



(XML::Parser .. XML::SAX)

Journées Perl 2004



- x **Agenda**
- x **XML**
- x **DOM / SAX**
- x **Flux + MI + Handle**
- x **XML::Parser**
- x **XML::Parser, reloaded**
- x **XML::SAX**
- x **Conclusion**
- x **Annexes**



XML¹

- x **XML est un méta-langage à balises imbriquées**
 - x XML
 - x **EX**tended **M**arkup **L**anguage
 - x Version *simplifiée* et *normalisée* du SGML
 - x Méta-langage
 - x Défini - par exemple - au travers d'une DTD (**D**ocument **T**ype **D**efinition)
 - x Permet une *validation* du document
 - x Balises imbriquées
 - x <BALISE attribut='valeur'>contenu</BALISE>
 - x Le contenu peut être constitué par d'autres balises
 - x Ou des données utilisateur
 - x Et ainsi de suite de façon hiérarchique



XML²

- x **Un document XML doit être bien formé et valide**

- x Bien formé

- x Respecte la syntaxe de la recommandation W3C
- x Déclaration d'en-tête correcte (e.g., `<?xml version='1.0' standalone='no'?>`)
- x Une balise ouvrante correspond à une balise fermante
- x Pas de recouvrement entre balises
- x Valeurs d'attributs entre guillemets ou cotes
- x Les caractères spéciaux hors balise sont remplacés par leur équivalent &

- x Valide

- x On vérifie que le document est bien conforme à sa DTD
- x Le validateur peut être un outil externe (i.e., rxp)
- x Mais ce n'est pas exhaustif
- x Utiliser sinon “XML Schema”, plus puissant mais aussi plus complexe



DOM / SAX

- x **2 approches principales, radicalement différentes**
 - x DOM : représentation “spatiale” / “pull”
 - x **D**ocument **O**bjet **M**odel
 - x Avaler tout le document (crée une *arborescence*)
 - x Se promener sur l'arbre à l'aide d'une *API*
 - x Fournit un accès *aléatoire* à un MI “privé” représentant le document
 - x SAX : représentation “temporelle” / “push”
 - x **S**imple **A**PI for **X**ML
 - x Travailler sur un *flux*
 - x Gérer des *événements* à l'aide de “callbacks”
 - x Construit un MI “maison” à partir du document *sérialisé*
- x **Nous allons étudier plus particulièrement SAX**



Flux + MI + Handle¹

- x **Le document devient un flux d'évènements**
 - x Début de balise (+ définition d'attributs)
 - x Caractères
 - x Fin de balise
- x **Déclenchant des fonctions (“callbacks”)**
- x **Utilisées pour peupler le modèle d'information**
 - x En insérant des données *incidentes* à l'endroit *courant* dans le MI
 - x Référencé par une *handle* (la dernière empilée dans le contexte)
 - x Pile de contexte mise à jour au fur et à mesure des évènements
 - x Les attributs et caractères sont fournis par les parsers
 - x Sous forme de hash ou string passées en paramètres aux callbacks



Flux + MI + Handle²

x **XML**

```
<foo nom='fou'>
  <bar>fred</bar>
</foo>
```

x **Notes**

- x On perd l'ordre (HoH, !LoL)
- x On doit gérer les collisions
- x FILS à la racine dû à META
- x META = pile de handles
- x Hiérarchie ~ à l'identique

x **MI**

```
$id = {
  'FILS' => {
    'foo' => { # $h1
      'ATTR' => {
        'nom' => 'fou'
      },
      'FILS' => {
        'bar' => { # $h2
          'CHAR' => 'fred'
        }
      }
    }
  }
};
```



Flux + MI + Handle³

x Flux d'évènements

- x Start (foo, %a)
- x Char("")
- x Start (bar)
- x Char ("fred")
- x End (bar)
- x End (foo)

x Etats & actions

- x [].....[\$h1]
peupler \$h1->{ATTR}
- x [\$h1].....[\$h1]
peupler \$h1->{CHAR}
- x [\$h1].....[\$h1, \$h2]
peupler \$h2->{ATTR}
- x [\$h1, \$h2].....[\$h1, \$h2]
peupler \$h2->{CHAR}
- x [\$h1, \$h2].....[\$h1]
remonter
- x [\$h1]..... []
remonter



XML::Parser¹

- x **Référencement des callbacks**

- x `my $p = new XML::Parser (Handlers => {
 Start => \&HStart, Char => \&HChar, End => \&HEnd });`

