



# Experience from a Real-World Application of Micro-Typography with pdfT<sub>E</sub>X

Hàn Th<sub>ế</sub> Thành

River Valley Technologies [hanthethanh@gmail.com](mailto:hanthethanh@gmail.com)

**KEYWORDS** micro-typography, margin kerning, font expansion, optimal page-breaking, semi-automatic pagination, automatic problem detection, paragraph stretching/shrinking.

**ABSTRACT** This article describes the experience from a real-world application of micro-typography with pdfT<sub>E</sub>X. The project involved typesetting a study edition of the Bible in Czech, where there was a lot of further information apart from the original text: footnotes, references, further annotation etc. The design was not complex, but the typographic requirements were very strict. This posed a real challenge to typesetting everything without conflicts. The biggest problem was how to achieve perfect page-breaking without changing a single word of the text. The solution was to use a semi-automated method: manual breaking of each page, with auto-detection of problematic pages (like orphan/widow, too little/too much space between body and footnotes, etc.). When such a problem occurred, it had to be fixed by changing the length of one or more paragraphs on that page until the problem disappeared. Since we are not allowed to modify the text, the length of a paragraph could be changed only by changing the formatting of the paragraph (`\looseness`), in order to make it longer or shorter. This usually results in poor-looking paragraphs. However, with the aid of micro-typography provided by pdfT<sub>E</sub>X, this was achieved without loss of quality.

## 1 Introduction

In 2000, a publisher in Czech Republic needed to re-publish a study edition of the Bible (only New Testament). The previous editions were typeset using Ventura. The publisher however was not very satisfied with the result, and wanted to achieve better result. They heard that a program called T<sub>E</sub>X can do that job nicely, so they started looking for someone willing to typeset the Bible in T<sub>E</sub>X. A T<sub>E</sub>X fellow introduced me to the publisher. At that time the micro-typographic extensions of pdfT<sub>E</sub>X were introduced and this was a perfect application for the new extensions. It had quite a good impact on pdfT<sub>E</sub>X development, since during this project many issues were found and then were fixed or improved. This was perhaps the first serious application of micro-typography in pdfT<sub>E</sub>X and was an important encouragement to me that the micro-typographic extensions can be really useful in practice.

In 2007 the publisher re-published the next edition of the Bible, and I did the typesetting again. Many things were easier than the previous time. However some problems remain difficult in principle and had to be solved with the “old tricks”. In this

article I would like to share the experience from the project, with the hope it might be interesting for the reader, and might be useful for someone who has to solve similar problems.

## 2 The workflow from the publisher's view

The workflow from the publisher's view (Figure 1) is rather simple and does not involve T<sub>E</sub>X. The publisher maintains primary data in an in-house format which is not known to the typesetter. The data are exported to XML and sent to the typesetter. Whatever the typesetter does with the XML, the publisher needs not know: only the PDF result is relevant.

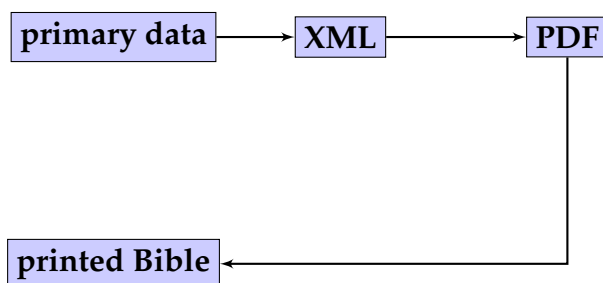


FIGURE 1. The workflow from the publisher's view

The choice of XML as the format to exchange data was a very good decision, since it frees the typesetter from knowing anything about the data format of the publisher. The publisher has their own tools to maintain and manipulate the data. The data format might be either very simple or complex, the tools might be either very reliable or buggy, but it is irrelevant to the typesetter. A DTD was developed for the project, and all what matters is whether the XML data from the publisher pass the validation against that DTD. If not, the problem is on the publisher's side and it is the responsibility of the publisher to check the data/tools and re-generate the XML data. If yes, then both sides know that the XML data are correct and if the PDF output is wrong, it is the typesetter's fault. So, the DTD acts as a contract for data exchange between the publisher and the typesetter. It helped to find out quickly who is responsible if a problem arises. Indeed, during the project many problems have been detected very early thanks to XML validation and that saved both sides from further communication problems.

## 3 The workflow from the typesetter's view

The workflow from the typesetter's view is more complicated (Figure 2).

When the typesetter receives the XML's, the following checks have to be done:

*Validation:* as mentioned above, validation is the first thing to do when the typesetter gets the XML data. If validation fails, the typesetter informs the publisher about the problem and waits for the updated data. The publisher is of course supposed

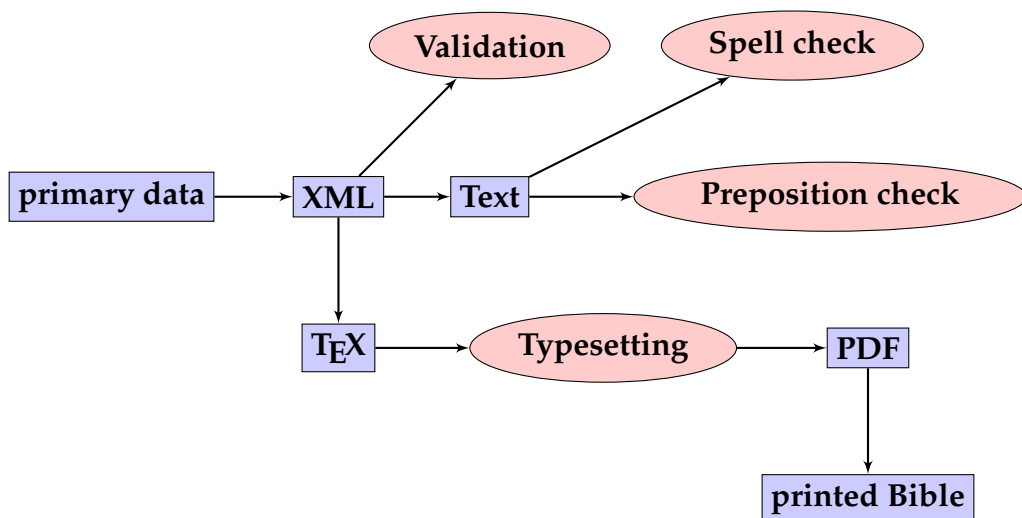


FIGURE 2. The workflow from the typesetter's view

to validate the XML data before sending them, but the typesetter should not rely on it. Instead, the typesetter must check and ensure the XML data are valid before processing further.

