

SWI-Prolog SGML/XML parser

1 Introduction

Markup languages have recently regained popularity for two reasons. One is document exchange, which is largely based on HTML, an instance of SGML and the other is for data-exchange between programs, which is often based on XML, which can be considered simplified and rationalised version of SGML.

James Clark's SP parser is a flexible SGML and XML parser. Unfortunately it has some drawbacks. It is very big, not very fast, cannot work under event-driven input and is generally

```
[],  
[ element(head,  
    [],  
    [ element(title,  
        [],  
        [ ' Demo'  
    ])  
    ]),  
element(body,  
    [],  
    [ '\n\n',
```

```
load_html_file(File, Term) :-  
    dtd(html, DTD),  
    load_structure(File, Term,  
        [ dtd(DTD),
```

[entity(alpha), ' < ', entity(beta)]

This is a special case of entity(*Code*), intended to handle special symbols by their name rather than character code.

The *Options* list controls the conversion process. Currently defined options are:

dtd(?*DTD*)

Reference to a DTD object. If specified, the <!DOCTYPE . . . > declaration is ignored and the document is parsed and validated against the provided DTD. If provided as a variable, the implicitly created DTD is returned. See section [3.4](#).

dialect(+*Dialect*)

Specify the parsing dialect. Supported are sgml (default), xml and smlns. See section [3.2](#) for details on the differences.

file(+*Name*)

Sets the name of the file on which errors are reported. Sets the linenumber to 1.

line(+*Line*)

Sets the starting line-number for reporting errors.

max_errors(+*Max*)

Sets the maximum number of errors. If this number is reached, an exception of the format below is raised. The default is 50.

~~max~~ 10.909 21-28(eci a f l / F2)n 0 1 5.45mi t4(The)-361.807 0 maxCode0454 0 cm BT70


```
....,
absolute_file_name(dtd(Type),
                    [ extensions([dtd]),
                      access(read)
```


dtd(*?DTD*)

If speci ed with an initialised DTD, this DTD is used for parsing the document, regardless of the document prologue. If speci ed using as a variable, a reference to the created DTD is returned. This DTD may be created from the document prologue or build implicitly from the document's content.

free_sgml_parser(*+Parser*)

Destroy all resources related to the parser. This does not destroy the DTD if the parser was created using the dtd(*DTD*) option.

set_sgml_parser(*+Parser*, *+Option*)

Sets attributes to the parser. Currently de ned attributes:

le(*File*)

source(+*Stream*)

An input stream that is read. Either this option or the goal (*Goal*) option must be provided.

goal(+*Goal*)

Goal is a callable term. The predicate `sgml_parse/2` opens an output stream to the parser and invokes `call(Goal, Stream)`, where *Goal*


```
        [ document(Content),  
          parse(content)  
        ],  
    process_rdf_description(element(Tag, Attr, Content))).
```

4 Processing Indexed Files

In some applications which process small portions of large SGML/XML/RDF

