# How Functional Programming leads to tight C-code at Runtime

**APACHE**
# Daffodil™

**Mike Beckerle,  Principal Engineer**
**Apache Daffodil PMC**
**mbeckerle at apache.org**

**OWL** Cyber Defense

**APACHECON**
**NORTH AMERICA**
**OCTOBER 3 - 6, 2022**
**NEW ORLEANS**
**LOUISIANA**
**WWW.APACHECON.COM**

1

# Outline/Agenda

- What is DFDL and What is Apache Daffodil?
- Daffodil's DFDL *schema compiler*
  - Functional programming
  - Object-Oriented Lazy Attribute Grammars (OOLAG).
- New runtime environment
  - C-code generator
    - C-language source code for standard C compilers
- Extras:
  - new Daffodil VSCode-based data debugger
  - EXI - dense binary "XML"

# What is DFDL
(Data Format Description Language)

## and

# What is Apache Daffodil?

# Got EDIFACT Data?

```
UNA:+.?*'
UNB+UNOC:4+5790000274017:14+5708601000836:14+990420:1137+17++INVOIC++++1'
UNH+30+INVOIC:D:03B:UN'
BGM+380+539602'
DTM+137:19990420:102'
RFF+CO:01671727'
NAD+BY+5708601000836::9'
RFF+VA:UK37499919'
NAD+SU++IBM UK'
RFF+VA:UK19430839'
RFF+ADE:00000767'
NAD+DP+++MyCompany+MyStreet+MyTown++1234+UK'
CUX+2:GBP:9'
LIN+1++V0370246:IN'
```

# Got bit-packed binary data ?

**Bytes are**

      **09 20 42 F0 0D B8 DD**

**Fields are decribed as:**

| | |
|---|---|
| **Message Number** | **XXXXXX00 00001xxx** |
| **FPI for Message Subtype** | **XXXXX0xx** |
| **FPI for File Name** | **XXXX0xxx** |
| **FPI for Message Size** | **XXX0XXXX** |
| **Operation Indicator** | **X01XXXXX** |
| **Retransmit Indictor** | **0XXXXXXX** |
| **Message Precendence Codes** | **XXXXX010** |
| **Security Classification** | **XXX00XXX** |

# Got NACHA Data?

```
101 121000248 1210002480608080107A094101WFB-W EDI CUST. DATA   WFB-E ACH SPINAT DATA

5200APD TX/FINCL SVC323413684               9666666606CCDAPD - TAX 0608020608022141021000024030649

62710700543200004001191477    0542151200614007046488KD8BRODY LABORATORIES INCFS002100024030840

82000000010010700543000542151200000000000000009666666606                        021000024030649

5200ARAR                            9000290001CCD PAYMENT         060807219102100027294149

62205100141200004945037059    00004036762394128       VIA LICENSING CORP      0021000027294283

820000000100051001410000000000000000004036769000290001                        021000027294149

5200ACME   WORLDWIDEDIRECT DEPOSIT       9954245682CCDA/P              060802214112200030000219

62210700543200004001191477    00000500001000142611008 100815047371712006       0122000036548030

62210700543200004001191477    00001475001000142611008 100815047375772006       0122000036548031

8200000002002140108600000000000000000001975009954245682                        122000030000219

5200BEST BANK NA   5046001042958       9560900031CCDEDI PAYMNT    0511181231021101100000014

6220510014149995275638771    0000037500504058967     XXXXXXXX BRANCH INC      1021101100001681

705RMR*OI*0140611**-1170.49*25*1170.36*

9999999999999999999999999999999999999999999999999999999999999999999999999999999999999999999
```

# Got ISO8583 Data?

1111FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF19123456789012345678
91234561234567890121234567890121234567890121231235959123
45678123456781234567812345699123112000099129912991231
12311231123412312312A1B2C3D4E5F6123123123412341991231
231234567890123456789012340B C0001500111234567890111234
56789012812345678901234567890123456782 6;11111111111111
11=1215=?1062;11222222222222222222222=1231123412341234121
1234561121212123 41?1A1B2C3D4E5F6A1B2C3123123A1A#A1A#A1#
A1#A1#A1#A1#15A1#A1#A1#A1#A1#15A1#A1#A1#A1#A1#35%A11111
111111111111^JOHNDOE^1215^?015A1#A1#A1#A1#A1#015A1#A1#A
1#A1#A1#015A1#A1#A1#A1#A1#ABCABCABC0800812345 6780080111
23011001100110011001100110011