*Text processing:* the next step is to convert XML data to text to check for possible problems with the text contents:

*Spell check:* it might look a bit unlogical that spell checking is done here. Spell checking should be done on the publisher's side within the primary data format and not by the typesetter. However, during the typesetting a few misspelt words were found. Since it is not much work for the typesetter to do a spell check, it was added to the process just as a doubled check.

*Preposition check:* the czech language has many single-letter prepositions: **k**, **o**, **s**, **u**, **v**, **z**. Those prepositions should not end up at the end of a line. The conventional way to deal with those single-letter prepositions is to use a non-breakable space after them (~). The publisher is responsible for ensuring that tilde is used instead of space in such cases. However it happens from time to time that a tilde is missing. For that reason, another doubled check was added for this case.

If there is any problem in this stage, the publisher is informed about the problem and must send again the XML data.

*Conversion to T<sub>E</sub>X:* if the XML data pass the previous checks, they are converted to T<sub>E</sub>X, ready for typesetting.

All the checking steps above are done by some scripts, so anytime the XML data are updated, the typesetter can tell immediately if there is any problem and provide feedbacks to the publisher. If no problem is found, the T<sub>E</sub>X data are then used for the next stage: typesetting.

## 4 Typographic requirements

The layout of this typesetting is rather simple, as it can be seen in Figure 3.

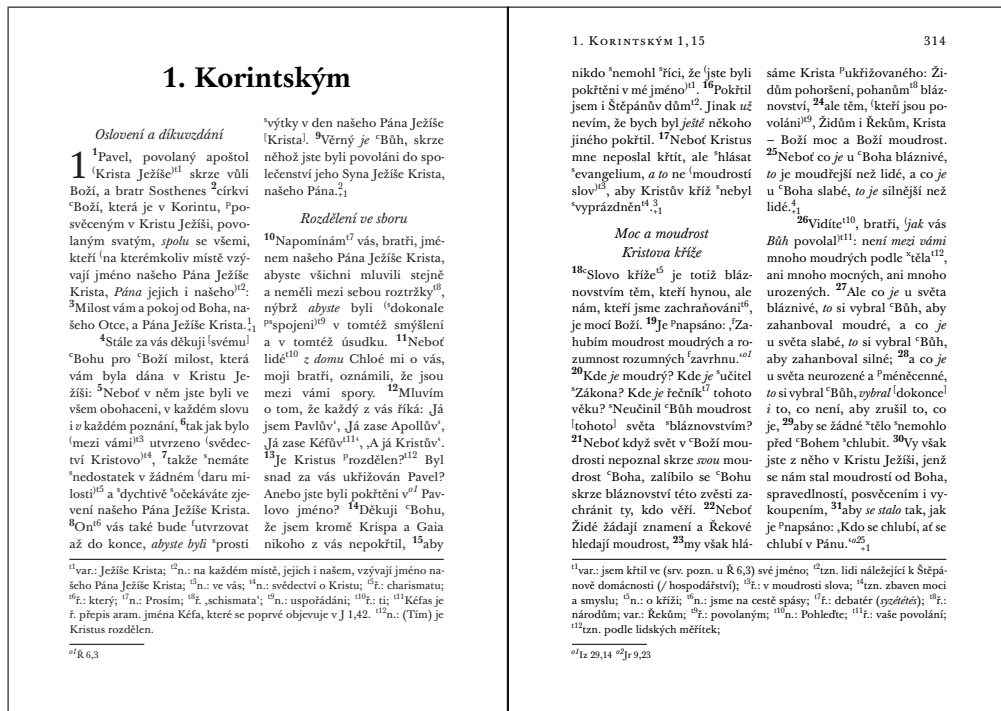


FIGURE 3. The sample layout

However, there are some requirements that are not that visible from the sample:

- Text contents must not be altered at all. As the XML data are converted to T<sub>E</sub>X, such data must be used for typesetting without a single manual change.
- The line spacing is fixed (i. e. all lines must sit on a “grid”).
- The text contents are rich with many elements apart from the main text: subtitles, footnotes, references, various marking to aid the reader, etc.
- There must not be any typographic “sin”:
  - no widow,
  - no orphan,
  - chapter number not end at page bottom (no “chapter orphan”),
  - not too much or too little gap between footnote and body,
  - the position where footnote ends is fixed,
  - always balanced columns.
- If there is any last-minute change in primary data, it must be reflected in the PDF quickly.

The conventional  $\text{\TeX}$  approach to solve problem with undesirable page breaks (to avoid widow/orphan) is to insert penalty at relevant places. However, this approach mostly leads to page breaking with very big gap between footnotes and body. To achieve optimal page breaking within that so restrictions, human decision and intervention is required.

## 5 How to achieve optimal page breaking?

When a document needs to be paginated optimally, the typesetter usually has to do two main things:

1. ensure that every problem is found (i. e. do not overlook any issue);
2. for a given issue, find an optimal way to solve the problem.

Solving a problematic case is usually done by changing the formatting of various elements around the problematic place. For example, the typesetter can change the gap around a figure, a paragraph, or a subtitle. Another example is to change the page dimensions (i. e. `\enlargethispage` in  $\text{\LaTeX}$ ), or the gap between body and footnote, etc. In difficult cases where the above didn't work, the typesetter might have to change the formatting of a paragraph, or even the wording.

