Extracting Information from (\mathbb{H})T_EX Source Files

Jean-Michel Hufflen

TUG 2022

24 July 2022

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion Presently Conclusion

Introduction

Description

Discussion

Presently

Conclusion

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion Presently Conclusion

< □ ▷ < 큔 ▷ < 토 ▷ < 토 ▷ = - 의 Q (~ 2/2/15



is a wonderful tool for typesetting texts \leftarrow well known;

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents

Introduction Description

Discussion

Presently

Conclusion

< □ > < □ > < □ > < 三 > < 三 > < 三 > 三 の Q (~ 3/3/15



is a wonderful tool for typesetting texts \Leftarrow well known; knows only its own formats \Leftarrow well known, too.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents

Introduction Description

Presently

Conclusion

◆□ → ◆□ → ◆ 三 → ◆ 三 → ◆ ○ へ ○
4/3/15



is a wonderful tool for typesetting texts \leftarrow well known; knows only its own formats \leftarrow well known, too. Some information belonging to (LA)T_EX source files may be usable for other purposes than typesetting, e.g., generating metadata for Web search engines. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Introduction Description Discussion

resently

Conclusior

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □



is a wonderful tool for typesetting texts \Leftarrow well known; knows only its own formats \Leftarrow well known, too. Some information belonging to (LA)TEX source files may be usable for other purposes than typesetting, e.g., generating metadata for Web search engines. (LA)TEX's commands can do such jobs, but this is *misuse* and complicates the writing of classes. TEX has not been designed for that, it is preferable to use modern programming languages, with more suitable structures. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Functions are first-class objects, as other data.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents

Introduction Description

Discussion

Presently

Conclusion



Functions are first-class objects, as other data. Functions can be results of a computation. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Introduction

Description

Discussion

Presently

Conclusior

< □ > < ⑦ > < ≧ > < ≧ > ≥ のへで 8/4/15

Functions are first-class objects, as other data. Functions can be results of a computation. So we can easily write *generators* of functions.



Jean-Michel Hufflen

Introduction Description Discussion Presently Conclusion

< □ > < ⑦ > < ≧ > < ≧ > ≥ のQ ペ 9/4/15

Functions are first-class objects, as other data. Functions can be results of a computation. So we can easily write *generators* of functions. Scheme \Leftarrow elegant, data and programs have the same format, as in any Lisp dialect. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents

Introduction Description Discussion Presently

Conclusior

Building a parsing function

```
(g-mk-tex-parsing-f directive ...)
```

All the *directives* are grouped, 'compiled' into a function ready to parse a source file.



Jean-Michel Hufflen

Contents Introduction Description

Discussion

Presently

Conclusion



Directives

(g-retain-command command-name arg-nb optional-arg? top-level? recursive? preamble? occ-nb-info function)

where:

command-name is the name of the command to be caught;

arg-nb is the argument number of this command;

optional-arg? is true if the first argument is optional, surrounded by square brackets, false otherwise;

- top-level? is true if this command is to be searched only at the top level, false otherwise;
- recursive? is used when \input commands are encountered: if it is true, corresponding files are searched recursively, otherwise such an \input command is just skipped;
 - preamble? stops searching after a preamble if it is bound to
 true; otherwise, goes on.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Other arguments

occ-nb-info may be bound to:

- 0 or the false value (#f): the command should not appear within the file, this is checked;
- a positive integer n: the first n occurrences of this command are processed, and following ones are ignored;
- the true value (#t): all the occurrences of this command are processed;
- function the Scheme function to call, it must have the same number of arguments than $\command-name$.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Directives (con'd)

is used when \command-name's arguments are expressed by means of a pattern, e.g., "#1\endcsname" for the \csname command. s is a string bound to such a pattern, the other arguments have the same meaning than g-retain-command's. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Directives (con'd)

strings in both cases.

is used when \command-name's arguments are expressed by means of a pattern, e.g., "#1\endcsname" for the \csname command. s is a string bound to such a pattern, the other arguments have the same meaning than g-retain-command's. The arguments of corresponding functions in Scheme are Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

g-mk-tex-parsing-f builds a function that applies to a filename and returns:

true in all other cases.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion

resentry

Conclusior

 g-mk-tex-parsing-f builds a function that applies to a filename and returns:

true in all other cases.

You have to update your own structures when a file is parsed. If an error occurs, they may be in an inconsistent state.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Destructuring

These functions can be used to destructure an optional argument—s, s_0 are strings—:



Jean-Michel Hufflen

Contents

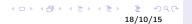
Introduction

Description

Discussion

Presently

Conclusion



Destructuring

These functions can be used to destructure an optional argument—s, s_0 are strings—:

(g-parse-to-alist s s₀) returns the successive associations of a comma-separated list whose elements are key=value pairs; if a key is given without a value, this missing value is replaced by s₀.

In both cases, the original order is preserved.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion Presently Conclusion

< □ > < ⑦ > < 注 > < 注 > 注 のへで 19/10/15

Destructuring

These functions can be used to destructure an optional argument—s, s_0 are strings—:

(g-parse-to-alist s s₀) returns the successive associations of a comma-separated list whose elements are key=value pairs; if a key is given without a value, this missing value is replaced by s₀.