# DFDL = Data Format Description Language

- A standard from Open Grid Forum (OGF)
- Started 2001, Ratified 2022
- Big - 200+ pages
- **DFDL → DaFfoDiL**

- DFDL is *mostly* not new ideas
  - *Standardizes existing practice* of data integration tools 1995 - 2010
- DFDL has some innovations
  - Especially for *unparsing binary data*

# Use Daffodil: NACHA as JSON Please...

```
{ "ACHFile":
{ "FileHeaderRecord":
{ "RecordTypeCode": "1",
"PriorityCode": "01",
"ImmediateDestination": " 123456789",
"ImmediateOrigin": " 987654321",
"FileCreationDate": "071030",
"FileCreationTime": "1634",
"FileIdModifier": "A",
"RecordSize": "094",
"ImmediateDestinationName": "TEST Destination ",
"ImmediateOriginName": "TEST Origination ", "ReferenceCode": " " },
"Batch": [ { …
```

# Use Daffodil: NACHA as XML Please...

```xml
<ACHFile xmlns="ach:2013">
<FileHeaderRecord>
<RecordTypeCode>1</RecordTypeCode>
 <PriorityCode>01</PriorityCode>
<ImmediateDestination> 123456789</ImmediateDestination>
<ImmediateOrigin> 987654321</ImmediateOrigin>
<FileCreationDate>071030</FileCreationDate>
<FileCreationTime>1634</FileCreationTime>
<FileIdModifier>A</FileIdModifier>
<RecordSize>094</RecordSize>
<ImmediateDestinationName>TEST Destination </ImmediateDestinationName>
<ImmediateOriginName>TEST Origination </ImmediateOriginName>
<ReferenceCode> </ReferenceCode>
</FileHeaderRecord> ...
```

# Introduction to

# Data Format Description Language

# aka DFDL

# Example – Delimited Text Data

<span style="color:red">rlimit=</span>5<span style="color:red">;rpngx=</span>-7.1E8

# Example – Delimited Text Data

Initiator

Separator

Initiator

$$rlimit=5;rpngx=-7.1E8$$

ASCII text integer

ASCII text floating point

- Separators, initiators (aka tags), & terminators are all examples in DFDL of *delimiters*

- Delimiters are one kind of *Framing.*

- DFDL divides the data into *content* (becomes values) and Framing (surrounds values)

# DFDL Schema

```
<complexType name="rPair">
  <sequence>
    <element name="rlim" type="xs:int"/>

    <element name="rpng" type="xs:float"/>

  </sequence>
</complexType>
```

Logical Elements

Simple Type

15

# DFDL Schema

Top level format declaration block applies to this entire schema *file*.

```
<annotation>
  <appinfo source="http://www.ogf.org/dfdl/">
    <dfdl:format representation="text"
             textNumberRep="standard" encoding="ascii"
             lengthKind="delimited" .../>
  </appinfo>
</annotation>

<complexType name="rPair">
  <sequence dfdl:separator=";">
    <element name="rLim" type="xs:int"
                 dfdl:initiator="rLimit=" />
    <element name="rpng" type="xs:float"
                 dfdl:initiator="rpngx=" />
  </sequence>
</complexType>
```

;

rLimit=5

rpngx=-7.1E8

DFDL properties

APACHE Daffodil

16

# DFDL Data and Infoset Lifecycle

**Data**

```
rLimit=5;rpngx=-7.1E8
```

**Parse**

**Unparse**

**DFDL Implementation**

**DFDL Schema**

**DFDL Implementation**

**Infoset**

*Element*
**Name:** rPair

*Element*
**Name:** rLimit
**Value:** 5
**Type:** Int

*Element*
**Name:** rpngx
**Value:** -7.1E8
**Type:** Double

**XML
or
EXI (binary XML)** *new!*
**or
JSON
or
Apache NiFi Records
or
Apache Drill ....
etc.**

# Internals of Apache Daffodil

# Apache Daffodil

- Daffodil contains
  - Full-blown compiler for DFDL schemas
  - JVM-based low-level runtime for parse/unparse
  - Big test suite (TDML)
  - New: C-code generating runtime for parse/unparse
- Written in Scala
  - Extensive use of Functional Programming