When we apply this approach to our case, it means that we must not miss anything from the issues mentioned in Section 4 and for each of them, find an optimal way to solve the problem.

However, in our case the possibility to change the formatting of an element is very limited: since grid typesetting is required, we cannot change spacing around any subtitle or paragraph. The gap between body and footnote cannot vary too much. Therefore, the only thing we can adjust is to change the formatting a particular paragraph, i. e. to make the formatting of a paragraph longer or shorter. Fortunately,  $\text{\TeX}$  provides a parameter to control this: `\looseness`. This is a  $\text{\TeX}$  parameter which means roughly the following: when `\looseness=1`,  $\text{\TeX}$  should try to format a paragraph to be one line longer than the default formatting. For example, if a paragraph is formatted to be 10 lines long by default, setting `\looseness=1` would cause the formatting of the paragraph to be 11 lines long.

Our approach to tune page breaking relies heavily on use of `\looseness`: when a problematic case is found, we change `\looseness` of relevant paragraphs to solve the problem. We can see this as a way to "stretch" or "shrink" a paragraph: stretching a paragraph is done by setting `\looseness=1`, and shrinking a paragraph is done by setting `\looseness=-1`. Setting `\looseness` does not always result in the desired effect: sometimes there is no way to make the formatting of a paragraph longer or shorter. So, a paragraph can be "stretchable" or "shrinkable", depending on whether changing `\looseness` gives the desired effect.

To illustrate how it works, let's look at an example. In Figure 4, there is a widow line. To avoid this widow line, the optimal way is to make some previous paragraph shorter so that the widow line would be moved to the left column. We do so by "shrinking" the second paragraph on the previous page, which means that we set `\looseness=-1` for that paragraph. The effect of this change is shown in Figure 5.

