

The L^AT_EX3 Programming Language— a proposed system for T_EX macro programming

© David Carlisle, Chris Rowley and
Frank Mittelbach
L^AT_EX3 project
`latex-l@urz.uni-heidelberg.de`

June 20, 2008

Abstract

This paper gives a brief introduction to a new set of programming conventions that have been designed to meet the requirements of implementing large scale T_EX macro programming projects such as L^AT_EX.

The main features of the system described are:

- classification of the macros (or, in L^AT_EX terminology, commands) into L^AT_EX functions and L^AT_EX parameters, and also into modules containing related commands;
- a systematic naming scheme based on these classifications;
- a simple mechanism for controlling the expansion of a function's arguments.

A system such as this is being used experimentally as the basis for T_EX programming within the L^AT_EX3 project. Note that the language is not intended for either document mark-up or style specification.

This paper is based on a talk given by David Carlisle in San Francisco, July 1997, but it describes the work of several people: principally Frank Mittelbach and Denys Duchier, together with Johannes Braams, David Carlisle, Michael Downes, Alan Jeffrey, Chris Rowley and Rainer Schöpf.

1 Introduction

This paper describes the conventions for a T_EX-based programming language which is intended to provide a more consistent and rational environment for the construction of large scale systems, such as L^AT_EX, using T_EX macros.

Variants of this language have been in use by The L^AT_EX3 Project Team since around 1990 but the syntax specification to be outlined here should *not* be considered final. This is an experimental language thus many aspects, such

as the syntax conventions and naming schemes, may (and probably will) change as more experience is gained with using the language in practice.

The next section shows where this language fits into a complete \TeX -based document processing system. We then describe the major features of the syntactic structure of command names, including the argument specification syntax used in function names.

The practical ideas behind this argument syntax will be explained, together with the semantics of the expansion control mechanism and the interface used to define variant forms of functions. The paper also discusses some advantages of the syntax for parameter names.

As we shall demonstrate, the use of a structured naming scheme and of variant forms for functions greatly improves the readability of the code and hence also its reliability. Moreover, experience has shown that the longer command names which result from the new syntax do not make the process of *writing* code significantly harder (especially when using a reasonably intelligent editor).

The final section gives some details of our plans to distribute parts of this system during the next year. More general information concerning the work of the \LaTeX 3 Project can be found in [4].

2 Languages and interfaces

It is possible to identify several distinct languages related to the various interfaces that are needed in a \TeX -based document processing system. This section looks at those we consider most important for the \LaTeX 3 system.

Document mark-up This comprises those commands (often called tags) that are to be embedded in the document (the `.tex` file).

It is generally accepted that such mark-up should be essentially *declarative*. It may be traditional \TeX -based mark-up such as \LaTeX 2 ϵ , as described in [3] and [2], or a mark-up language defined via SGML or XML.

One problem with more traditional \TeX coding conventions (as described in [1]) is that the names and syntax of \TeX 's primitive formatting commands are ingeniously designed to be 'natural' when used directly by the author as document mark-up or in macros. Ironically, the ubiquity (and widely recognised superiority) of logical mark-up has meant that such explicit formatting commands are almost never needed in documents or in author-defined macros. Thus they are used almost exclusively by \TeX programmers to define higher-level commands; and their idiosyncratic syntax is not at all popular with this community. Moreover, many of them have names that could be very useful as document mark-up tags were they not pre-empted as primitives (e.g., `\box` or `\special`).

Designer interface This relates a (human) typographic designer's specification for a document to a program that 'formats the document'. It should

ideally use a declarative language that facilitates expression of the relationship and spacing rules specified for the layout of the various document elements.

This language is not embedded in document text and it will be very different in form to the document mark-up language. For SGML-based systems the DSSSL language may come to play this role. For L^AT_EX, this level was almost completely missing from L^AT_EX2.09; L^AT_EX 2_ε made some improvements in this area but it is still the case that implementing a design specification in L^AT_EX requires far more ‘low-level’ coding than is acceptable.

Programmer interface This language is the implementation language within which the basic typesetting functionality is implemented, building upon the primitives of T_EX (or a successor program). It may also be used to implement the previous two languages ‘within’ T_EX, as in the current L^AT_EX system.