**Lines = 293K Total**

# Apache Daffodil

*If you download it, what do you get?*

- Jar libraries – runs on JVM
  - **DFDL Schema Compiler**, runtime, utilities, TDML runner
  - Signed Jars available from Maven Central
  - Java & Scala API with documentation
  - new: C-generator backend
    - (today: handles small subset of DFDL)

- Command Line Interface
  - Interactive CLI debugger and trace
  - XML, JSON, and EXI (new!) for parse-output, unparse-input

- New: Can also get the Daffodil VSCode Extension

# Any Compiler



Construct a different representation

Abstract Syntax Tree

Stay within a representation

Organized into passes

textual syntax

AST

Next

Next

Final

output

# Daffodil Schema Compiler



Construct a different representation, lazily

Each new representation is a pure function (lazy) of the prior representation

DFDL syntax

DSOM

GRAM

Back -end

output

DFDL Schema Object Model (DSOM)

Within each representation each member (aka 'attribute') is a pure function (lazy) of the existing 'tree' and attributes.

Ask for the output. Everything happens by Lazy Evaluation

# DSOM

- Similar to XML Schema Object Model (XSOM)
- Scala, lazy functional programming

# OOLAG - Object Oriented Lazy Attribute Grammars

- Functional Programming Idiom for Compilers
  - Johnsson, Thomas. (1995). *Attribute Grammars as a Functional Programming Paradigm*. LNCS. 274. 10.1007/3-540-18317-5_10.
- Attribute grammars - a grammar with 'attribute' computations
  - Wikipedia "Attribute Grammar"
  - synthetic (bottom up)
  - inherited (top down) (Not the OO notion of inherit)
- Object-Orientation
  - Mixins (via Scala traits), Inheritance
- Powerful pattern for Rich Transformations (like compilers)

# OOLAG - Object Oriented Lazy Attribute Grammars

- Lazy Evaluation - Avoids organizing compiler into 'passes'
- All values are special OOLAGValue that allow the answer of a computation to be
  - an ordinary value
  - a set of diagnostic objects, one or more of which are errors
  - both (a value and diagnostics that are only warnings)
- Code is structured into OOLAGValue calculations (using the LV idiom) and OOLAGHost objects
- OOLAGHost objects carry a list of required evaluations that must be evaluated to insure they are 'done' and can answer the isError test.

# Lazy Evaluation in Scala

**on class Choice (extends Term)**

```scala
lazy val hasKnownRequiredSyntax: Boolean = LV {
 hasFraming || groupMembers.forall { _.hasKnownRequiredSyntax }
}.value
```

**on class Term (a supertype of Choice)**

```scala
def hasKnownRequiredSyntax: Boolean

lazy val hasFraming = LV {
   hasInitiator || hasTerminator || !hasNoSkipRegions
}.value

lazy val hasNoSkipRegions = LV {
 leadingSkip == 0 && trailingSkip == 0
}.value
```

# Lazy Evaluation in Scala

```
on class Sequence (extends Term)

lazy val hasKnownRequiredSyntax: Boolean = LV {
 hasFraming ||
 groupMembers.exists { m =>
   m.isRequired &&
   m.hasKnownRequiredSyntax
 }
}.value
```

# OOLAG Host and OOLAG Value

- Error accumulation
  - Gathering more than one error before giving up
  - Avoiding duplicates
  - Associating file and line number with the right object to 'blame' from the DFDL schema
- Circular Definition Detection

# Daffodil Schema Compiler

# Gram - Grammar Trees

- Data Grammar
    - Based on concepts *Scala Combinator Parsers*
        - Yet Another Functional Programming Pattern


- Optimizations done on Grammar Trees
    - Back-end independent

# Grammar Rules in Scala

**on trait ModelGroupGrammarMixin**

```scala
lazy val termContentBody = prod {
  startGroupStmts ~ groupLeftFraming ~ groupContentWithDelims ~
  groupRightFraming ~ endGroupStmts
}

lazy val groupLeftFraming = prod {
  LeadingSkipRegion() ~ AlignmentFill()
}

lazy val groupContentWithDelims = prod {
  initiatorRegion ~ groupContent ~ terminatorRegion
}
```