In both cases, the original order is preserved. **Remark** Note that g-parse-to-list, g-parse-to-alist, g-retain-command and g-retain-match are functions, whereas g-mk-tex-parsing-f is a *macro*. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Example

- the used options of the babel package,
- the title,

the number of occurrences of the \emph command. (Show.) Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents

Introduction

Description

Discussion

Presently

Conclusion

< □ > < ♂ > < 差 > < 差 > 差 の < ⊙ < ⊙ < ○ 21/11/15

When I realised $MI{\rm BiBT}_{\!E\!}X$, I needed to know which options of the babel package were used. Only the preamble was to be searched.



Jean-Michel Hufflen

Contents Introduction Description Discussion Presently Conclusion

< □ > < ⑦ > < ≧ > < ≧ > ≧ シ うへで 22/12/15

When I realised $\text{MIB}{\rm IB}{\rm T}_{E}\!X$, I needed to know which options of the babel package were used. Only the preamble was to be searched.

The same to determine which encoding was used.



Jean-Michel

Hufflen



When I realised $\text{MIB}{\rm IB}{\rm T}_{E}\!X$, I needed to know which options of the babel package were used. Only the preamble was to be searched.

The same to determine which encoding was used. One year and a half ago, I became the new editor of the journal of the French-speaking T_{EX} user group.



Jean-Michel Hufflen



When I realised $\text{MIB}{\rm IB}{\rm T}_{E}\!X$, I needed to know which options of the babel package were used. Only the preamble was to be searched.

The same to determine which encoding was used. One year and a half ago, I became the new editor of the journal of the French-speaking T_{EX} user group. I decided to revise the class used for this journal and discovered that the previous class was used to build other files, such as metadata for Web search engines.



Jean-Michel Hufflen

Contents Introduction Description Discussion Presently Conclusion

When I realised $\text{MIB}{\rm IB}{\rm T}_{E}\!X$, I needed to know which options of the babel package were used. Only the preamble was to be searched.

The same to determine which encoding was used. One year and a half ago, I became the new editor of the journal of the French-speaking TEX user group. I decided to revise the class used for this journal and discovered that the previous class was used to build other files, such as metadata for Web search engines. From my point of view, TEX is intended for typesetting, and other tasks should be delegated to more modern programming languages. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

When I realised $\text{MIB}{\rm IB}{\rm T}_{E}\!X$, I needed to know which options of the babel package were used. Only the preamble was to be searched.

The same to determine which encoding was used. One year and a half ago, I became the new editor of the journal of the French-speaking TEX user group. I decided to revise the class used for this journal and discovered that the previous class was used to build other files, such as metadata for Web search engines. From my point of view, TEX is intended for typesetting, and other tasks should be delegated to more modern programming languages. I could have used Lua... Extracting Information from (IA)T_EX Source Files

Jean-Michel Hufflen

When I realised $\text{MIB}{\rm IB}{\rm T}_{E}\!X$, I needed to know which options of the babel package were used. Only the preamble was to be searched.

The same to determine which encoding was used. One year and a half ago, I became the new editor of the journal of the French-speaking TEX user group. I decided to revise the class used for this journal and discovered that the previous class was used to build other files, such as metadata for Web search engines. From my point of view, TEX is intended for typesetting, and other tasks should be delegated to more modern programming languages.

I could have used Lua... but LuaTEX was unable to process some texts designed for pdfTEX or X3LATEX.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Principles

Simple extractions—e.g., a title—must be simple.



Jean-Michel Hufflen

Contents Introduction Description Discussion

Presently

Conclusion



Principles

Simple extractions—e.g., a title—must be simple. If addressing occurrences of a T_EX command inside a source text, it should be possible to catch it.



Jean-Michel Hufflen

Contents Introduction Description Discussion Presently

Conclusior

Principles

Simple extractions—e.g., a title—must be simple. If addressing occurrences of a TEX command inside a source text, it should be possible to catch it. But more experience will be needed. Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion

Conclusior

< □ > < @ > < ≧ > < ≧ > 差 ● のへで 31/13/15

Current state

Some bugs to fix for the pattern-matching of TEX command's arguments.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description

Discussion

Presently

Conclusior

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Some bugs to fix for the pattern-matching of TEX command's arguments. The other points seems to be OK.



Jean-Michel Hufflen

Contents Introduction Description Discussion Presently

Conclusion



Some bugs to fix for the pattern-matching of TEX command's arguments. The other points seems to be OK. Available as a Scheme library.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion Presently

Conclusior

Ending

For many years, we have seen that in addition to T_EX 's works, many other tasks are more related to 'classical' programming. In particular, that is why Lua T_EX emerged.

Extracting Information from (IA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion Presently Conclusion

< □ > < @ > < ≧ > < ≧ > ≧ のQ () 35/15/15

Ending

For many years, we have seen that in addition to T_EX 's works, many other tasks are more related to 'classical' programming. In particular, that is why Lua T_EX emerged. We relate our work to avoiding *information redundancy*.

Extracting Information from (LA)T_EX Source Files

Jean-Michel Hufflen

Contents Introduction Description Discussion Presently Conclusion

< □ > < ⑦ > < ≧ > < ≧ > ≧ > うへで 36/15/15