

An Introduction to XML::Twig

Growing Twigs

Content

- What is XML::Twig
 - A description of the module, why use it
- Working with XML::Twig
 - Resources, Installation, Example code
- Behind the scenes
 - Development process, why open-source

What is XML::Twig

- XML::Twig: XML, The Perl Way
 - a Perl module
 - to process XML
 - hybrid processing model, perlish API
- Alternatives
 - Perl Modules: XML::LibXML, XML::Simple, XML::SAX
 - XSLT
 - Java, Python, Ruby...

A Perl Module

- a Perl Module is a library that can be used from a perl program
- most perl modules (several 1000s) can be found on CPAN (http://cpan.org)
- like a lot of modules, XML::Twig is Object Oriented:

```
use XML::Twig;
my $twig= XML::Twig->new( @arguments);
```

Processing XML

- XML::Twig can parse XML and process it
- I use it to:
 - generate XML from IEEE Standards in FrameMaker
 - generate XHTML from IEEE Standards in XML
 - extract definitions from IEEE Standards in XML and store them into a data base
 - move data between databases on different OSs
 - power the templating system for my wife's website

XML Processing Models

- Stream Mode (SAX)
 - during parsing, call methods for each parsing event (open tag, text, close tag)
 - low memory usage, complex to use
- Tree Mode (DOM)
 - load the XML in memory, as a tree of objects
 - select nodes using navigation or queries (Xpath)
 - transform using delete, move, insert methods
 - needs more memory, easier to use

XML Tree Model



Growing Twigs

XML::Twig Processing Model

- Tree mode, using a simplified DOM
- Possibility to add handlers to elements
 - selected by element name, or complex condition,
 - called when the element is finished parsing
 - handler has access to the tree for the element
- Possibility to build the tree only for certain elements
 - other elements are ignored or output as-is

Other XML:: Twig Features

XML::Twig is designed to be practical

- whitespace handling
- comment/processing instructions handling
- encoding handling
- rich API (over 500 methods)

Other Perl Modules

• XML::LibXML

- based on libxml2 (http://xmlsoft.org)
- very powerful, fast, supports Xpath, DOM and lots of other W3C standards

• XML::Simple

- converts data-oriented XML to a Perl data structure

• XML::SAX

- event-driven, low-level
- lots of helper modules

XSLT

- W3C's language for processing XML
- Works well
- The code is in XML
- You don't get CPAN!

Using XML::Twig

- Installing XML::Twig
 - installing the pre-requisites: expat, XML::Parser
- Resources
 - finding information on how to use the module
- Examples
 - using XML::Twig with data-oriented XML
 - using XML::Twig with document-oriented XML

Installing XML::Twig

- Pre-requisites:
 - perl! (5.005 minimum, 5.8.3+ recommended)
 - expat: the low-level XML parsing library
 - XML:: Parser: the Perl wrapper for expat
 - optional Perl modules (XML::XPath, LWP, HTML::Entities)

Installing Perl Modules

the old-fashioned way

tar zxvf XML-Twig-3.22.tar.gz
perl Makefile.PL
make
make test
make install

• cpan / cpanplus

cpan XML::Twig

• distribution packages urpmi perl-XML-Twig

Resources

• The README file

install instructions, dependencies, links

• perldoc XML::Twig

reference doc

• http://xmltwig.com

docs, tutorial, FAQ, examples, development version

• http://perlmonks.org

The Most Important Slide

- Always, ALWAYS, check the data first:
 - parse the XML before doing anything with it
 - if you can, refuse the XML if it is not valid
 - if you cannot, write code to fix it, then validate it
- It doesn't matter who generated the XML, another company, another department, your department, YOU
- ONLY work on clean data
- You WILL hate character encodings

. . .

Examples Introduction

- 2 Examples: data XML and document Data-oriented vs Document-oriented XML
- Text is messy, data is simpler!
- Data has more structure
- Data has no mixed-content
- Differences in usage
- Some tools work best (or only!) for dataoriented XML
- Most XML these days is data-oriented

Data-oriented XML

- Data Base dumps/extracts
- Standard Documents
- Serialized objects
- XML-RPC
- log files
- Configuration files

Time for a quick break!

Growing Twigs

Stop Using <XML> Everywhere Please!

XML is Everywhere

 Documents Configuration files Data Serialized objects

Configuration Files

- XML is turned into a Perl Data Structure
- Works reasonably

but... XML is UGLY!

Growing Twigs

XML Version

```
<config logdir="/var/log/foo/"
        debugfile="/tmp/foo.debug">
  <server name="sahara" osname="solaris"</pre>
          osversion="2.6">
     <address>10.0.101</address>
     <address>10.0.1.101</address>
   </server>
   <server name="gobi" osname="irix"</pre>
           osversion="6.5">
      <address>10.0.102</address>
      <address>10.0.103</address>
   </server>
</config>
```

YAML Version

```
Debugfile: '/tmp/foo.debug'
logdir: '/var/log/foo/'
server:
  gobi:
    address:
      -10.0.0.102
      -10.0.0103
    osname: irix
    osversion: 6.5
  sahara:
    address:
      -10.0.0101
      -10.0.1.101
    osname: solaris
    osversion: 2.6
```

Try this at home

perl -MYAML -MXML::Simple \ -e 'print Dump XMLin "conf.xml"'

XML for Data

Data lives in...

Data Bases!

Growing Twigs

Data Bases

•Fast Reliable •Multi-user Scalable!

XML for data

It's just like text files... ...only **Slower**!

Growing Twigs

Exporting XML

• XML::Generator::DBI

• XML::Handler::YAWriter

```
use DBI;
use XML::Generator::DBI;
use XML::Handler::YAWriter;
my $dbh= DBI->connect(...);
my $ya = XML::Handler::YAWriter->new(AsFile => "-");
my $generator = XML::Generator::DBI->new(
Handler => $ya, dbh => $dbh
);
$generator->execute('SELECT * FROM data');
```

Conclusion

- Use XML when it makes sense
- Don't use it just because it's a buzzword

THINK!

Growing Twigs

Typical Use of Data-oriented XML

XML is an EXCHANGE format

- Extract data
 - Put it in a Data Base
- Fix the data
- Add data
- Avoid XML transformations!
 - SQL vs XPath

Example 1: XML Catalog

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

```
<plant id="id_001">
```

<common>Bloodroot</common>

<botanical>Sanguinaria canadensis/botanical>

<zone>4</zone>

```
<light>Mostly Shady</light>
```

```
<price>$2.44</price>
```

```
<availability>2005-03-05</availability>
```

```
</plant>
```

Example 1: Store the XML in a DB

- Store records from the catalog in a table in a database
- The catalog file can be very large, so do not load it in memory all at once

Example 1: Code

```
#!/usr/bin/perl
use strict;
use warnings;
use DBI;
use XML::Twiq;
my $CATALOG FILE = "plant catalog.xml";
my $DB FILE = "plant catalog.db";
my $dbh= DBI->connect("dbi:SQLite:dbname=$DB FILE","","");
my $sth= $dbh->prepare( "INSERT into plant
              (id, common, botanical, zone, light, price, availability)
        VALUES(?, ?, ?, ?, ?, ?, ?)");
XML::Twig->new( twig handlers => { plant => \&store plant })
        ->parsefile( $CATALOG FILE);
sub store plant
  { my( $t, $plant) = @ ;
    # the map... returns the list of text of all children of plant
    $sth->execute( $plant->id, map { $ ->text } $plant->children);
    $t->purge;
```

Example 2: Convert Currency

Convert the prices in dollars to prices in euros, add the currency as an attribute:

<price>\$3.02</price>

becomes

<price currency="EUR">2.58</price>

The code will be a filter that will only update the price elements and leave everything else untouched.

Example 2: code

```
#!/usr/bin/perl
use strict;
use warnings;
use XML::Twig;
my $CATALOG FILE = "plant catalog.xml";
# a rather silly example of extracting information from a web page
# nparse creates a new twig and parses the argument (an html file here)
my $RATE = XML::Twig->nparse( "http://www.x-rates.com/index.html")
                    ->first elt( 'a[@href="/d/USD/EUR/graph120.html"]')
                    ->text;
warn "rate: 1 EUR = SRATE USD\n";
my $catalog= XML::Twig->new( twig_roots => { price => \&price, },
                              twig print outside roots => 1,
                            );
$catalog->parsefile( $CATALOG FILE);
exit;
```

Example 2: code (cont.)

```
sub price
  { my( $twig, $price) = @ ;
   my $value= $price->text;
    if( svalue = ~ /^ (s(.*)s/)
      { my $dollar value= $1;
        my $euro value= sprintf( "%5.2f", $dollar_value / $RATE);
        $price->set text( $euro value);
        $price->set att( currency => "EUR");
    else
      { die "wrong dollar value '$value'\n"; }
    $price->print;
  }
```

Example 3: Update the data

- Update the catalog file with data from an other file
- 2 input XML files:
 - catalog
 - updates
- Output: updated catalog file
- The update file can be loaded in memory, not the main catalog

Example 3: XML update

<updates>

```
<plant id="id_013">
   <price>$7.22</price>
</plant>
<plant id="id_033">
   <availability>2005-05-28</availability>
</plant>
<plant id="id_021">
   <price>$4.20</price>
   <availability>2005-05-08</availability>
```

</plant>

</updates>

Example 3: code

```
#!/usr/bin/perl
use strict;
use warnings;
use XML::Twiq;
my $CATALOG FILE = "plant catalog.xml";
my $UPDATE FILE = "updates.xml";
my $updates= XML::Twig->new->parsefile( $UPDATE FILE);
my $catalog= XML::Twig->new(  # element => subroutine
               twig handlers => { plant => \&plant, },
               pretty print => 'indented',
                            );
$catalog->parsefile( $CATALOG FILE);
$catalog->flush;
exit;
```

Example 3: code (cont.)

```
sub plant
  \{ my( \ stwig, \ splant) = @; \
   my $id= $plant->att( 'id');
   my $update= $updates->elt id( $id); # updates is global
    if( $update)
      { foreach my $updated ( $update->children)
          { my $field = $updated->tag;
            my $original = $plant->first child( $field);
            $original->replace with( $updated);
            warn "updating $id - $field: ", $original->text,
                                     " => ", $updated->text, "\n";
          }
    $twig->flush; # prints the XML so far, and frees the memory
```

Document-oriented XML

- Important in publishing
- Allows:
 - independence from vendors
 - re-purposing of documents or parts of documents
- Often includes embedded data (part numbers, bibliographical items...)
- Often used to generate HTML or PDF

Processing Document-oriented XML

- Need to be able to work at 4 levels:
 - document level: to grab cross-references, number clause titles... often in a separate pass,
 - complex element level: tables, lists with internal references, chapter,
 - simple element processing: change a tag into an other tag (often adding the initial element as a class attribute),
 - within text: generate links from URLs, or from text elements, parse element text or attribute values.

Document Example

```
<?xml version="1.0" encoding="utf-8"?>
 <plant id="id 001">
    <common>Bloodroot</common>
    <botanical>Sanguinaria canadensis</botanical>
    <zone>4</zone>
    <light>Mostly Shady</light>
    <price>$2.44</price>
    <available>2005-03-05</available>
    <desc>A perennial <i>native</i> with a solitary white
        flower with golden stamens around a solitary pistil on a smooth
        stalk. 5-10 inches tall, this early plant has a reddish-orange
        juice down to the root (hence the name). The large blue/grey
        to green basal leaf is palmately scalloped into 5-9 lobes. See
        http://www.main.nc.us/naturenotebook/plants/bloodroot.html and
       http://en.wikipedia.org/wiki/Bloodroot
    </desc>
  </plant>
```

HTML Generation Code

```
#!/usr/bin/perl
use strict;
use warnings;
use XML::Twiq;
use Regexp::Common 'URI';
my $PLANT FILE="plant.xml";
my $twig= XML::Twig->new(
  twig handlers => {
    # here handlers are declared as anonymous subs, $ is the element
    common => sub { $_->set_tag( 'h1') }, # $_ is the element
   botanical => sub { $_->set_tag_class( 'p') }, # set tag to 'p' and
    zone => sub { $_->set_tag_class( 'p'); # class to the tag
                      $ ->prefix( "Grows in zone ");
    light => sub { $ ->set tag class( 'p'); },
   price => sub { $ ->set tag class( 'p'); },
    available => sub { $ ->set tag class( 'span');
                      $ ->prefix( ", available ");
                      $ ->move( last child => $ ->prev sibling);
                    },
```

HTML Generation Code (cont.)

```
desc
            => sub { $ ->set tag class( 'p');
                       $ ->insert new elt( before => h2 => "Description");
                       $ ->subs text( gr/($RE{URI}{HTTP})/,
                                     '&elt( a =>{ href => $1 }, $1)'
                                 );
           => sub { $_->set_tag( 'body'); },
   plant
  pretty print => 'indented',
                        );
$twig->parsefile( $PLANT FILE);
# add the html "wrapping"
my $html= $twig->root->wrap in( 'html');
my $head= $html->insert new elt( first child => 'head');
my $name= $twig->first elt( 'h1')->text;
$head->insert new elt( first child => 'title', $name);
$twig->print;
```

Behind the Scenes

The history of XML::Twig

- Why did I write XML::Twig?
- Why is it Open-Source?
- Development Process
- ToDo list

Why did I write XML::Twig

- Timeline:
 - 1998-02-10: the XML recommendation is published
 - 1998-03-??: XML::Parser published on CPAN
 - 1998-10-??: XML::Twig development starts
 - 1998-10-21: XML::DOM on CPAN
 - 1999-10-04: XML::Twig 1.6 on CPAN
 - 2005-10-14: XML::Twig 3.22 on CPAN
- In 1998 there were no XML module that would do what I wanted, I had to write my own!

Why is XML::Twig Open-Source?

- Instead of having to find the bugs, people (sometimes!) find them for me
- A good way to give back to the Open Source community that gave me Linux, Apache, PostgreSQL, SQLite, vi, Firefox, OpenOffice... and Perl!
- It's fun!

Development Process

- It has evolved with time:
 - in 1998 there was no Test Driven Development

• Now:

- revision control (CVS)
- tests added for every bug and new feature (Devel::Cover used to check coverage)
- Still very much a Cathedral, not a Bazaar

ToDo List

•Write a proper Xpath parser

needs to be used both in streaming mode, to trigger handlers and in normal mode, on an element or document

•Add "multi-parsing"

start several parsers (in threads) and allow them to rendez-vous to perform actions on all of them example: merging sorted XML files



Questions?

Growing Twigs

Grazie