.13. Sections 4.1.2, 11. DFDL currently does not adequately describe how to handle decoding and encoding errors.

When parsing, these are the errors that can occur when decoding characters into Unicode/ISO 10646. The behavior in these cases is controlled by encodingErrorPolicy.

- The data is broken invalid byte sequences that don't match the definition of the encoding are encountered.
- Not enough bytes are found to make up the entire encoding of a character. That is, a fragment of a valid encoding is found.

The Unicode <u>replacement character</u> (U+FFFD) is substituted for the offending errors, one replacement character for any incorrect fragment of an encoding.

It is suggested that if a DFDL modeler 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.

The following specific situations involving UTF-16, UTF-16LE, and UTF-16BE do not cause an encoding error:

- UTF-16 and unpaired surrogate code-point
- UTF-16 and out-of-order surrogate code-point pair
- UTF-16 encoding utf16Width='fixed' and a surrogate code point is encountered
- UTF-16 byte-order-marks found not at the beginning of the data

In all these cases the code-point becomes part of the DFDL Information Item for the string.

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 2 bytes for a DBCS, but only 1 byte remains in the available length.

The behavior is determined by the encoding specification.

Each encoding has a replacement/substitution character specified by the ICU. These can be found conveniently in the ICU Converter Explorer (http://demo.icu-project.org/icu-bin/convexp). This character is substituted for the unmapped character or the character that has too large an encoding.

It is a processing error if it is not possible to output the replacement character because there is not enough room for its representation.

It is a processing error if a character encoding does not provide a substitution/replacement character definition and one is needed. (This would be rare, but could occur if a DFDL implementation allows many encodings beyond the minimum set.)

The current wording in section 4.1.2 description of [dataValue] is revised.

The above language describes what happens for values of simple elements that are of type string. This is the only time when replacement characters are used. For other data types, and for delimiters, failure to encode/decode is a processing error.

There are several details:

If lengthUnits is 'characters', then the replacement character counts as 1 character. Any character that is skipped counts as 0.

When lengthKind is 'delimited', then characters are decoded during the scan for delimiters <u>but</u> no replacement takes place. Any decoding error ends a match in that it causes a successful match to end prior to the error, or it causes a match that still has required characters to fail. When there are multiple possible delimiters, then some may fail due to decoding errors while others are still able to succeed. For example, a long delimiter may not match because it overshoots the end of a representation region in the data, and begins erroneously decoding binary data following the representation region, and this can result in a decoding error, which will be interpreted to mean that the error-causing decode cannot contribute to the match of that delimiter. That delimiter could still be a successful match that ends just prior to the decoding error location. If not, then, some other shorter delimiter that shares a prefix with the longer one may succeed to match, thereby indicating the proper end of the representation.

In the following, the match is against data where <u>any character replacement</u> has already taken place.

- when lengthKind is 'pattern', the scan for data matching the pattern
- when matching nilValue and the nilKind is literalCharacter or literalValue

When trimming padding characters, it is assumed that <u>any character replacement</u> has already taken place, so the trimming takes place on data that contains replacement characters.

Adding a new property to control behaviour after encode/decode error (eg, replace or skip or ignore) was considered in an earlier version of this Errata document but has been withdrawn.

3.14. Section 14.2. To better describe the property and its behaviour, property separatorPolicy is renamed to separatorSuppressionPolicy, and its enums renamed as follows:

```
'required' -> 'never'
'suppressed' -> 'anyEmpty'
'suppressedAtEndLax' -> 'trailingEmpty'
'suppressedAtEndStrict -> 'trailingEmptyStrict'.
```

The behaviour associated with each enum does not change, so 'anyEmpty' and 'trailingEmpty' behaviours are both lax.

3.15. Section 15. A new mechanism is introduced for resolving choices, the motivation being to make the cost of resolution close to constant time for choices with large numbers of branches where the branch to take is known in advance of parsing the choice.

A new element property is added called elementID of type 'DFDL String Literal'. This provides a single alternative identifier for the element. Allowed on local element, global element and element reference only.

A new dfdl:choice property is added called choiceBranchRef of type 'DFDL Expression'. The expression must evaluate to an xs:string. The resultant string must match (case insensitive) the elementID property value of one of the element branches of the choice, and if so discriminates in favour of that branch. The parser then goes straight to that branch, ignoring schema order.

Rules:

Because the branch is 'known to exist' no backtracking takes place if a processing error subsequently occurs.

Both properties behave like ref and hiddenGroupRef in that it is not possible to set a value in scope by a dfdl:format annotation, and is only set at its point of use. This is because there is nothing sensible that could be set in scope. Empty string is not an allowed value. It is allowable for an element reference to override the elementId of a global element.