- x **Paramètres des callbacks**

- x `sub HStart { my($expat, $element, %attribute) = @_ }
sub HChar { my($expat, $string) = @_ }
sub HEnd { my($expat) = @_ }`

- x **Mise à jour du MI et de la handle par effet de bord**

- x Depuis les callbacks (pas très propre, n'est-ce pas ?)
 - x La handle (dernier élément de la pile) pointe le noeud courant
 - x On descend sur “start” et on remonte sur “end”



XML::Parser²

```
x  #!/usr/bin/perl -w
    use strict;

x  use XML::Parser;
    my $parser = new XML::Parser(
        Handlers => {
            Start => \&HStart,
            Char  => \&HChar,
            End   => \&HEnd,
        }
    );

x  our $id;
    $id->{META} = [];
    $parser->parsefile( 'foo.xml' );
    delete $id->{META};

x  use Data::Dumper;
    $Data::Dumper::Indent    = 1;
    $Data::Dumper::Sortkeys = 1;
    print Data::Dumper->Dump( [ $id ], [ 'id' ] );
```

```
x  sub HStart {
    my( $expat, $element, %attribute ) = @_;
    our $id;

    my $h = handle($id)->{FILS}->{$element} = {};
    while ( my( $k, $v ) = each %attribute ) {
        $h->{ATTR}->{$k} = $v;
    }
    push( @{$id->{META}}, $h );
}

x  sub HChar {
    my( $expat, $string ) = @_;
    our $id;

    unless ( $string =~ /\s+$/ ) {
        handle($id)->{CHAR} = $string;
    }
}

x  sub HEnd {
    my( $expat, $element ) = @_;
    our $id;

    pop @{$id->{META}};
}

x  sub handle {
    my( $id ) = @_;

    @{$id->{META}} ? $id->{META}->[-1] : $id;
}
```



XML::Parser, reloaded¹

- x **Le MI doit être un des paramètres des callbacks**
 - x Afin de ne plus utiliser de référence “globale” comme MI
 - x Ou d'employer le MI de \$expat
- x **Utilisation d'une closure**
 - x Start => sub { Hstart (\$id, @_) }
 - Char => sub { HChar (\$id, @_) }
 - End => sub { HEnd (\$id, @_) }
- x **Mais ce n'est pas suffisant**
 - x L'API n'est pas normalisée, donc dépendante du back-end
 - x On ne peut pas facilement faire de pipelines
 - x Le back-end (<eXpat/>) ne permet pas la validation (à dessein)



XML::Parser, reloaded²

```
x #!/usr/bin/perl -w
use strict;

x my $id;

x use XML::Parser;
my $parser = new XML::Parser(
    Handlers => {
        Start => sub { HStart ( $id, @_ ) },
        Char  => sub { HChar  ( $id, @_ ) },
        End   => sub { HEnd   ( $id, @_ ) },
    }
);

x $id->{META} = [];
$parser->parsefile( 'foo.xml' );
delete $id->{META};

x use Data::Dumper;
>Data::Dumper::Indent    = 1;
>Data::Dumper::Sortkeys = 1;
print Data::Dumper->Dump( [ $id ], [ 'id' ] );
```

```
x sub HStart {
    my( $id, $expat, $element, %attribute ) = @_;

    my $h = handle($id)->{FILS}->{$element} = {};
    while ( my( $k, $v ) = each %attribute ) {
        $h->{ATTR}->{$k} = $v;
    }
    push( @{$id->{META}}, $h );
}

x sub HChar {
    my( $id, $expat, $string ) = @_;

    unless ( $string =~ /\s+$/ ) {
        handle($id)->{CHAR} = $string;
    }
}

x sub HEnd {
    my( $id, $expat, $element ) = @_;

    pop @{$id->{META}};
}

x sub handle {
    my( $id ) = @_;

    @{$id->{META}} ? $id->{META}->[-1] : $id;
}
```



XML::SAX¹

- x **Décorrélation du front-end et du back-end**
 - x Le front-end est constitué de méthodes d'une classe (e.g., start_element, character, end_element, etc.)
 - x Le back-end est au choix : PurePerl, <eXpat/>, libxml2, etc.
- x **Approche orientée objet**
 - x Le MI est maintenant celui de l'objet appelant la méthode
 - x Permet de définir des constructeurs et destructeurs adaptés
 - x Et bien d'autres méthodes privées
- x **Permet de construire des pipelines**
 - x `use XML::SAX::Machines`