<p>1. KORINTSKÝM 2, 16</p> <p style="text-align: right;">316</p> <p style="text-align: center;"><i>Projevy nezralosti</i></p> <p><b>3</b> <sup>1</sup>A já, bratři, jsem k vám nemohl mluvit jako k duchovním, ale jako k <sup>2</sup>tělesným, jako k nedospělým v Kristu. <sup>3</sup>Dal jsem vám <sup>4</sup>pit mléko, pokrm <i>jem vám nedal, neboť ten jste ještě nemohli sníst</i>. Ale ani teď ještě nemůžete, <sup>5</sup>neboť jste <i>stále ještě</i> <sup>6</sup>tělesní. Vždyť pokud je mezi vámi žárlivost, rozbroj <sup>7</sup>a rozdělení<sup>1</sup>, zdali nejste <sup>8</sup>tělesní a <sup>9</sup>nežijete jen po lidsku<sup>11</sup>? <sup>4</sup>Když jeden říká: Já jsem Pavlův<sup>1</sup> a druhý: Já Apollův<sup>1</sup>, nejste <i>jako jiní lidé</i><sup>12</sup>? <sup>3</sup>Kdo je Apolos? A kdo je Pavel? Služebníci, skrze něž jste uvěřili, jak každému dal Pán. <sup>6</sup>Já jsem zasadil, Apolos zalil, ale <sup>7</sup>Bůh <sup>8</sup>dával <sup>9</sup>růst. <sup>7</sup>Tedy ani ten, kdo sází, ani ten, kdo zalévá, nejsou <sup>10</sup>něčím <i>zvláštním</i>, ale Bůh, který dává růst. <sup>8</sup>Ten, kdo sází i ten, kdo zalévá, jsou jedno, každý však <sup>9</sup>dostane <i>svou vlastní mzdu</i> podle <sup>10</sup>své<sup>13</sup> námahy. <sup>9</sup>Neboť jsme Boží spolupracovníci; <sup>10</sup>vy jste Boží pole, Boží stavba.<sup>11</sup> <sup>10</sup>Podle <sup>11</sup>Boží<sup>12</sup> milosti, která mi byla dána, jsem jako</p> <p>moudrý stavitel položil základy a jiný <sup>12</sup>na <sup>13</sup>něm <sup>14</sup>staví. Každý <sup>15</sup>at <sup>16</sup>si <sup>17</sup>dává <sup>18</sup>pozor, jak <sup>19</sup>na <sup>20</sup>něm <sup>21</sup>staví. <sup>11</sup>Neboť nikdo nemůže položit jiný základ nežli ten, který je již položen,<sup>12</sup> a tím je Ježíš Kristus. <sup>12</sup>Jestliže někdo na <sup>13</sup>tomto <sup>14</sup>základě <sup>15</sup>staví <sup>16</sup>ze zлата, <sup>17</sup>stříbra, <sup>18</sup>drahých kamenů, <sup>19</sup> dřeva, <sup>20</sup>seny <sup>21</sup>nebo <sup>22</sup>slámy, <sup>13</sup>dílo každého se stane zjevným. Ten den <sup>23</sup>je <sup>24</sup>ukáže, neboť se zjeví<sup>14</sup> v ohni, a oheň vyzkouší dílo každého, jaké je. <sup>14</sup>Jestliže někdo <sup>15</sup>na <sup>16</sup>tomto <sup>17</sup>základě <sup>18</sup>vystaví<sup>15</sup> dílo a ono <sup>19</sup>mu <sup>20</sup>zůstane<sup>16</sup>, dostane odměnu. <sup>15</sup>Jestliže <sup>16</sup>mu <sup>17</sup>jeho <sup>18</sup>dílo <sup>19</sup>shoří, <sup>20</sup>utrpí <sup>21</sup>škodu; sám se sice zachrání, ale jako skrze oheň. <sup>16</sup>Nevtí, že jste Boží svatyně a <sup>17</sup>že <sup>18</sup>ve<sup>17</sup> vás <sup>19</sup>bydlí <sup>20</sup>“Duch <sup>21</sup>“Boží? <sup>17</sup>Ničí-li někdo <sup>18</sup>Boží <sup>19</sup>svatyni, <sup>20</sup>zničí <sup>21</sup>Bůh <sup>22</sup>jej. Neboť <sup>23</sup>Boží <sup>24</sup>svatyně je <sup>25</sup>svatá, <sup>26</sup>a <sup>27</sup>ta <sup>28</sup>svatyní<sup>18</sup> <sup>29</sup>jste <sup>30</sup>vy.<sup>1-11</sup> <sup>18</sup>At <sup>19</sup>nikdo <sup>20</sup>neklame <sup>21</sup>sám sebe. Jestliže si někdo mezi vámi myslí, <sup>22</sup>že <sup>23</sup>je <sup>24</sup>v <sup>25</sup>tomto <sup>26</sup>věku<sup>19</sup> <sup>27</sup>moudrý, <sup>28</sup>ať <sup>29</sup>se <sup>30</sup>stane <sup>31</sup>bláznem, aby se stal moudrým. <sup>19</sup>Neboť <sup>20</sup>moudrost <sup>21</sup>tohoto <sup>22</sup>světa <sup>23</sup>je <sup>24</sup>před <sup>25</sup>mou</p> <p><sup>1</sup>ř.; nechodíte podle člověka; <sup>2</sup>var.: tělesní; <sup>3</sup>ř.: vlastní; <sup>4</sup>ř.: se odhaluje; <sup>5</sup>ř. přistavě (na ten základ); podobně ve v. 10 a 12; <sup>6</sup>var.: zůstává; <sup>7</sup>n.: mezi vámi; <sup>8</sup>ř.: kterou; <sup>9</sup>řzn. podle měřiček tohoto světa;</p> <p><sup>11</sup>řz 28,16</p>	<p>317</p> <p style="text-align: right;">1. KORINTSKÝM 4, 13</p> <p>hem bláznovstvím. Vždyť je <sup>1</sup>napsáno: „On<sup>11</sup> chytá moudré v jejich chytráctví.“<sup>12</sup> <sup>20</sup>A opět: „Pán zná myšlenky moudrých a ví, že jsou marné.“<sup>21</sup> <sup>21</sup>A tak <sup>22</sup>at se <sup>23</sup>nikdo <sup>24</sup>nechlubí <sup>25</sup>lidmi. Vždyť všechno je vaše: <sup>22</sup>at Pavel nebo Apolos nebo Kéfas, <sup>23</sup>at svět nebo život nebo smrt, <sup>24</sup>at <sup>25</sup>věci <sup>26</sup>přítomné nebo budoucí – všechno je vaše, <sup>27</sup>vy <sup>28</sup>pak <sup>29</sup>jste <sup>30</sup>Kristovi <sup>31</sup>a <sup>32</sup>Kristus <sup>33</sup>Boží.<sup>10</sup></p> <p style="text-align: center;"><i>Služebníci Kristovi</i></p> <p><b>4</b> <sup>1</sup>At <sup>2</sup>o <sup>3</sup>nás <sup>4</sup>každý<sup>2</sup> <sup>5</sup>smýšlí jako <sup>6</sup>služebníci <sup>7</sup>Kristovi a <sup>8</sup>správcích <sup>9</sup>Božích <sup>10</sup>tajemství. <sup>2</sup>A od <sup>3</sup>správčů <sup>4</sup>se <sup>5</sup>nakonec <sup>6</sup>vžaduje, <sup>7</sup>aby <sup>8</sup>byl <sup>9</sup>každý <sup>10</sup>shledán <sup>11</sup>věrným. <sup>3</sup>Pro <sup>4</sup>mne <sup>5</sup>je <sup>6</sup>“pramálo <sup>7</sup>důležité, <sup>8</sup>zda <sup>9</sup>mě <sup>10</sup>posuzuje <sup>11</sup>vy nebo <sup>12</sup>lidský <sup>13</sup>soud<sup>13</sup>. <sup>4</sup>Ale <sup>5</sup>ani <sup>6</sup>já <sup>7</sup>sám <sup>8</sup>sebe <sup>9</sup>neposuzuji. <sup>4</sup>Ničeho, <sup>5</sup>co <sup>6</sup>by <sup>7</sup>svědčilo <sup>8</sup>proti <sup>9</sup>mně, <sup>10</sup>si <sup>11</sup>sice <sup>12</sup>nejsem <sup>13</sup>vědom, <sup>14</sup>ale <sup>15</sup>tím <sup>16</sup>nejsem <sup>17</sup>ospravedlněn; <sup>10</sup>ten, <sup>11</sup>kdo <sup>12</sup>mne <sup>13</sup>posuzuje, <sup>14</sup>je <sup>15</sup>Pán. <sup>3</sup>Proto <sup>4</sup>nic <sup>5</sup>nesudte <sup>6</sup>předčasně, <sup>7</sup>dokud <sup>8</sup>nepřijde <sup>9</sup>Pán, <sup>10</sup>kteřý <sup>11</sup>osvítí <sup>12</sup>věci <sup>13</sup>skryté <sup>14</sup>ve <sup>15</sup>tmě <sup>16</sup>a <sup>17</sup>zjeví <sup>18</sup>úmysly <sup>19</sup>srdcí. <sup>4</sup>A <sup>5</sup>tehdy <sup>6</sup>se <sup>7</sup>každému <sup>8</sup>dostane <sup>9</sup>pochvaly <sup>10</sup>od <sup>11</sup>Boha.<sup>11</sup></p> <p><sup>6</sup>Toto <sup>7</sup>jsem, <sup>8</sup>bratři, <sup>9</sup>vztáhl na <sup>10</sup>sebe <sup>11</sup>a <sup>12</sup>na <sup>13</sup>Apolla <sup>14</sup>kvůli <sup>15</sup>vám, <sup>16</sup>abyste <sup>17</sup>se <sup>18</sup>na <sup>19</sup>nás <sup>20</sup>naučili <sup>21</sup>“smýšlet<sup>1</sup> <sup>22</sup>ne <sup>23</sup>nad <sup>24</sup>to, <sup>25</sup>co <sup>26</sup>je <sup>27</sup>“napsáno“, <sup>28</sup>abyste <sup>29</sup>se <sup>30</sup>kvůli <sup>31</sup>jednomu <sup>32</sup>učiteli <sup>33</sup>nenadávali <sup>34</sup>jeden <sup>35</sup>na <sup>36</sup>druhého. <sup>7</sup>Vždyť <sup>8</sup>kdo <sup>9</sup>“<sup>10</sup>ti <sup>11</sup>dává <sup>12</sup>výsledek<sup>14</sup>? <sup>13</sup>Co <sup>14</sup>z <sup>15</sup>toho, <sup>16</sup>co <sup>17</sup>máš, <sup>18</sup>ji <sup>19</sup>nedostal? <sup>14</sup>A <sup>15</sup>když <sup>16</sup>jsi <sup>17</sup>to <sup>18</sup>dostal, <sup>19</sup>proč <sup>20</sup>se <sup>21</sup>chlubíš, <sup>22</sup>jako <sup>23</sup>bys <sup>24</sup>to <sup>25</sup>nedostal? <sup>15</sup>Už <sup>16</sup>jste <sup>17</sup>“nasytzeni, <sup>18</sup>už <sup>19</sup>jste <sup>20</sup>zbohatli, <sup>21</sup>bez <sup>22</sup>nás <sup>23</sup>“jste <sup>24</sup>začali <sup>25</sup>“kralovat. <sup>16</sup>Kéž <sup>17</sup>byste <sup>18</sup>kralovali, <sup>19</sup>ale <sup>20</sup>tak, <sup>21</sup>abychom <sup>22</sup>i <sup>23</sup>my <sup>24</sup>“kralovali <sup>25</sup>“spolu <sup>26</sup>“s <sup>27</sup>vámi. <sup>9</sup>Zdá <sup>10</sup>se <sup>11</sup>mi, <sup>12</sup>že <sup>13</sup>nás, <sup>14</sup>apostoly, <sup>15</sup>Bůh <sup>16</sup>postavil <sup>17</sup>jako <sup>18</sup>poslední, <sup>19</sup>jako <sup>20</sup>“odsouzené <sup>21</sup>na <sup>22</sup>“smrt, <sup>23</sup>neboť <sup>24</sup>jsme <sup>25</sup>se <sup>26</sup>stali <sup>27</sup>divadlem <sup>28</sup>světu, <sup>29</sup>anđlům <sup>30</sup>i <sup>31</sup>lidem. <sup>10</sup>My <sup>11</sup>jme <sup>12</sup>blázní <sup>13</sup>kvůli <sup>14</sup>Kristu, <sup>15</sup>ale <sup>16</sup>vy <sup>17</sup>jste <sup>18</sup>rozumní <sup>19</sup>v <sup>20</sup>Kristu; <sup>10</sup>my <sup>11</sup>jme <sup>12</sup>slabí, <sup>13</sup>vy <sup>14</sup>vsak <sup>15</sup>silní; <sup>16</sup>vy <sup>17</sup>slavní, <sup>18</sup>my <sup>19</sup>“beze <sup>20</sup>čti. <sup>11</sup>Až <sup>12</sup>do <sup>13</sup>této <sup>14</sup>hodiny <sup>15</sup>hladovíme <sup>16</sup>a <sup>17</sup>žízňame, <sup>18</sup>jsme <sup>19</sup>šora <sup>20</sup>naží, <sup>21</sup>“dostáváme <sup>22</sup>řány, <sup>23</sup>jsme <sup>24</sup>“bez <sup>25</sup>domova, <sup>12</sup>opotomne <sup>13</sup>je <sup>14</sup>práci <sup>15</sup>vlastních <sup>16</sup>rukou. <sup>17</sup>Když <sup>18</sup>“nám <sup>19</sup>špilají, <sup>20</sup>žehnáme, <sup>21</sup>“když <sup>22</sup>“nás <sup>23</sup>pronásledují, <sup>24</sup>snášáme <sup>25</sup>to, <sup>26</sup>“když <sup>27</sup>“jsme <sup>28</sup>“haněni, <sup>29</sup>domlouváme. <sup>12</sup>Stali <sup>13</sup>jsme <sup>14</sup>se <sup>15</sup>ja-</p> <p><sup>1</sup>ř.: Ten, který; <sup>2</sup>ř.: člověk; <sup>3</sup>ř.: den (tzn. zasedání tribunálu; srv. 3.13: zde den určený lidmi); <sup>4</sup>ř.: té odlišuje (od ostatních);</p> <p><sup>11</sup>Ju 5,13 <sup>12</sup>2.94,11</p>
--	--

FIGURE 4. A problematic case: a widow line at the right page.

Of course this approach is very time-consuming and usually works only for short documents. Applying this approach to a book of a few hundred pages is not trivial and requires some tricks to speed up the process.

## 6 A semi-automatic approach

The way to tune page breaking in the previous paragraph can be improved in two places:

- problem detection: checking the result by eye is very tiresome and it is easy to overlook problems. It would be much better if we get a report of problems, saying for example *there is a widow line on page 317*.
- adjustment of paragraph formatting: we need an easy way to find which paragraph can be tweaked, i. e. to be formatted longer or shorter.

The workflow looks as in Figure 6.

To be able to refer to a paragraph, we number all paragraphs and typeset the number of each paragraph during tuning phase. In the final run they are left out.

### 6.1 Automatic problem detection

The method to detect problematic cases is to use a trick called absolute positioning. pdfT<sub>E</sub>X allows recording the absolute position of a point using the primitive

<p>1. KORINTSKÝM 2, 16</p> <p>316</p> <p><i>Projev nezralosti</i></p> <p><b>3</b> <sup>1</sup>A já, bratři, jsem k vám nemohl mluvit jako k duchovním, ale jako k <sup>2</sup>tělesným, jako k nedospělým v Kristu. <sup>3</sup>Dal jsem vám <sup>4</sup>pit mléko, pokrm <i>jem vám nedat</i>, <i>neboť ten</i> jste ještě nemohli <i>snést</i>. Ale ani teď ještě nemůžete, <sup>5</sup>neboť jste <i>stále</i> ještě <sup>6</sup>tělesní. Vždyť pokud je mezi vámi žárlivost, rozbroj <sup>7</sup>a rozdělení, zdali nejste <sup>8</sup>tělesní a <sup>9</sup>nežijete jen po lidsku<sup>10</sup>? <sup>4</sup>Když jeden říká: Já jsem Pavlův <sup>11</sup>a druhý: Já Apollův<sup>12</sup>, nejste <sup>13</sup>jako jiní lidé<sup>12</sup>? <sup>5</sup>Kdo je Apollós? A kdo je Pavel? Služebníci, skrze něž jste uvěřili, jak každému dal Pán. <sup>6</sup>Já jsem zasadil, Apollós zalil, ale <sup>7</sup>Bůh <sup>8</sup>dával <sup>9</sup>růst. <sup>7</sup>Tedy ani ten, kdo sází, ani ten, kdo zalévá, nejsou <sup>10</sup>něčím <i>zvláštním</i>, ale Bůh, který dává <sup>11</sup>růst. <sup>8</sup>Ten, kdo sází i ten, kdo zalévá, jsou jedno, každý však <sup>12</sup>dostane <sup>13</sup> svou vlastní mzdu podle <sup>14</sup>své<sup>13</sup> námahy. <sup>9</sup>Neboť jsme Boží spolupracovníci; <sup>10</sup>vy jste Boží pole, Boží stavba. <sup>11</sup><sup>10</sup>Podle <sup>12</sup>Boží<sup>11</sup> milosti, která mi byla dána, jsem jako</p> <p>moudrý stavitel položil základy a jiný <sup>12</sup>na <sup>13</sup>něm <sup>14</sup>stavi. Každý <sup>15</sup>at <sup>16</sup>si <sup>17</sup>dává <sup>18</sup>pozor, jak <sup>19</sup>na <sup>20</sup>něm <sup>21</sup>stavi. <sup>11</sup>Neboť nikdo nemůže položit jiný základ nežli ten, který je již položen, <sup>12</sup>a tím je Ježíš Kristus. <sup>13</sup>Jestliže někdo <sup>14</sup>na <sup>15</sup>tomto <sup>16</sup>základě <sup>17</sup>staví <sup>18</sup>ze <sup>19</sup>zlata, <sup>20</sup>stříbra, <sup>21</sup>drahých kamenů, dřeva, <sup>22</sup>seny <sup>23</sup>nebo <sup>24</sup>slámy, <sup>25</sup>dílo každého se stane <sup>26</sup>zjevným. Ten den <sup>27</sup>je <sup>28</sup>ukáze, <sup>29</sup>neboť se <sup>30</sup>zjeví<sup>28</sup> v ohni, a oheň vyzkouší <sup>31</sup>dílo každého, jaké je. <sup>14</sup>Jestliže někdo <sup>15</sup>na <sup>16</sup>tomto <sup>17</sup>základě <sup>18</sup>vystaví<sup>15</sup> dílo a ono <sup>19</sup>mu <sup>20</sup>zůstane<sup>16</sup>, dostane <sup>21</sup>odměnu, <sup>15</sup>jestliže <sup>16</sup>mu <sup>17</sup>jeho <sup>18</sup>dílo <sup>19</sup>shoří, <sup>20</sup>utrpí <sup>21</sup>škodu; sám se sice zachrání, ale jako skrze oheň. <sup>16</sup>Nevíte, že jste Boží svatyně a že <sup>17</sup>ve<sup>17</sup> vás <sup>18</sup>bydlí <sup>19</sup>„Duch <sup>20</sup>Boží“? <sup>17</sup>Ničili někdo <sup>18</sup>Boží <sup>19</sup>svatyni, <sup>20</sup>zničí <sup>21</sup>„Bůh jej. Neboť <sup>22</sup>Boží <sup>23</sup>svatyně je <sup>24</sup>svatá, <sup>25</sup>a <sup>26</sup>ta <sup>27</sup>svatyně<sup>18</sup> je <sup>28</sup>vy<sup>18</sup>. <sup>18</sup>At <sup>19</sup>nikdo <sup>20</sup>neklame <sup>21</sup>sám sebe. Jestliže si někdo mezi vámi myslí, <sup>22</sup>že <sup>23</sup>je <sup>24</sup>v <sup>25</sup>tomto <sup>26</sup>věku<sup>19</sup> moudrý, <sup>27</sup>at <sup>28</sup>se <sup>29</sup>stane <sup>30</sup>bláznem, aby se <sup>31</sup>stal <sup>32</sup>moudrým. <sup>20</sup>Neboť <sup>21</sup>moudrost <sup>22</sup>tohoto <sup>23</sup>světa <sup>24</sup>je <sup>25</sup>před <sup>26</sup>„Bohem <sup>27</sup>bláznovstvím. Vždyť je <sup>28</sup>na-</p> <p><sup>11</sup>ř: nechodíte podle člověka; <sup>12</sup>var.: tělesní; <sup>13</sup>ř: vlastní; <sup>14</sup>ř: se odhaluje; <sup>15</sup>ř: přistavěl (na ten základ); podobně v. 10 a 12; <sup>16</sup>var.: zůstává; <sup>17</sup>n.: mezi vámi; <sup>18</sup>ř: kterou; <sup>19</sup>ř: zrn. podle měřítka tohoto světa;</p> <p><sup>10</sup>ř 28,16</p>	<p>317</p> <p>1. KORINTSKÝM 4, 13</p> <p>psáno: „On<sup>11</sup> chytá moudré v jejich chytráctví.“<sup>10</sup> <sup>20</sup>A opět: „Pán zná myšlenky moudrých a ví, že jsou marné.“<sup>22</sup> <sup>21</sup>A tak <sup>22</sup>at <sup>23</sup>se <sup>24</sup>nikdo <sup>25</sup>nechlubí <sup>26</sup>lidmi. Vždyť všechno je vaše: <sup>22</sup>at <sup>23</sup>Pavel <sup>24</sup>nebo <sup>25</sup>Apollós <sup>26</sup>nebo <sup>27</sup>Kéfas, <sup>28</sup>at <sup>29</sup>svět <sup>30</sup>nebo <sup>31</sup>život <sup>32</sup>nebo <sup>33</sup>smrt, <sup>34</sup>at <sup>35</sup>o <sup>36</sup>věci <sup>37</sup>přítomné <sup>38</sup>nebo <sup>39</sup>budoucí – všechno je vaše, <sup>22</sup>vy <sup>23</sup>pak <sup>24</sup>jste <sup>25</sup>Kristovi <sup>26</sup>a <sup>27</sup>Kristus <sup>28</sup>Boží.<sup>10</sup></p> <p><i>Služebníci Kristovi</i></p> <p><b>4</b> <sup>1</sup>At <sup>2</sup>o <sup>3</sup>nás <sup>4</sup>každý<sup>2</sup> smýšlí jako <sup>5</sup>o <sup>6</sup>služebnících <sup>7</sup>tajemství. <sup>2</sup>A <sup>3</sup>od <sup>4</sup>správců <sup>5</sup>se <sup>6</sup>nakonec <sup>7</sup>vyzaduje, <sup>8</sup>aby <sup>9</sup>byl <sup>10</sup>každý <sup>11</sup>shledán <sup>12</sup>věrným. <sup>3</sup>Pro <sup>4</sup>mne <sup>5</sup>je <sup>6</sup>„pramálo <sup>7</sup>důležité, <sup>8</sup>zda <sup>9</sup>mě <sup>10</sup>posuzujete <sup>11</sup>vy <sup>12</sup>nebo <sup>13</sup>lidský <sup>14</sup>soud<sup>3</sup>. Ale <sup>15</sup>ani <sup>16</sup>já <sup>17</sup>sám <sup>18</sup>sebe <sup>19</sup>neposuzuji. <sup>4</sup>Ničeho, <sup>5</sup>co <sup>6</sup>by <sup>7</sup>svědčilo <sup>8</sup>proti <sup>9</sup>mně, <sup>10</sup>ale <sup>11</sup>ten <sup>12</sup>nejsem <sup>13</sup>ospravedlněn; <sup>14</sup>ten, <sup>15</sup>kdo <sup>16</sup>mne <sup>17</sup>posuzuje, <sup>18</sup>je <sup>19</sup>Pán. <sup>5</sup>Proto <sup>6</sup>nice <sup>7</sup>nesudte <sup>8</sup>předčasně, <sup>9</sup>dokud <sup>10</sup>nepřijde <sup>11</sup>Pán, <sup>12</sup> který <sup>13</sup>osvítí <sup>14</sup>všech <sup>15</sup>skryté <sup>16</sup>ve <sup>17</sup>tmě <sup>18</sup>a <sup>19</sup>zjeví <sup>20</sup>úmysly <sup>21</sup>srdce. <sup>6</sup>A <sup>7</sup>tehdy <sup>8</sup>se <sup>9</sup>každému <sup>10</sup>dostane <sup>11</sup>pochovaly <sup>12</sup>„Boha.“<sup>11</sup></p> <p><sup>11</sup>ř: Ten, který; <sup>12</sup>ř: člověk; <sup>13</sup>ř: den (tzn. zasedání tribunálu; srv. 3:13; zde den určený lidmi); <sup>14</sup>ř: t: odlišuje (od ostatních); <sup>15</sup>ř: může zn. ty, kteří nesou nečistotu světa (Lv 16,10);</p> <p><sup>10</sup>Jn 5,19 <sup>22</sup>94,11</p>
--	---

FIGURE 5. Solving the problem by making a paragraph shorter.

\pdfsavepos and its friends. To see how it can help to find problematic cases, let's have a close look at widow detection. A widow line is the last line of a paragraph that ends up at the beginning of a page or column. We detect widow lines as follows:

1. For each paragraph, record the absolute position of the last line of the paragraph. It can be done by appending the primitive \pdfsavepos to all relevant paragraphs.
2. For each page, record the absolute position where the page body starts.
3. For each page, compare the absolute position of the page start against the absolute position of the last line of each paragraph. If they are close enough, then a widow line is found.

In the example in the previous section, we have the following lines recorded:

```
\bodytopos{317}{28146752}
\endparpos{11}{317}{27423153}
```

The first line was generated for the body top of page 317, the second line was generated for the last line of the paragraph number 11 on page 317. The last numbers in those commands are the absolute positions (in scaled point) of corresponding elements. Since they are close enough (28146752 – 27423153 = 723599 = 11pt), this indicates that there is a widow line from paragraph 11 on page 317.

Other problematic cases are detected in a similar manner.

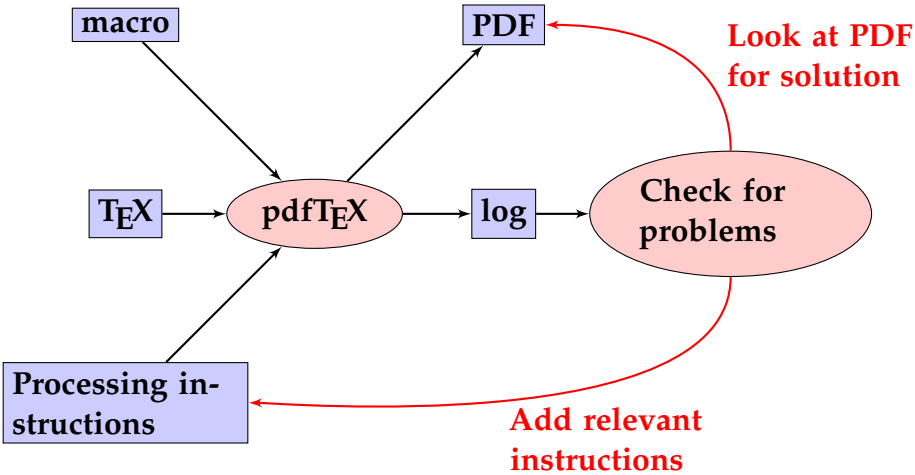


FIGURE 6. The process to tune page breaking

## 6.2 Adjustment of paragraph formatting

Once a problematic case is detected, we need to tweak some paragraph to solve the problem. Let's continue with the previous example. To make the widow line disappear, we have two choices:

- move the widow line to the left column, i. e. to “shrink” one of the previous paragraph by one line; or
- move another line from the left column to the right column, so that there are 2 lines of the paragraph 11 on the right column. In this case, we need to “stretch” one of the previous paragraphs.

The question here is to find which paragraph can be “shrunk” or “stretched”. Since we already number all paragraphs, this turned out to be simple: together with the paragraph number, we also show whether the paragraph can be “stretched” or “shrunk”. If a paragraph can be stretched, a sign +1 is shown below the paragraph number. If it can be shrunk, -1 is shown too. For example, the paragraph containing the widow line in our example is paragraph number 11, and it can be stretched (+1 is shown) but cannot be shrunk (-1 is not shown). Usually most of paragraphs can be stretched, but only a few can be shrunk.

Having all these numbers displayed, finding a suitable paragraph to tweak is quite easy. In our example, we see that paragraph number 9 can be shrunk, so we add a command saying “shrink the paragraph number 9 by one line” and rerun pdfT<sub>E</sub>X. The result is shown in Figure 5.

This tuning process must be done as the last step, since any change in the data might require retuning page breaking.



## 7 Use of micro-typography

### 7.1 Margin kerning

Margin kerning is heavily used in this project. This is essential to make the margins look smooth, since too often it happens that non-letter elements ends up at the beginning or the end of a line. Those elements include:

- hyphen (very frequent, due to narrow columns),
- punctuations,
- footnote marks and references,
- verse number,
- various text marks.

How to apply margin kerning to hyphen, punctuations and other characters is not difficult and is described well in pdfTeX manual, microtype manual and various articles. However it is not clear how to apply margin kerning to elements that are not a character from the main font, so it is worthwhile to take a close look at such an example.

Suppose we want to apply margin kerning to the verse numbers. Those numbers often end up at the left margin. Since those verse numbers are typeset as superscript at a smaller size than the body, their shape actually take less “ink” than a regular character. So when they end up at the left margin, they will cause the visual effect that the actual line looks slightly shorter than other lines. If we leave out margin kerning for them, the left margin would look a bit uneven. However, we cannot apply margin kerning to those verse numbers in the same way as for hyphen or punctuations, since a verse number is typeset as a (raised) hbox, and margin kerning applies only to characters.

We make use of the fact that the amount of marginal kerning for all verse numbers are approximately the same. So, our trick is to insert before each verse number a “virtual” character. Such a character has zero dimensions and no “shape”. This can be done using a special virtual font, where each character definition is empty. Since the character is “empty”, it is not visible but we can apply margin kerning to such a character and therefore to the element after it, namely the verse number in this case. Other superscript elements are handled in a similar manner.

About 50% of lines has right margin kerning and 20% of lines has left margin kerning.

### 7.2 Font expansion

Font expansion is even more important than margin kerning for this task. Since we rely mainly on tweaking `\looseness` to stretch or shrink a paragraph, it is important that there are as many stretchable/shrinkable paragraphs as possible. Font expansion gives us exactly what is needed in this case: it gives more room to line-breaking, therefore makes a paragraph more likely stretchable or shrinkable.

Also, when a paragraph is stretched or shrunk without font expansion, usually it looks very ugly due to loose spacing between words. Font expansion is critical to compensate this: it makes a stretched or shrunk paragraph look reasonable.

Font expansion was used at the limit  $\pm 2\%$ , using `autoexpand` and without tuning expansion factor of individual characters. Going further than this limit is not recommended, since the effect of font expansion might be visible in some cases.

## 8 The result

The job took about 3 days to tune page breaking for about 500 pages. There are only 2 places in the book where there is an orphan line in the left column, since it is impossible to avoid them without breaking other rules. The total number of processing instructions (to stretch or shrink a paragraph) is about 250. About 15 paragraphs were shrunk, the rest was stretched. Almost all requirements are fulfilled, and the publisher was satisfied with the result.

## 9 What is still imperfect

Despite our effort, the typesetting was not perfect. We will take a quick look at issues that could be improved.

*Hyphen at page breaks:* it was not required by the publisher to prohibit hyphenation at page breaks, however it is better if we could do so. It is possible to detect such a problem automatically. Unfortunately it is not easy to fix such problem without changing the text at all. Due to lack of time, this issue was left unsolved.

*Subtitles break at undesired places:* Some subtitles are longer than one line and hence must be broken into 2 lines. However, sometimes they are broken at undesired place. This cannot be auto-detected, and cannot be fixed without understanding the actual subtitle. So, this issue must be fixed in the primary data by the publisher.

*Last line in paragraph too short:* we did put no restrictions on the length of the last line in a paragraph. Due to frequent stretching/shrinking paragraphs, there are many cases when the last line of a paragraph contains a single word (or even worse, a fragment from a hyphenated word). One could argue that this is tolerable for narrow columns, but it is better if we could have avoided it.

## 10 Conclusions

- T<sub>E</sub>X is still a very powerful typesetting system if properly used. It is not suitable for every application, however in certain domains it is really hard to beat T<sub>E</sub>X, despite the fact that it is a 30-years old program.
- The micro-typographic extensions of pdfT<sub>E</sub>X can be very helpful in extreme cases, where we need to give a little extra flexibility to paragraph formatting.
- High-quality typesetting requires human decision and intervention.