Only the last of these three interfaces is covered by this paper, which describes a system aimed at providing a suitable basis for coding large scale projects in T_EX (but this should not preclude its use for smaller projects). Its main distinguishing features are summarised here.

- A consistent naming scheme for all commands, including T_EX primitives.
- The classification of commands as L^AT_EX functions or L^AT_EX parameters, and also their division into modules according to their functionality.
- A simple mechanism for controlling argument expansion.
- Provision of a set of core L^AT_EX functions that is sufficient for handling programming constructs such as queues, sets, stacks, property lists.
- A T_EX programming environment in which, for example, all white space is ignored.

3 The naming scheme

The naming conventions for this programming language distinguish between *functions* and *parameters*. Functions can have arguments and they are executed. Parameters can be assigned values and they are used in arguments to functions; they are not directly executed but are manipulated by mutator and accessor functions. Functions and parameters with a related functionality (for example accessing counters, or manipulating token-lists, etc.) are collected together into a *module*.

Note that all these terms are only L^AT_EX terminology and are not, for example, intended to indicate that the commands have these properties when considered in the context of basic T_EX or in any more general programming context.

3.1 Examples

Before giving the details of the naming scheme, here are a few typical examples to indicate the flavour of the scheme; first some parameter names.

`\l_tmpa_box` is a local parameter (hence the `l_` prefix) corresponding to a box register.
`\g_tmpa_int` is a global parameter (hence the `g_` prefix) corresponding to an integer register (i.e., a `TEX` count register).
`\c_empty_toks` is the constant (`c_`) token register parameter that is for ever empty.

Now here is an example of a typical function name.

`\seq_push:Nn` is the function which puts the token list specified by its second argument onto the stack specified by its first argument. The different natures of the two arguments are indicated by the `:Nn` suffix. The first argument must be a single token which ‘names’ the stack parameter: such single-token arguments are denoted `N`. The second argument is a normal `TEX` ‘undelimited argument’, which may either be a single token or a balanced, brace-delimited token list (which we shall here call a *braced token list*): the `n` denotes such a ‘normal’ argument form.

`\seq_push:cn` would be similar to the above, but in this case the `c` means that the stack-name is specified in the first argument by a token list that expands, using `\csname...`, to a control sequence that is the *name* of the stack parameter.

The names of these two functions also indicate that they are in the module called `seq`.

3.2 Formal syntax of the conventions

We shall now look in more detail at the syntax of these names. The syntax of parameter names is as follows:

$$\backslash \langle access \rangle _{ \langle module \rangle } \langle description \rangle _{ \langle type \rangle }$$

The syntax of function names is as follows:

$$\backslash \langle module \rangle _{ \langle description \rangle } : \langle arg-spec \rangle$$

3.3 Modules and descriptions

The syntax of all names contains

$$\langle module \rangle \text{ and } \langle description \rangle :$$

these both give information about the command.

A *module* is a collection of closely related functions and parameters. Typical module names include `int` for integer parameters and related functions, `seq` for sequences and `box` for boxes.

Packages providing new programming functionality will add new modules as needed; the programmer can choose any unused name, consisting of letters only, for a module.

The *description* gives more detailed information about the function or parameter, and provides a unique name for it. It should consist of letters and, possibly, `_` characters.

3.4 Parameters: access and type

The *<access>* part of the name describes how the parameter can be accessed. Parameters are primarily classified as local, global or constant (there are further, more technical, classes). This *access* type appears as a code at the beginning of the name; the codes used include:

c constants (global parameters whose value should not be changed);

g parameters whose value should only be set globally;

l parameters whose value should only be set locally.

The *<type>* will normally (except when introducing a new data-type) be in the list of available *data-types*; these include the primitive `TeX` data-types, such as the various registers, but to these will be added data-types built within the `LaTeX` programming system.

Here are some typical data-type names:

int integer-valued count register;

toks token register;

box box register;

fint ‘Fake-integer’: (or fake-counter) a data type created to avoid problems with the limited number of available count registers in (standard) `TeX`;

seq ‘sequence’: a data-type used to implement lists (with access at both ends) and stacks;

plist property list

When the *<type>* and *<module>* are identical (as often happens in the more basic modules) the *<module>* part is often omitted for aesthetic reasons.

3.5 Functions: argument specifications

Function names end with an *<arg-spec>* after a colon. This gives an indication of the types of argument that a function takes, and provides a convenient method of naming similar functions that differ only in their argument forms (see the next section for examples).

The $\langle arg-spec \rangle$ consists of a (possibly empty) list of characters, each denoting one argument of the function. It is important to understand that ‘argument’ here refers to the effective argument of the \LaTeX function, not to an argument at the \TeX -level. Indeed, the top level \TeX macro that has this name typically has no arguments. This is an extension of the existing \LaTeX convention where one says that `\section` has an optional argument and a mandatory argument, whereas the \TeX macro `\section` actually has zero parameters at the \TeX level, it merely calls an internal \LaTeX command which in turn calls others that look ahead for star forms and optional arguments.

The list of possible argument specifiers includes the following.

- n** Unexpanded token or braced token list.
This is a standard \TeX undelimited macro argument.
- o** One-level-expanded token or braced token list.
This means that the argument is expanded one level, as by `\expandafter`, and the expansion is passed to the function as a braced token list. Note that if the original argument is a braced token list then only the first token in that list is expanded.
- x** Fully-expanded token or braced token list.
This means that the argument is expanded as in the replacement text of an `\edef`, and the expansion is passed to the function as a braced token list.
- c** Character string used as a command name.
The argument (a token or braced token list) must, when fully expanded, produce a sequence of characters which is then used to construct a command name (via `\csname`, `\endcsname`). This command name is the single token that is passed to the function as the argument.
- N** Single token (unlike **n**, the argument must *not* be surrounded by braces).
A typical example of a command taking an **N** argument is `\def`, in which the command being defined must be unbraced.
- O** One-level-expanded single token (unbraced).
As for **o**, the one-level expansion is passed (as a braced token list) to the function.
- X** Fully-expanded single token (unbraced).
As for **x**, the full expansion is passed (as a braced token list) to the function.
- C** Character string used as a command name then one-level expanded.
The form of the argument is exactly as for **c**, but the resulting token is then expanded one level (as for **O**), and the expansion is passed to the function as a braced token list.

- p** Primitive \TeX parameter specification.
This can be something simple like `#1#2#3`, but may use arbitrary delimited argument syntax such as: `#1,#2\q_stop#3`.
- T,F** These are special cases of **n** arguments, used for the true and false code in conditional commands.

There are two other specifiers with more general meanings:

- D** This means: **Do not use**. This special case is used for \TeX primitives and other commands that are provided for use only while bootstrapping the \LaTeX kernel. If the \TeX primitive needs to be used in other contexts it will be given an alternative, more appropriate, name with a useful argument specification. The argument syntax of these is often weird, in the sense described next.
- w** This means that the argument syntax is ‘weird’ in that it does not follow any standard rule. It is used for functions with arguments that take non standard forms: examples are \TeX -level delimited arguments and the boolean tests needed after certain primitive `\if...` commands.

4 Expansion control

4.1 Simpler means better

Anyone who programs in \TeX is frustratingly familiar with the problem of arranging that arguments to functions are suitably expanded before the function is called. To illustrate how expansion control can bring instant relief to this problem we shall consider two examples copied from `latex.ltx`.

```
\global
\expandafter
  \expandafter
\expandafter
  \let
\expandafter
  \reserved@a
\csname \curr@fontshape \endcsname
```

This first piece of code is in essence simply a global `\let`. However, the token to be defined is obtained by expanding `\reserved@a` one level; and, worse, the token to which it is to be let is obtained by fully expanding `\curr@fontshape` and then using the characters produced by that expansion to construct a command name. The result is a mess of interwoven `\expandafter` and `\csname` beloved of all \TeX programmers, and the code is essentially unreadable.

Using the conventions and functionality outlined here, the task would be achieved with code such as this:

```
\glet:0c \g_reserved_a_tlp
      \l_current_font_shape_tlp
```

The command `\glet:0c` is a global `\let` that expands its first argument once, and generates a command name out of its second argument, before making the definition. This produces code that is far more readable and more likely to be correct first time.

Here is the second example.

```
\expandafter
  \in@
\csname sym#3%
  \expandafter
  \endcsname
\expandafter
  {%
  \group@list}%
```

This piece of code is part of the definition of another function. It first produces two things: a token list, by expanding `\group@list` once; and a token whose name comes from ‘`sym#3`’. Then the function `\in@` is called and this tests if its first argument occurs in the token list of its second argument.

Again we can improve enormously on the code. First we shall rename the function `\in@` according to our conventions. A function such as this but taking two normal ‘`n`’ arguments might reasonably be named `\seq_test_in:nn`; thus the variant function we need will be defined with the appropriate argument types and its name will be `\seq_test_in:c0`. Now this code fragment will be simply:

```
\seq_test_in:c0 {sym#3} \l_group_seq
```

Note that, in addition to the lack of `\expandafter`, the space after the `}` will be silently ignored since all white space is ignored in this programming environment.

4.2 New functions from old

For many common functions the L^AT_EX3 kernel will provide variants with a range of argument forms, and similarly it is expected that extension packages providing new functions will make them available in all the commonly needed forms.

However, there will be occasions where it is necessary to construct a new such variant form; therefore the expansion module provides a straightforward mechanism for the creation of functions with any required argument type, starting from a function that takes ‘normal’ T_EX undelimited arguments.

To illustrate this let us suppose you have a ‘base function’ `\demo_cmd:nnn` that takes three normal arguments, and that you need to construct the variant `\demo_cmd:cnx`, for which the first argument is used to construct the *name* of

a command, whilst the third argument must be fully expanded before being passed to `\demo_cmd:nnn`. To produce the variant form from the base form, simply use this:

```
\exp_def_form:nnn {demo_cmd} {nnn} {cnx}
```

This defines the variant form so that you can then write, for example:

```
\demo_cmd:cnx {abc} {pq} {\rst \xyz }
```

rather than ... well, something like this!

```
\def \tempa {{pq}}%
\edef \tempb {\rst \xyz}%
\expandafter
  \demo@cmd
\csname abc%
  \expandafter
    \expandafter
      \expandafter
        \endcsname
  \expandafter
    \tempa
  \expandafter
    {%
  \tempb
  }%
```

Another example: you may wish to declare a function `\demo_cmd_b:xcxcx`, a variant of an existing function `\demo_cmd_b:nnnnn`, that fully expands arguments 1, 3 and 5, and produces commands to pass as arguments 2 and 4 using `\csname`. The definition you need is simply

```
\exp_def_form:nnn
  {demo_cmd_b} {nnnnn} {xcxcx}
```

This extension mechanism is written so that if the same new form of some existing command is implemented by two extension packages then the two definitions will be identical and thus no conflict will occur.

5 Parameter assignments and accessor functions

5.1 Checking assignments

One of the advantages of having a consistent scheme is that the system can provide more extensive error-checking and debugging facilities. For example, an accessor function that makes a *global* assignment of a value to a parameter can check that it is not passed the name of a *local* parameter as that argument: it does this by checking that the name starts with `\g_`.

Such checking is probably too slow for normal use, but the code can have hooks built in that allow a format to be made in which all functions perform this kind of check.

A typical section of the source¹ for such code might look like this (recall that all white space is ignored):

```
%<!*check>
\let_new:Nn
  \toks_gset:Nn \tex_global:D
%</!check>
%<*check>
\def_new:Npn
  \toks_gset:Nn #1
  {
    \chk_global:N #1
    \tex_global:D #1
  }
%</check>
```

In the above code the function `\toks_gset:Nn` takes a single token (`N`) specifying a token register, and globally sets it to the value passed in the second argument.

A typical use of it would be:

```
\toks_gset \g_xxx_toks {<some value>}
```

In the normal definition, `\toks_gset` can be simply `\let` to `\global` because the primitive `TEX` token register does not require any explicit assignment function: this is done by the `%<!*check>` code above.

The alternative definition first checks that the argument passed as `#1` is the name of a global parameter and raises an error if it is not. It does this by taking apart the command name passed as `#1` and checking that it starts `\g_`.

5.2 Consistency

The primitive `TEX` syntax for register assignments has a very minimal syntax and, apart from box functions, there are no explicit functions for assigning values to these registers.

This makes it impossible to implement alternative data-types with a syntax that is both consistent and at all similar to the syntax for the primitives; moreover, it encourages a coding style that is very error prone.

As in the `\toks_gset:Nn` example given above, all `LATEX` data-types are provided with explicit functions for assignment and for use, even when these have essentially empty definitions. This allows for better error-checking as described above; it also allows the construction of further data-types with a similar interface, even when the implementation of the associated functions is very complex.

¹This code uses the `docstrip` system described in [2], Section 14.3.

For example, the ‘fake-counter’ (`fint`) data-type mentioned above will appear at the \LaTeX programming level to be exactly like the data-type based on primitive count registers; however, internally it makes no use of count registers. Typical functions in this module are illustrated here.

```
\fint_new:N \l_tmpa_fint
```

This declares the local parameter `\l_example_fint` as a fake-counter.

```
\fint_add:Nn \l_example_fint \c_thirty_two
```

This increments the value of this fake-counter by 32.

6 The experimental distribution

The initial implementations of a \LaTeX programming language using this kind of syntax remain unreleased (and not completely functional); they partly pre-date \LaTeX 2 ϵ ! The planned distribution will provide a subset of the functionality of those implementations, in the form of packages to be used on top of \LaTeX 2 ϵ .

The intention is to allow experienced \TeX programmers to experiment with the system and to comment on the interface. This means that *the interface will change*. No part of this system, including the name of anything, should be relied upon as being available in a later release. Please do *experiment* with these packages, but do *not* use them for code that you expect to keep unchanged over a long period.

In view of the intended experimental use for this distribution we shall, in the first instance, produce only a few modules for use with \LaTeX 2 ϵ . These will set up the conventions and the basic functionality of, for example, the expansion mechanism; they will also implement some of the basic programming constructs, such as token-lists and sequences. They are intended only to give a flavour of the code: the full \LaTeX 3 kernel will provide a very rich set of programming constructs so that packages can efficiently share code, in contrast with the situation in the current \LaTeX where every large package must implement its own version of queues, stacks, etc., as necessary.

In the first release of this experimental system at least the following modules will be distributed.

l3names This sets up the basic naming scheme and renames all the \TeX primitives. If it is loaded with the option `[removeoldnames]` then the old primitive names such as `\box` become *undefined* and are thus available for user definition. Caution: use of this option will certainly break existing \TeX code!

l3basics This contains the basic definition modules used by the other packages.

l3chk A module that provides functionality comparable to \LaTeX ’s `\newcommand` and `\renewcommand`, and also the extra level of checking described above in section 5.1.

- l3tlp** This implements a basic data-type, called a *token-list pointer*, used for storing named token lists: these are essentially \TeX macros with no arguments.
- l3expan** This is the argument expansion module discussed above.
- l3quark** A ‘quark’ is a command that is defined to expand to itself! Therefore they must never be expanded as this will generate infinite recursion; they do however have many uses, e.g., as special markers and delimiters within code.
- l3seq** This implements data-types such as queues and stacks.
- l3prop** This implements the data-type for ‘property lists’ that are used, in particular, for storing key/value pairs.
- l3int** This implements the integer and ‘fake integer’ data-types.
- l3toks** A data-type corresponding to \TeX ’s primitive token registers.
- l3io** A module providing low level input and output functions.
- l3precom** A ‘pre-compilation’ module that provides functions dealing with pointer creation and handling, and using external files to record the state of the current definitions.

This distribution will also contain the \LaTeX source for the latest version of this document, a docstrip install file and three small test files.

In later releases we plan to add further modules and a full-fledged example of the use of the new language: a proto-type implementation for the ideas described in the article ‘Language Information in Structured Documents: A Model for Mark-up and Rendering’ [5].

References

- [1] Donald E Knuth *The \TeX book*. Addison-Wesley, Reading, Massachusetts, 1984.
- [2] Goossens, Mittelbach and Samarin. *The \LaTeX Companion*. Addison-Wesley, Reading, Massachusetts, 1994.
- [3] Leslie Lamport. *\LaTeX : A Document Preparation System*. Addison-Wesley, Reading, Massachusetts, second edition, 1994.
- [4] Frank Mittelbach and Chris Rowley. The \LaTeX 3 Project. *TUGboat*, ????? ??? 1998.
- [5] Frank Mittelbach and Chris Rowley. Language Information in Structured Documents: A Model for Mark-up and Rendering. *TUGboat*, ????? ??? 1998.