XML::SAX²

```
x #!/usr/bin/perl -w
x use strict;
x
x use XML::SAX;
x use MySAXHandler;
x my $parser = XML::SAX::ParserFactory->parser(
x     Handler => my $id = MySAXHandler->new );
x
x open( my $fh, 'foo.xml' );
x $parser->parsefile( $fh );
x close $fh;
x
x use Data::Dumper;
x $Data::Dumper::Indent = 1;
x $Data::Dumper::Sortkeys = 1;
x print Data::Dumper->Dump( [ $id ], [ 'id' ] );
x
x package MySAXHandler;
x use base ( 'XML::SAX::Base' );
x
x sub new { bless {} }
x
x sub handle {
x     my( $self ) = @_;
x     @{$self->{META}} ?
x         $self->{META}->[-1] : $self;
x   }
x
x sub xml_decl {}
```

```
x sub start_document {
x     my( $self ) = @_; $self->{META} = [];
x }
x
x sub end_document {
x     my( $self ) = @_; delete $self->{META};
x }
x
x sub start_element {
x     my( $self, $hash ) = @_;
x     my $h = $self->handle->{FILS}->{$hash->{Name}} = {};
x     for my $v ( values %{$hash->{Attributes}} ) {
x         $h->{ATTR}->{$v->{Name}} = $v->{Value};
x   }
x     push( @{$self->{META}}, $h );
x }
x
x sub characters {
x     my( $self, $hash ) = @_;
x     unless ( $hash->{Data} =~ /\s+$/ ) {
x         $self->handle->{CHAR} = $hash->{Data};
x   }
x }
x
x sub end_element {
x     my( $self, $hash ) = @_;
x     pop @{$self->{META}};
x }
x
x 1;
```



Conclusion

- x **Parser du XML avec SAX, ce n'est pas compliqué !**
- x **Avantages**
 - x Normalisé et reconnu par de multiples langages
 - x Faible consommation de mémoire
 - x Un front-end unique mais plusieurs back-end
 - x Permet de filtrer les données à la volée
- x **Inconvénients**
 - x Pas d'accès aléatoire
 - x Pas de requêtes complexes “à la” XPath
 - x Pas d'accès à la DTD



Annexes¹

- x **Pourquoi un MI “maison” ?**
 - x Compatibilité avec le format de l'ère pré-XML
 - x Un nouveau parser est nécessaire
 - x L'API est en effet conservée
 - x Get / set, itérateurs, etc.
 - x Vérifications qualitatives et quantitatives du MI
 - x Translation dans d'autres formats propriétaires
 - x Qualification du nouveau parser uniquement
- x **Utiliser du DOM à la place ?**
 - x Réécriture de l'API à l'aide de celle de DOM
 - x Requalification de toute l'API



Annexes²

x Un exemple d'utilisation de XML::Writer

```
x use IO::File;
$fh = IO::File->new( '> foo2.xml' );
use XML::Writer;
my $writer = XML::Writer->new( OUTPUT => $fh, DATA_MODE => 1, DATA_INDENT => 2 );

x BuildNode( $writer, $id->{FILS} );

x $writer->end();
close $fh;

x sub BuildNode {
    my( $writer, $node ) = @_;

    for my $e ( keys %$node ) {
        my $eh = $node->{$e};

        $writer->startTag( $e, map { ( $_, $eh->{ATTR}->{$_} ) } keys %{$eh->{ATTR}} );

        if ( $eh->{CHAR} ) {
            $writer->characters( $eh->{CHAR} );
        }
        if ( $eh->{FILS} ) {
            BuildNode( $writer, $eh->{FILS} );
        }
        $writer->endTag( $e );
    }
}
```




Annexes³

x Liens utiles (non exhaustif)

- x XML.....<http://www.w3.org/XML/>
.....<http://www.xml.com/>
- x Perl & XML.....<http://xml.com/pub/q/perlxml>
- x SAX.....<http://www.saxproject.org/>
- x XML Schemas.....<http://www.w3schools.com/schema/>
- x XSLT / Xpath.....<http://www.w3.org/Style/XSL>
- x <eXpat/>.....<http://expat.sourceforge.net/>
- x Gnome libxml2.....<http://xmlsoft.org/>
- x RXP.....<http://www.cogsci.ed.ac.uk/~richard/rxp.html>
- x CPAN.....<http://search.cpan.org>