**on class InitiatedTerminatedMixin**

```scala
lazy val terminatorRegion =
  prod(hasTerminator) { delimMTA ~ Terminator() }
```

# C-code Runtime
## aka "Runtime 2"

Mostly a contribution of John Interrante of GE Research

# Daffodil Schema Compilation

# Runtime 2: Different Goals

- Accepts Restrictions on DFDL for simplicity and performance
  - Deterministic - no backtracking
  - Must fit in memory - no streaming
  - Focus on binary data, mostly fixed-length fields
- Generates C code
  - Fast, small footprint
  - Statically allocate everything possible

- Future: generate VHDL for FPGA hardware realization

# Runtime 2 Infoset

- More efficient than what programmers would typically create

- Generality needed to handle all of DFDL infoset

- Walkable: metadata connected

- Cache/Prefetch friendly - localized/linearized (not pointers)
  - Similar concepts to Java 'Valhalla' JVM design goals

# Runtime 2 Infoset - Localized



base → meta

A → meta

B → meta

C → meta

A

B

C

malloc/heap oriented, pointer intensive

linear, cache-friendly, prefetch friendly, hardware friendly

# Runtime for Daffodil Runtime 2



Conveniently plugs into all our test infrastructure

Makes the infoset tangible

# C Back-end (aka Runtime 2)

- Two simple tuples

```xml
<complexType name="FooType">
  <sequence>
    <element name="a" type="xs:int"/>
    <element name="b" type="xs:int"/>
    <element name="c" type="xs:int"/>
  </sequence>
</complexType>

<complexType name="BarType">
  <sequence>
    <element name="x" type="xs:double"/>
    <element name="y" type="xs:double"/>
    <element name="z" type="xs:double"/>
  </sequence>
</complexType>
```

# C Back-end (aka Runtime 2)

- Tagged union of two tuples

```
<complexType name="NestedUnionType">
  <sequence>
    <element name="tag" type="xs:int32"/>
    <element name="data">
      <complexType>
        <choice dfdl:choiceDispatchKey="{xs:string(../tag)}">
          <element name="foo" type="idl:FooType"
              dfdl:choiceBranchKey="1 2"/>
          <element name="bar" type="idl:BarType"
              dfdl:choiceBranchKey="3 4"/>
        </choice>
      </complexType>
    </element>
  </sequence>
</complexType>
```
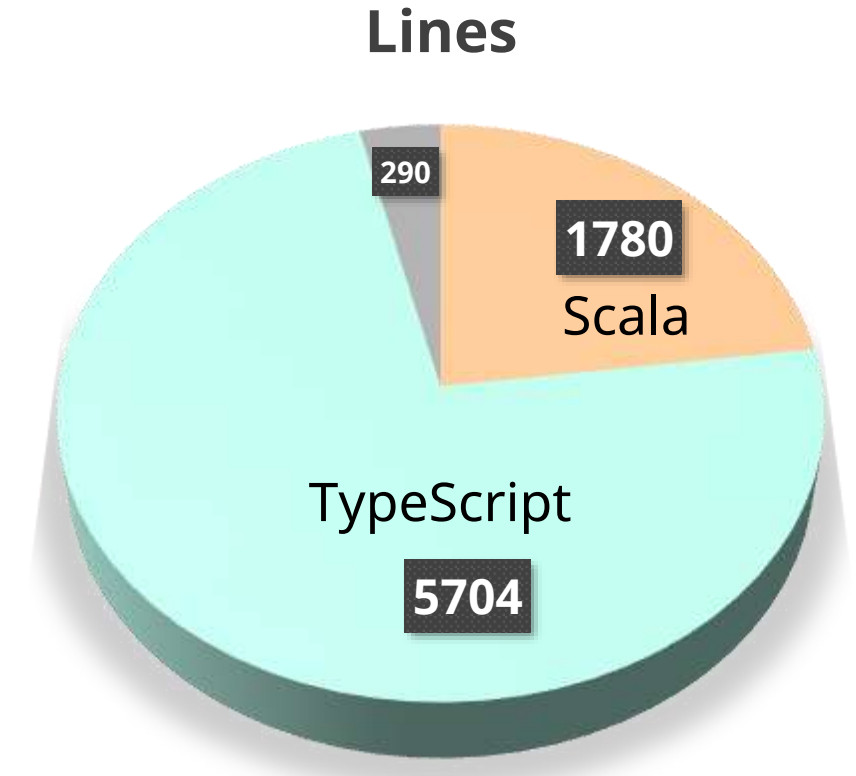
# Generated C Highlights

- the tuples

```c
typedef struct foo_data_NestedUnionType_
{
    InfosetBase _base;
    int32_t      a;
    int32_t      b;
    int32_t      c;
} foo_data_NestedUnionType_;

typedef struct bar_data_NestedUnionType_
{
    InfosetBase _base;
    double       x;
    double       y;
    double       z;
} bar_data_NestedUnionType_;
```

# Generated C Highlights

```c
typedef struct data_NestedUnionType_
{
    InfosetBase _base;
    size_t      _choice; // choice of which union field to use
    union
    {
        foo_data_NestedUnionType_ foo;
        bar_data_NestedUnionType_ bar;
    };
} data_NestedUnionType_;

typedef struct NestedUnion_
{
    InfosetBase _base;
    int32_t     tag;
    data_NestedUnionType_ data;
} NestedUnion_;
```

# Generated C Highlights

```c
static void
data_NestedUnionType__unparseSelf(const data_NestedUnionType_ *instance, UState *ustate)
{
    static Error error = {ERR_CHOICE_KEY, {0}};

    ustate->error = instance->_base.erd->initChoice(&instance->_base, rootElement());
    if (ustate->error) return;

    switch (instance->_choice)
    {
    case 0:
        foo_data_NestedUnionType__unparseSelf(&instance->foo, ustate);
        if (ustate->error) return;
        break;
    case 1:
        bar_data_NestedUnionType__unparseSelf(&instance->bar, ustate);
        if (ustate->error) return;
        break;
    default:
        // Should never happen because initChoice would return an error first
        ...
        return;
    }
}
```

```c
static void
foo_data_NestedUnionType__unparseSelf(
    const foo_data_NestedUnionType_ *instance,
    UState *ustate)
{
    unparse_be_int32(instance->a, 32, ustate); if (ustate->error) return;
    unparse_be_int32(instance->b, 32, ustate); if (ustate->error) return;
    unparse_be_int32(instance->c, 32, ustate); if (ustate->error) return;
}
```

# C-code Generator Status

- Still Partial
  - Needs strings, variable-length arrays, expressions

# More cool stuff…

# Apache Daffodil VSCode Debugger

- Data Format Debugger
  - Eventually a full Data-Format-Oriented IDE

- Extension to VSCode
  - Front-end - Typescript
    - Strongly typed Javascript
  - Back-end server - Scala
    - Uses the Daffodil library (Scala backend)
    - More functional programming idioms: typelevel FS2 & Cats Effect

**Lines**

290

1780

Scala

TypeScript

5704

# Apache Daffodil VSCode Extension

# Apache Daffodil VSCode Extension

# EXI - Dense Binary XML Alternative

- EXI = Efficient XML Interchange Format (W3C)
- Coming in Daffodil 3.4.0 (soon)
- Wrings all the redundancy and inefficiency out of XML text

Example: Aircraft messaging data format

- Original Message: 174 bytes
- Daffodil -> XML Text Infoset: 1493 bytes
- Daffodil -> EXI infoset: 160 bytes

# Conclusion/Review

- Quick Intro to DFDL and Apache Daffodil

- Daffodil Schema Compiler - Functional Programming & Scala
  - Useful idioms for DFDL compilation
  - Enables Code-generation for fast runtimes

- More Cool Stuff
  - VSCode Data Debugger/IDE
  - EXI dense binary alternative to XML text

- What am I working on this week?
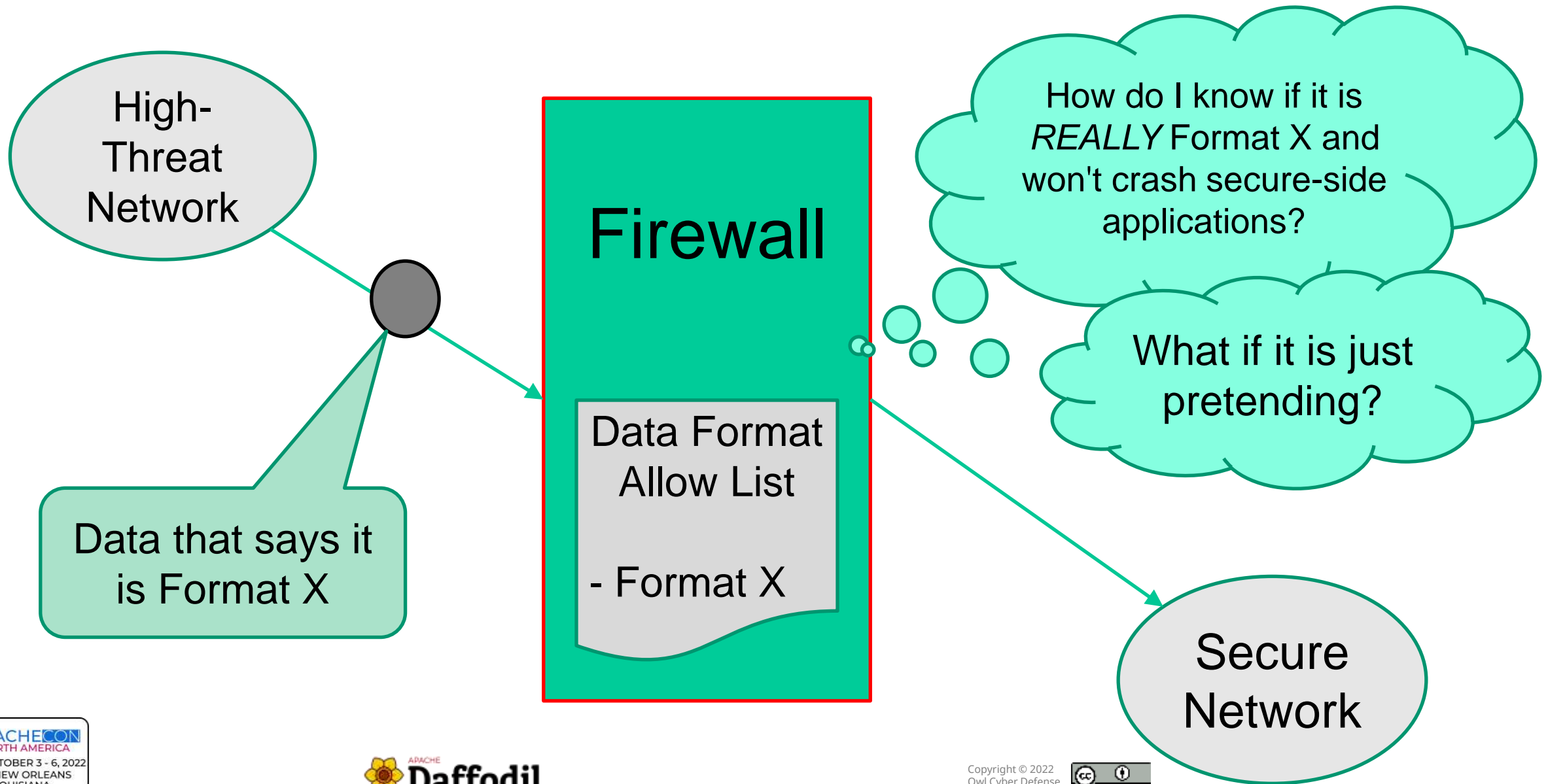  - Integration of Daffodil with Apache Drill !

# END

Notice:
This work is licensed under a
<u>Creative Commons Attribution 4.0
International License</u>

50

Use Case

# DFDL AND CYBER SECURITY

# Cyber-Security Use Case: Bad Data DoS Attack

# Cyber-Security Use Case: Full Protocol Break

**OWL** Cyber Defense

**Firewall**

DFDL Schema

Format X

High-Threat Net

**Daffodil Parse**

**Infoset**

Element
Name: rPair

Element
Name: rLimit
Value: 5
Type: Int

Element
Name: rpngx
Value: -7.1E8
Type: Double

**Daffodil Unparse**

Secure Net

Data says it is Format X

Infoset can be XML, JSON, EXI, or Daffodil's internal tree

Data is *proven* to be Format X, *by construction*

APACHECON
NORTH AMERICA
OCTOBER 3 - 6, 2022
NEW ORLEANS
LOUISIANA
WWW.APACHECON.COM

APACHE
**Daffodil**

Copyright © 2022
Owl Cyber Defense.

53