Jost Gippert

# Linking methods as a basis for cross-linguistic text retrieval: Problems and solutions

1.   Text segmentation principles:

1.1. Internal vs. external principles of segmentation of texts:

1.1.1.   Purely linguistic segmentation (sequence of sentences, words within sentences) vs.

1.1.2.   content-based segmentation (e.g., books, chapters, paragraphs, strophes, verses) vs.

1.1.3.   segmentation based on "surface" representation (pages, lines of a given edition).

1.2. Superiority of content-based segmentation with respect to larger texts:

1.2.1.   More easy referencing by not using too large numbers (e.g., sentence no. 25387);

1.2.2.   does not depend on a given printed edition the arrangement of which is always more or less "accidental";

1.2.3.   is the only reasonable basis for a cross-textual retrieval (see below).


2.   Problems of establishing a solid content-based segmentation for Buddhist texts:

2.1. Example 1: The Pāli Vinayapiṭaka

2.1.1.   Electronic version (public domain) produced by the "Sri Lanka Buddha Jayanti Tripitaka Series" (version 15.1.97) contains referencing to two printed editions, the BJT edition and the PTS edition, but for the latter, only book numbers and pages are indicated. There is no referencing to the content-based segmentation used traditionally (cp. Box 1: this would be Vin. I, 6, 7-10).

2.1.2.   Electronic version produced by Mahidol University Computing Center (Bangkok, Thailand, 1.8.1997: "BUDSIR IV") contains referencing to one printed edition, and for this too, only book numbers and pages are indicated. There is no referencing to the content-based segmentation used traditionally (cp. Fig. 1).
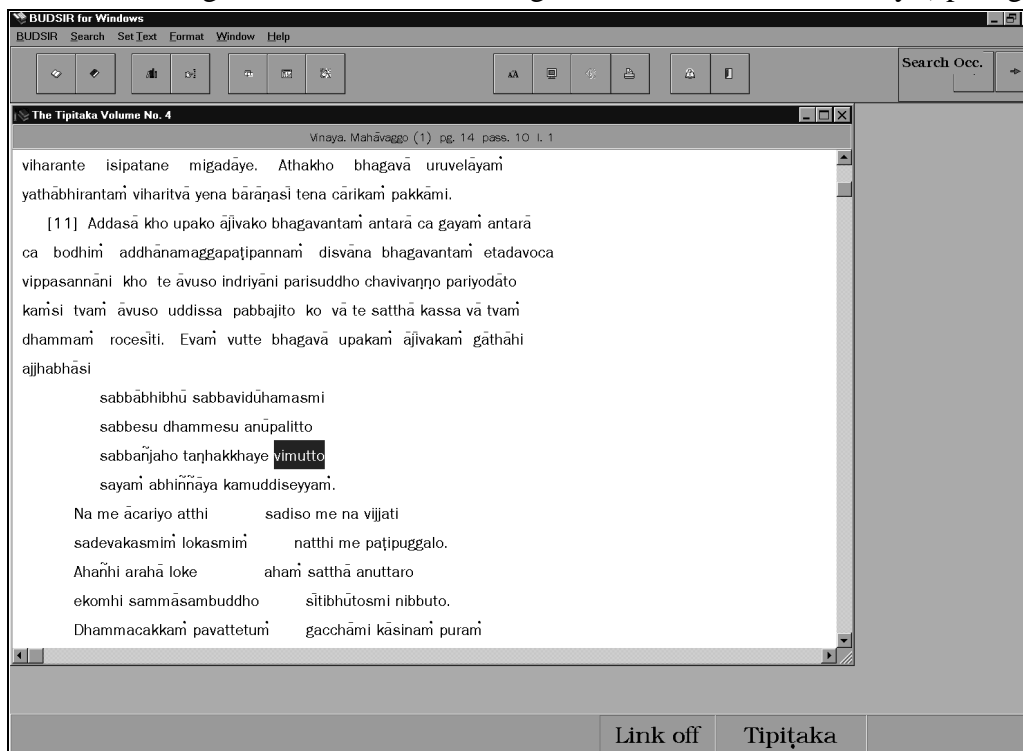


**Fig. 1**:                                   Mahidol version of Vin. I, 6, 7-8 (screen output)

<table>
<tr><td>

[BJT Page 016.]

4. Addasā kho upako ājivako1- bhagavantaṃ antarā ca gayaṃ antarā ca bodhiṃ aññamagga-paṭipannaṃ. Disvāna bhagavantaṃ etadavoca: "vippasannāni kho te āvuso indriyāni. Parisuddho chavivaṇṇo pariyodāno. Kaṃ si tvaṃ āvuso uddissa pabbajito? Ko vā te satthā? Kassa vā tvaṃ dhammaṃ rovesī?"Ti. Evaṃ vutte bhagavā apakaṃ ājivakaṃ gāthāhi ajjhabhāsi: -

"Sabbābhibhu sabbavidu'hamasmi
Sabbesu dhammesu anupalitto,
Sabbañajabho taṇhakkhaye vimutto
Sayaṃ abhiññāya kamuddiseyyaṃ.
Na me ācariyo atthi sadiso me na vijjati,
Sadevakasmiṃ lokasmiṃ natthi me paṭipuggalo.
Ahaṃ hi arahā loke ahaṃ satthā anuttaro,
Eko'mhi sammāsambuddho sītibhutosmi nibbuto.
Dhammacakkaṃ pavattetuṃ gacchami kāsinaṃ paraṃ,
Andhabhūtasmiṃ lokasmiṃ āhañachaṃ2-amatadundubhi"ti.
Yathā kho tvaṃ āvuso paṭijānāsi, arahasi anantajinoti.
"Mā disā ve jinā honti ye pattā āsavakkhayaṃ,
Jinā me pāpakā dhammā tasmā'haṃ upakā jino"ti.

5. Evaṃ vutte upako ājivako ājivako "huveyyapāvuso"ti3- vatvā sīsaṃ okampetvā ummaggaṃ gahetvā pakkāmi.

6. Atha kho bhagavā anupubbena cārikaṃ caramāno yena bārāṇasī isipatanaṃ migadāyo , yena pañcamaggiyā bhikkhu, tenupasaṃkami. Addasaṃsu. Kho pañcavaggiyā bhikkhu bhagavantaṃ durato'va āgacchantaṃ. Disvāna aññamaññaṃ saṇṭhapesuṃ: "ayaṃ āvuso samaṇo gotamo āgacchati bāhuliko [PTS Page 009] paṭhānavibbhanto āvatto bāhullāya. So neva abhivādetabbo. Na paccuṭhātabbo. Tassa pattacīvaraṃ paṭiggahetabbaṃ api ca kho āsanaṃ ṭhapetabbaṃ, sace ākaṅkhissati, nisīdissati"ti.

1. "Xjiviko ma. Nu. Pu;a. Ma vi; 2. "Ahañachiṃ āhañachuṃ ityapi
3. "Huveyyāvuso' - katthavi

Box **1**:   BJT version of Vin. I, 6, 7-10 (raw text)
</td><td>

|b16.    [BJT Page 016.]
|c4. |xVin_I_6,_7
|s1   Addasā kʰo upako ājivako1- bