
HINT: Reflowing T_EX output

Martin Ruckert

Introduction

Current implementations of T_EX produce `.pdf` (portable document format) or `.dvi` (device independent) files. These formats are designed for printing output on physical paper where the paper size and perhaps even the output resolution is known in advance. If these conditions are met, T_EX, in spite of its age, still produces results of unsurpassed quality.

Due to improvements in display size, resolution, and technology over the past decades, it has become common practice to read T_EX output on screen not only before printing but also instead of printing. For viewing T_EX output before printing, excellent programs [4, 5] for “pre-viewing” are available. The prefix “pre” indicates that these programs intend to provide the user with a view that matches, as close as possible, the “final” appearance on paper. If, however, there is no intention of printing, for example if we read during a train ride on a mobile device, then matching the appearance on paper is of no importance, and we would rather prefer that the T_EX output instead adapt to the size and resolution of our mobile device. Anyone who has been forced to read a PDF file designed for output on letter paper on a 5" smartphone screen knows the problem.

For this reason, web browsers or ebooks use a reflowable text format. The HTML format, however, was never designed as a format for book printing, and `epub`, the ebook file format based on it, has inherited its deficiencies. Microsoft’s PDF reflow solution — converting PDF files to Word documents — is an indication of the need for reflowable file formats but is a proprietary surrogate at best.

Considering that the T_EX engine is able to reflow whole documents just by assigning new values to `hsize` and `vsize`, it seems long overdue to put this engine to use for that purpose.

The HINT project does just that. It defines a file format and provides two utilities: `HiTEX`, a special version of T_EX to produce such files, and `HINT`, a standalone viewer to display them.

What is HINT?

Adopting the usual free software naming convention, HINT is a recursive acronym for “HINT is not T_EX”. But then, what is it? One answer could be: It’s 90% T_EX and the rest is a mixture of good and bad luck. So let me start explaining the details.

A first overview can be obtained by looking at Figures 1–3. The first figure is a simplified depiction of T_EX’s structure: A complex input processing part translates T_EX input files into lists of 16-bit integers, called tokens, which form the machine language of T_EX. The main loop of T_EX is an interpreter that executes these programs, which eventually produce lots of boxes — most of them character boxes — and glue (and a few other items) that end up on the so-called contribution list. Every now and then, the page builder will inspect the contribution list and moves items to the current page. As soon as it is satisfied with the current page, it will invoke the (user-defined) output routine, again a token list, which can inspect the proposed page, change it at will, add insertions like footnotes, floating images, page headers and footers, even store it for later use, and eventually “ship out” the page to a `.dvi` file.

HINT splits this whole machinery into two separate parts: frontend and backend. The backend is the HINT viewer. The design goal is to reduce the processing in the backend as much as possible because we expect the viewer to run on small mobile devices where reduced processing implies reduced energy consumption and thus longer battery life. The frontend is the `HiTEX` version of T_EX which is prevented from doing the full job of T_EX because it does not know the values of `hsize` and `vsize`. As a first approximation of this split, `HiTEX` can write the contribution list to a file and HINT can read this file and feed it to the page builder as shown in Figures 2 and 3.

As an overall design goal, the `HiTEX` and HINT combination should produce exactly the same rendering as T_EX for a given `hsize` and `vsize`.

A closer look at Figure 3 reveals that the “Output” arrow, representing the user’s output routine, has disappeared; instead a new arrow, labeled “Templates”, has taken its place. Keeping the full power of T_EX’s output routines would imply keeping the full T_EX interpreter, all the token lists generated from the T_EX input file, and possibly even the files that such an output routine might read or write in the viewer. This seemed to be too high a price and therefore output routines have been replaced by the template mechanism described below. This was the single most important design decision guided by the desire to allow lightweight viewers to run efficiently with a minimum amount of resources.

Several iterations were necessary to arrive at a suitable file format that was compact, easy to digest, and sufficiently expressive to provide the necessary information to the viewer. Finally, many smaller

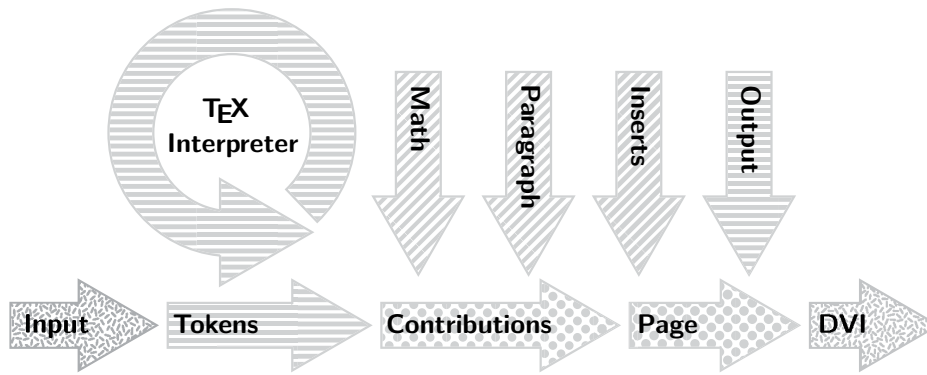
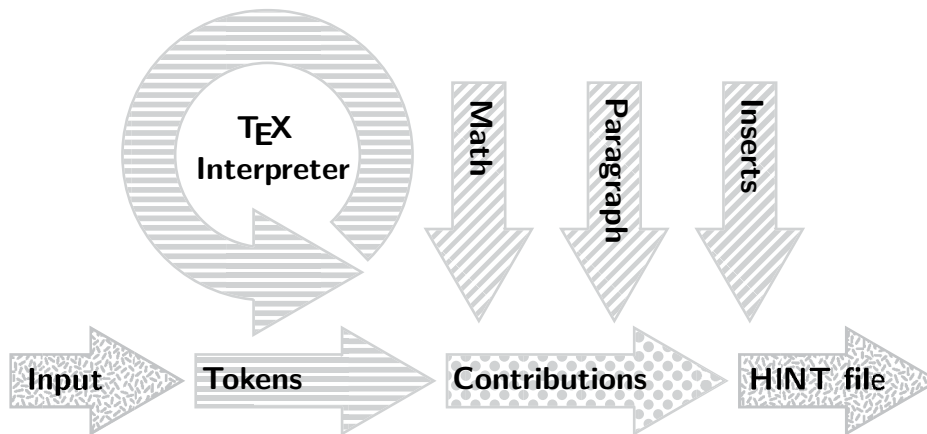
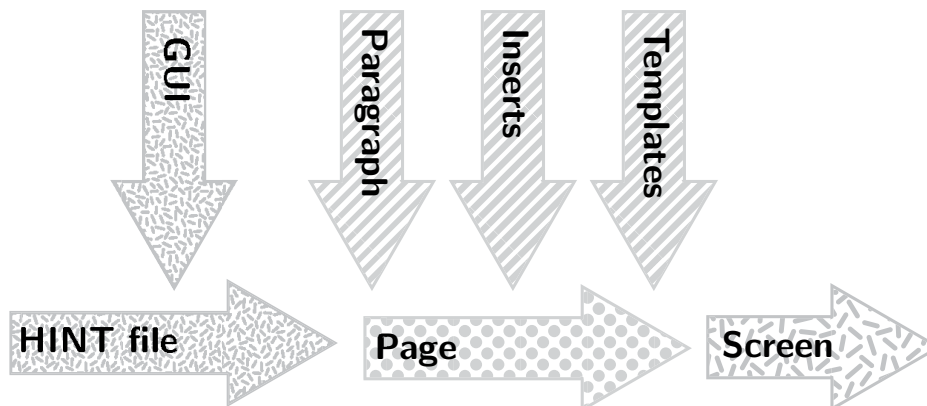
Figure 1: The structure of \TeX Figure 2: The structure of HiTeX 

Figure 3: The structure of HINT

components of \TeX needed to be moved back and forth between front- and backend before a satisfactory separation was accomplished.

Before I begin to describe these in more detail, I want to emphasize that the current state of the file format and the two utilities is not the end-point of development but a starting point. While I hope that the current specification provides enough func-

tionality to attract a first small community of users, I see it more as a test-bed for experimentation with reflowable \TeX output leading to better concepts, better formats, and better implementations.

Further, I consider the HINT viewer and its file format, while derived from \TeX , as \TeX independent. Why should not, for example, OpenOffice have a plug-in producing HINT output files?

```

HINT 1.0
<directory 4 (lists resources)
  <section 3 'TeXfonts/cmr10.tfm'>
  <section 4 'TeXfonts/cmr10.600pk'>
>
<definitions (lists definitions)
  <max <font 0>> (using just font 0)
  <font *0 'cmr10' 3 4
    <glue *13 (space skip)>
    <hyphen "-" 0 (default hyphen)>>
>
<content (a paragraph showing a kern)
  <par *0 "Hello w<kern -0x0.471D pt>orld!">
>

```

Figure 4: Example HINT file in long format

File formats

There are two file formats: a short form that represents HINT files as a compact byte stream for the viewer and a long form that represents HINT files in a readable form for editing and debugging. Figure 4 gives an example of the latter. Note the hexadecimal floating point notation in the kern node which is an exact representation of TeX’s “scaled points”.

After reading a HINT file, we have a byte stream in memory. This stream contains the directory, the definitions, the content stream, and finally resources. In the definition part, we define fonts and associate them with font-numbers for compact reference and do similar things for glues and other units that are used frequently. We supplement the definitions by setting suitable defaults. Then follows a content stream of at most 4GB. The latter restriction ensures that positions inside the content stream can be stored in 32 bits. The content stream consists of a list of nodes; each node representing a glue, a kern, a ligature, a discretionary hyphen, . . . , or a box. Of course the content of boxes is again a list of nodes. After the content stream, we store file resources, for example image and font files.

If we want the viewer to support changing the page size while moving around in the stream — going to the next or previous page, following a link or using an index — practically any position in the stream can be the start of a page. This makes precomputing page starts impossible.

As a consequence, we need to be able to parse the content stream forward and backward. A node in the content stream therefore has a start byte, from which the parser can infer the structure and size of the node, and the same byte again as an end byte. Given an arbitrary position in the stream, it is possible to check if the current byte is a start byte or an end byte by computing the node-length from it and check the stream at the computed position

for a matching byte. To be sure that the match is not a coincidence, the process can be repeated for a sequence of several nodes.

Start and end bytes contain a 5-bit “kind” and a 3-bit “info” field. This allows for 32 different kinds of nodes. The info bits can be used for small parameters or flags, or indicate the absence of certain fields in the node.

Lists. A special case is nodes describing lists of nodes. The method described above to distinguish start and end bytes is not feasible for a list of nodes because it is not possible to compute the size of the list from the start or end byte. Therefore, we store the size of the list content after the start byte and before the end byte. The three info bits are used to indicate whether the size is stored as 0, 1, 2, 3, or 4 bytes. This scheme enables a parser to find the corresponding start or end byte. Specifying 0 bytes for the size implies an empty list.

Texts. Because many lists consist mostly of characters, there is a special list format optimized for storing character nodes. We call such a list a “text”. The start and end bytes of a text are like those of ordinary lists, but they are of kind “text”. Only forward parsing is supported for a text node. Using the size information, we can skip easily to the beginning of a text.

A text can be thought of as a list of integers. Small integers in the range 0 to 127 are stored as single bytes; for larger integers the multi-byte encoding from UTF-8 is used. The integers from 0 to 32 are considered control codes, and all other integers are considered character codes — or rather glyph-numbers, to be more precise. The control codes are used for a variety of purposes. For example, a glyph-number in the range 0 to 32 can be specified by using the control code 0x1D followed by the glyph-number; or arbitrary nodes can be inserted in the text after the control code 0x1E.

A glyph-number references a specific glyph in the current font; the current font in a text is given implicitly. The control codes 0x00 to 0x07 can be used to select the 8 most common fonts; other fonts can be selected by using the control code 0x08 followed by the font number.

hsize and vsize

TeX treats `hsize` and `vsize` like any other dimension register; you can set them to any value and do all kinds of computations with them. HiTeX is more restrictive. At the global level, you cannot change `hsize` and `vsize` at all, because they denote the

dimensions given in the viewer. \TeX , however, allows local modifications of dimension registers; for instance you can say `\vbox{\hsize = 0.5\hsize \advance \hsize by -8pt ...}` to obtain a vertical box, and inside this box, the value of `hsize` is just a bit smaller than half its global size. Hence, paragraphs inside this box are broken into lines that are almost half a page wide. The value of `hsize` will return to its old value once the box is completed. To make this possible, HiTeX treats dimensions as linear functions $\alpha + \beta \cdot \text{hsize} + \gamma \cdot \text{vsize}$, where α , β , and γ are constants. Computations are allowed as long as they stay inside the set of linear functions. For example `\multiply \hsize by \hsize` would not work. For lack of a better name, such a linear function is called an “extended dimension”. The good news is that the viewer can convert an “extended dimension” immediately to a normal dimension since in the viewer `hsize` and `vsize` are always known.

Paragraphs

Breaking paragraphs into lines is \TeX ’s most sophisticated and complex function. Fortunately the implementation is very efficient (it used to run fairly smoothly on my 8MHz 80286). It needs to be present in the frontend and in the backend. If `hsize = α` is a known constant (with $\beta = \gamma = 0$), the frontend can perform the line breaking; if $\beta \neq 0$ or $\gamma \neq 0$, line breaking must be performed in the backend.

On the other hand, we do not want the backend to perform hyphenation. Hyphenation is an expensive operation; it requires hyphenation tables to be present; and then it would be impossible for an author or editor to check the correctness of hyphenations. Therefore HiTeX will always insert all the discretionary hyphens that \TeX would compute normally in the second pass of its line breaking algorithm. To reproduce the exact behavior of \TeX ’s line breaking algorithm, the discretionary hyphens found in this way are marked and are used only during the second pass in the viewer. This gives preference to line breaks that do not use hyphenation (or only user specified discretionary hyphens) in the same way as \TeX .

The paragraph shape is controlled by the variables `hangindent`, `hangafter`, and `parshape` which can be used to specify an individual indentation and length for any line in the paragraph. Obviously, these computations must be performed in the viewer. A complication arises if the viewer needs to start a page in the middle of a paragraph: the line number of the first line on the new page, and with it its indentation and length, then depends on how

the previous page was formatted. This might not be known, for example if paging backward or if the page size has changed since the viewer had formatted the previous page. It remains an open question what gives the best user experience in such a situation.

Packing boxes and alignment

\TeX knows two kinds of boxes: horizontal boxes, where the reference points of the content are aligned along the baseline; and vertical boxes, where the content is stacked vertically. Let’s look at horizontal boxes; vertical boxes are handled similarly.

\TeX produces horizontal boxes with the function `hpack`. The function traverses the content list and determines its total natural height, depth, and width. Furthermore, it computes the total stretchability and shrinkability. From these numbers it computes a glue ratio such that stretching or shrinking the glue inside the box by this ratio will make the box reach a given target width. HiTeX faces two problems: It might not be possible to determine the natural dimensions of the content, because, for example, the depth of a box can depend on how the line breaking algorithm forms the last line of a paragraph. In this case packing the box with `hpack` must be done in the viewer. But even if the natural dimensions of the content can be determined, a target width that depends on `hsize` will prevent HiTeX from computing a glue ratio. Therefore the `HINT` format knows three kinds of horizontal boxes: those that are completely packed, those that just need the computation of a glue ratio, and those that need a complete traversal of the box content.

Handling \TeX ’s alignments introduces a little extra complexity. When \TeX encounters a horizontal alignment, it packs the rows into `unset` boxes adding material from the alignment template and the appropriate `tabskip` glue. After all rows are processed, \TeX packs the rows using the `hpack` function. At that point HiTeX can use the mechanisms just described for ordinary calls of `hpack`.

Baseline skips

When \TeX builds vertical stacks of boxes, typically lines of text, it tries to keep the distances between the baselines constant, that is: independent of the actual depth of descenders or height of ascenders. Three parameters govern the insertion of glue between two boxes in vertical mode: \TeX will insert glue to make the distance between baselines equal to `baselineskip` unless this would make the glue smaller than `lineskiplimit`; in the latter case, the

glue is set to `lineskip`. Additional white space between boxes, for instance a `\vskip 2pt`, does not interfere with this computation. Instead, \TeX uses the variable `prev_depth`, containing the depth of the last box added to the list, for the computation. This offers a convenient lever for authors and macro designers to manipulate \TeX 's baseline calculations. For example setting `prev_depth` to the value `ignore_depth` will suppress the generation of a `baselineskip` for the next box on the list. This is of course a fact that the viewer should know about.

The HINT format is designed to be “stateless”, that is: given the position of a page break in the stream, it is possible to read, understand, and format the page starting at that position or the page ending at that position. This turns the insertion of baseline skips into an interesting problem: In simple cases, when all relevant information is at hand, $\text{Hi}\TeX$ can insert the correct glue directly. If some information is missing, a baseline node is generated. To process such a baseline node, the current values of the parameters mentioned before are required, and these parameters do change occasionally.

Storing the current values in every baseline node would require up to 54 bytes per node. HINT uses a more space-efficient approach: It defines default values that are constant for the entire stream. A baseline node using the defaults does not need to specify parameters. Further, the definition part of the stream can specify up to 256 baseline definitions, each defining the full set of parameters; such a parameter set can be used by specifying its number in a single byte. Only in the rare case that these two mechanisms are not sufficient must the baseline node contain the necessary values directly. The same approach is used for glues, extended dimensions, paragraphs, and displays. It can be generalized to arbitrary parameter sets.

Displayed equations

The positioning of displayed equations in \TeX is no simple task. Usually the formula is centered on the line, but if `hsize` is so small that the formula would come too close to the equation number, it is centered in the remaining space between equation number and margin; if `hsize` is even smaller, the equation number will be moved to a separate line. Vertical spacing around the formula depends on the length of the last line preceding the display, which in turn depends on the outcome of the line breaking algorithm. If the line is short enough, \TeX will use the `abovedisplayshortskip` glue, otherwise it uses `abovedisplayskip`. Of course there is also

`belowdisplayshortskip` and `belowdisplayskip` to go with them. In addition, the variables controlling the paragraph shape influence the positioning of the displayed equation. The required computations must be done in the viewer; they are not very expensive but the code is complicated. HINT uses display nodes to describe displayed formulas. Fortunately, none of the math mode processing need be done in the viewer.

Images

Native \TeX does not define a mechanism for including images, instead providing a generic extension mechanism. For the HINT viewer to be able to open and display any correct HINT file, we need to specify the image types that a viewer is required to support and the exact format of the image nodes. Image files are included in the resource part of the HINT file and are referenced by defining an image number, its position, and its size in the definition part.

For simplicity, the HINT viewer will not do any image manipulation except scaling. Scaling will be necessary to display the same HINT file on a wide variety of devices in a user friendly way. Various designs for the syntax and semantics of image nodes are possible and only the experience of real users will tell what is good or useless.

At present, images are treated like two dimensional glue: you can specify a width or a height, a stretchability, and a shrinkability. If neither width nor height are given, the natural width and height will be taken from the image file. When an image is part of the content of a box, it will stretch or shrink together with other glue to achieve the target size of the box. This mechanism works surprisingly well in practice; the image and the white space surrounding it scale in a consistent way to fill the space that is assigned to it by the enclosing box.

Page building

\TeX 's page builder starts at the top of a new page and collects vertical material, keeping track of its natural height, stretchability, and shrinkability until the page is so full that possible page breaks can only get worse. Then it uses the best page break found so far and moves remaining material back to the contribution list. Of course it also accounts for the size of inserts, and it uses the penalties found to estimate the goodness of a page break. HINT uses the same algorithm, complementing it with a reverse version that starts at the bottom of a new page. The reverse version is used when paging backward.

At the point where $\text{T}_{\text{E}}\text{X}$ calls the output routine, a new mechanism is needed, because (as mentioned above) we want the viewer to be simple, thus precluding the use of the $\text{T}_{\text{E}}\text{X}$ interpreter that would be necessary to execute a general output routine. **HINT** replaces output routines by page templates, but before we can describe this mechanism, it is necessary to see how **HINT** handles insertions.

Insertions. The $\text{T}_{\text{E}}\text{X}$ page builder identifies different insertions by their insertion number. It accounts for the contribution of inserted material to the total page height by weighting the insertion's natural height by the insertion scaling factor. There is also a constant overhead that needs to be added if the insertion is nonempty, for example the space occupied by a footnote rule and the space surrounding it.

HINT uses the concept of content streams for this. Stream number zero is used for the main page content; other stream numbers are defined in the definition part of the **HINT** file along with stream parameters such as the insertion scaling factor and the maximum vertical extent e that the stream content is allowed to occupy on the page. **HiT_EX** maps insertion numbers to stream numbers and appends the insertion nodes to the content stream.

Streams have some more parameters: a list b of boxes that is used before and a list a that is used after the inserted material if it is not empty; the topskip glue g that is inserted between b and the first box of inserted material reduced by the height of this box; a stream number p , where the material from this stream should go if there is still space available for stream p ; a stream number n , where the material from this stream should go if there is no more space available for the stream but still space available for stream n ; a split ratio r that, if positive, specifies how to split the material of the stream between streams p and n .

The latter stream parameters are new and offer a mechanism to organize the flow of insertions on the page. For example, when plain $\text{T}_{\text{E}}\text{X}$ encounters a floating insertion, it decides whether there is still enough space on the current page and if so makes a mid-insert; otherwise a top-insert. **HiT_EX** needs to postpone this decision. It will channel such an insertion to a stream with $e = 0$, $p = 0$, and n equal to the stream of top-inserts. When such an insertion arrives at the **HINT** page builder, it will check whether there is still space on stream 0, the main page, and if so moves the insertion there. Otherwise, setting the maximum extent e to zero forces the page builder to move the insertion to the stream n of top-inserts.

If the split ratio r is nonzero, the splitting of the stream will be postponed even further: The page builder will collect all contributions for the given stream and will split it in the given ratio between streams p and n just before assembling the final page. For example it is possible to put all the footnotes in one stream with an insertion scaling factor of 0.5 and split the collected footnotes into two columns using a split ratio of 0.5; with a cascade of splits, three or more columns are also possible.

Marks. $\text{T}_{\text{E}}\text{X}$ implements marks as token lists, and the output routine has access to the top, first, and bottom mark of the page. Sophisticated code can be written to execute these token lists producing very flexible headers or footers. In **HINT** we cannot use token lists but only boxes. Consequently, **HINT** uses the stream concept, developed for insertions, and extends it slightly. A flag can be added to a stream designating it as a “first” or “last” stream. Such a stream will retain at most one insertion per page. Now a package designer can open a stream for first marks and a stream for bottom marks, put $\text{T}_{\text{E}}\text{X}$'s marks into boxes, and add them into both streams. The implementation of top marks is difficult because it requires processing the preceding page. Top marks are not part of the present implementation.

Templates. Once the main page and all insertions are in place, **HINT** needs to compose the page. For this purpose it is possible in **HiT_EX** to specify one or more page templates. A page template is just a **vbox** with arbitrary content: boxes, glue, rules, alignments, . . . , and, most importantly, inserts. **HiT_EX** will store the output template in the definition part together with its valid range of stream positions. When **HINT** needs to compose the page, it will search for an output template that includes the stream position of the current page in its range. It makes a copy of the template replacing each insert node by the material accumulated for it — insert node 0 will be replaced by the content of the main page. Material given as parameters a and b of an insert stream will be copied as necessary. After repacking the resulting **vbox** and all its subboxes, the **vbox** will be rendered on the display.

Implementation

For the work described above, I needed to make substantial changes to the $\text{T}_{\text{E}}\text{X}$ source code. The common tool chain from $\text{T}_{\text{E}}\text{X}$ Live uses **tangle** to convert **tex.web** into Pascal code (**tex.pas**) which is then translated by **web2c** [6] into C code. Already the translation to Pascal code expands all macros

and evaluates constant expressions, because neither is supported by Pascal. As a result, the generated Pascal code, let alone the further translation to C, becomes highly unreadable and cannot be used as a basis for any further work. So I wrote a translator converting the original WEB source code of T_EX into cweb source code [2, 3]. This cweb source is the basis of the development of HiT_EX and HINT.

For the implementation of HiT_EX and HINT, I had only limited time at my disposal: my sabbatical during the fall semester of 2017/2018. As a consequence, I often moved on as soon as the current research problem had changed — in my view — into an engineering problem. This allowed me to make fast progress but left lots of “loose ends” in the code.

The current prototype has the functionality of Knuth’s T_EX with the adaptations described above, and without added features like search paths for input files or PDF specials. It is capable of generating format files for plain T_EX or L^AT_EX and it can handle even large files. The code for paging backwards is buggy because I occasionally implemented new features in the forward page builder and neglected to update the backwards page builder accordingly.

Open questions and future work

Conditionals. It seems reasonable to implement different output templates depending on screen size and aspect ratio. Also conditional content, for example a choice between a small and a wide table layout, might be useful. For a whole list of ideas, see [1].

Macros for L^AT_EX support. Since the input part of HiT_EX is taken directly from T_EX, basic L^AT_EX is supported. But since L^AT_EX uses complex output procedures, many macros might need changes with variable page sizes now in mind. Templates are still an experimental feature of HINT that might need changes to better support L^AT_EX.

Usage of control codes. Three control codes used in texts are indispensable: based on the bytes that follow, one control code switches to any of 256 possible fonts, one specifies an arbitrary character code, and one specifies an arbitrary node. The remaining 30 control codes provide plenty of room for experiments. Currently 8 of them are dedicated to font selection, 8 to reference globally predefined nodes, and 14 to reference font-specific predefined nodes. This should allow a convenient and compact encoding that can accomplish the most common operations with a single byte and use two or more bytes for less common operations. To decide whether the current dedication is optimal in this respect is an

open question. A statistical analysis using a large collection of T_EX documents should give an answer.

Images. The implementation of glue-like images is experimental. Another obvious idea is the specification of background (and foreground) properties of boxes. The background could be a color (making rules a special case of boxes), a shading, or an image that can be stretched, or tiled, or positioned to fill the box. Certainly this would extend the capabilities of HINT beyond the necessities for T_EX. Is this a direction worth considering?

Because it was easy to implement, currently only Windows bitmaps are supported. A full implementation should certainly support also JPEG and PNG files, and some form of vector graphic, probably SVG. I think it is better to have a small collection of formats that is well supported across all implementations than a long list of formats that enjoy only limited support. But how about sound and video? Should there be support? How could an extension mechanism look that keeps the HINT format open for future development?

Platforms. Currently the HINT viewer is written for the Windows platform just because this was convenient for me. Since HINT targets mobile devices, a HINT viewer for Android would be a next logical step. I also think that ebook readers deserve a better rendering engine and HINT would be a candidate.

References

- [1] H. Hagen. Beyond the bounds of paper and within the bounds of screens; the perfect match of T_EX and Acrobat. In *Proceedings of the Ninth European T_EX Conference*, vol. 15a of *MAPS*, pp. 181–196. Elsevier Science, September 1995. nlg.nl/maps/15a/09.pdf
- [2] M. Ruckert. Converting T_EX from WEB to cweb. *TUGboat* 38(3):353–358, 2017. tug.org/TUGboat/tb38-3/tb120ruckert.pdf
- [3] M. Ruckert. *web2w: Converting T_EX from WEB to cweb*, 2017. ctan.org/pkg/web2w
- [4] C. Schenk. Yap: Yet another previewer. miktex.org
- [5] P. Vojta. Xdvi. math.berkeley.edu/~vojta/xdvi.html
- [6] *Web2C: A T_EX implementation*. tug.org/web2c

◇ Martin Ruckert
Hochschule München
Lothstrasse 64
80336 München, Germany
[ruckert \(at\) cs dot hm dot edu](mailto:ruckert@cs.hm.edu)