Both properties are only used when parsing.

When choiceBranchRef is present, all choice branches must be local elements or element references. It is a schema definition error otherwise.

It is a processing error if the resolved value of choiceBranchRef does not match one of the branches.

It is a schema definition error if individual elementID values are not unique across all elements that are branches of a choice that carries choiceBranchRef, or if the elementID values of global elements are not unique within a given namespace.

It is a schema definition error if both initiatedContent and choiceBranchRef are provided on the same choice.

It is not a schema definition error if either initiatedContent or choiceBranchRef is provided on a choice and a discriminator exists on a choice branch.

Raw byte entities and character classes are not allowed for elementId.

- **3.16**. Section 14.2. Property documentFinalSeparatorCanBeMissing is removed as it is redundant. A postfix separator where the final separator can be missing can be modelled as an infix separator with documentFinalTerminatorCanBeMissing on the parent element.
- **3.17**. Section 21. The list of optional DFDL features is extended to make it easier for implementers to create minimal and extended conforming DFDL processors.

Feature	Detection
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"
BCD calendars	dfdl:binaryCalendarRep= <u>"bcd"</u> _
BCD numbers	dfdl:binaryNumberRep="bcd"
Multiple schemas	xs:include or xs:import in xsd
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 or attribute value
End of parent	dfdl:lengthKind = "endOfParent"

Existing optional feature 'Variables' is clarified to be dependent on optional feature 'Expressions'.

3.18. Section 9.2. The DFDL grammar productions are revised to make clear the distinction between the different allowable representations that an element can have and to enforce the correct use of the terms 'content', 'value' and 'representation'.

This has a significant effect on the grammar is shown in Appendix A of this document.

All sections of the specification are updated to ensure that 'content', 'value' and 'representation' are used correctly and consistently.

3.19. Sections 7.7. Additions and clarifications for the define Variable annotation.

A defaultValue expression must be evaluated before processing the data stream. It is a schema definition error otherwise.

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.

<u>If a defaultValue expression references another variable and this causes a circular reference,</u> it is a schema definition error

If the type of variable is a user-defined simple type restriction, it is a schema definition error

3.20. Sections 7.8. Additions and clarifications for the newVariableInstance annotation.

Only allowed as an annotation on sequence, choice or group reference. It is a schema definition error otherwise.

The resolved set of annotations for a component may contain multiple newVariableInstance statements. They must all be for unique variables, it is a schema definition error otherwise.

However, the order of execution among them is not specified. Schema authors can insert sequences to control the timing of evaluation of statements more precisely.

3.21. Sections 7.9. Additions and clarifications for the setVariable annotation.

Additionally allowed as an annotation on a simple type.

Not allowed as an annotation on a complex element or element reference to such.

The resolved set of annotations for a component may contain multiple setVariable statements. They must all be for unique variables, it is a schema definition error otherwise. However, the order of execution among them is not specified. Schema authors can insert sequences to control the timing of evaluation of statements more precisely.

<u>**3.22**</u>. *New appendix*. Add an explanation of the rationale behind the current variables design, covering why DFDL has adopted a write-once read-many behaviour for variables.

3.23. Sections 7.3.1. Additions and clarifications for the assert annotation.

Asserts can be placed as annotations on sequence, choice, group references, local and global element declarations, element references, and simple type definitions.

Replace "More than one dfdl:assert may be used at an annotation point. The dfdl:asserts will be evaluated in the order defined in the schema." with "If the resolved set of annotations for a schema component contain 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. Schema authors can insert sequences to control the timing of evaluation of statements more precisely."

3.24. Sections 7.3.1. Additions and clarifications for the discriminator annotation.

<u>Discriminators can be placed as annotations on sequence, choice, group references, local and global element declarations, element references, and simple type definitions.</u>

Replace "Any one annotation point can contain only a single dfd:discriminator or one or more dfdl:asserts, but not both. It is a schema definition error otherwise." with "The resolved set of annotations for a schema component can contain only a single dfd:discriminator or one or more dfdl:asserts, but not both. It is a schema definition error otherwise."

3.25. Section 9. Evaluation Order for Statement Annotations

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

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

For elements and element refs:

- dfdl:discriminator or dfdl:assert(s) with testKind='pattern' (parsing only)
- 2. dfdl:element following property scoping rules
- dfdl:setVariable(s)
- 4. dfdl:discriminator or dfdl:assert(s) with testKind='expression' (parsing only)

For sequences, choices and group refs:

- dfdl:discriminator or dfdl:assert(s) with testKind='pattern' (parsing only)
- 2. dfdl:newVariableInstance(s)
- 3. dfdl:setVariable(s)
- 4. dfdl:sequence or dfdl:choice or dfdl:group following property scoping rules
- 5. dfdl:discriminator or dfdl:assert(s) with testKind='expression' (parsing only)

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.

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.

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.

4. Appendix A. Grammar

This appendix provides a consolidated grammar that incorporates the several errata in this document.

Productions

Document = UnicodeByteOrderMark DocumentElement

DocumentElement = SimpleElement | ComplexElement

SimpleElement = SimpleLiteralNilElementRep | SimpleEmptyElementRep |
SimpleNormalRep

SimpleEnclosedElement = SimpleElement | AbsentElementRep

ComplexElement = ComplexLiteralNilElementRep | ComplexNormalRep |

ComplexEmptyElementRep

ComplexEnclosedElement = ComplexElement | AbsentElementRep

EnclosedElement = SimpleEnclosedElement | ComplexEnclosedElement

AbsentElementRep = Absent

SimpleEmptyElementRep = EmptyElementLeftFraming EmptyElementRightFraming ComplexEmptyElementRep = EmptyElementLeftFraming EmptyElementRightFraming

EmptyElementLeftFraming = LeadingAlignment EmptyElementInitiator PrefixLength
EmptyElementRightFraming = EmptyElementTerminator TrailingAlignment

SimpleLiteralNilElementRep = NilElementLeftFraming [NilLiteralCharacters |
NilElementLiteralContent] NilElementRightFraming

ComplexLiteralNilElementRep = NilElementLeftFraming NilElementRightFraming

NilElementLeftFraming = LeadingAlignment NilElementInitiator PrefixLength
NilElementRightFraming = NilElementTerminator TrailingAlignment

NilElementLiteralContent = LeftPadding NilLiteralValue RightPadOrFill

SimpleNormalRep = LeftFraming PrefixLength SimpleContent RightFraming
ComplexNormalRep = LeftFraming PrefixLength ComplexContent *ElementUnused*RightFraming

<u>LeftFraming = LeadingAlignment *Initiator*</u>
<u>RightFraming = *Terminator* TrailingAlignment</u>

<u>PrefixLength = SimpleContent | PrefixPrefixLength SimpleContent</u> PrefixPrefixLength = SimpleContent

SimpleContent = LeftPadding [NilLogicalValue | SimpleValue] RightPadOrFill

ComplexContent = Sequence | Choice

Sequence = LeftFraming SequenceContent RightFraming

SequenceContent = [PrefixSeparator EnclosedContent [Separator EnclosedContent]*

PostfixSeparator]

Choice = LeftFraming ChoiceContent RightFraming

ChoiceContent = [EnclosedContent] ChoiceUnused

EnclosedContent = [EnclosedElement | Array | Sequence | Choice]

Array = [EnclosedElement [Separator EnclosedElement]* [Separator StopValue]]
StopValue = SimpleElement

<u>LeadingAlignment = LeadingSkip AlignmentFill</u> TrailingAlignment = *TrailingSkip* RightPadOrFill = *RightPadding* | *RightFill*

5. Authors and Contributors

Stephen M. Hanson, IBM Software Group, Hursley, Winchester,UK smh@uk.ibm.com

Michael J. Beckerle, MA, USA mbeckerle.dfdl@gmail.com

We greatly acknowledge the contributions made to this document by the following people.

Tim Kimber, IBM Software Group, Hursley, UK Stephanie Fetzer, IBM Software Group, Charlotte, USA Richard Schofield, IBM Software Group, Hursley, UK Suman Kalia, IBM Software Group, Markham, Toronto, Canada

6. 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.

7. 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.

8. Full Copyright Notice

Copyright (C) Open Grid Forum (2011, 2012). All 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 on all such copies and derivative works. However, this document itself 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 Grid Recommendations in which case the procedures for copyrights defined in the OGF Document process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the OGF or its successors or assignees.

9. References

[DFDL] DFDL 1.0 http://www.ogf.org/documents/GFD.174.pdf/

C:\Data\Common

 $\underline{Transformation} \\ \underline{DFDL} \\ \underline{Specification} \\ \underline{V1.0} \\ \underline{http:\\www.ogf.org\\ \underline{documents} \\ \underline{GFD.174.p} \\ \underline{df} \\ \underline{}$

[GFD] OGF Document Process and Requirements

http://www.ogf.org/documents/GFD.174.pdf/

C:\Data\Common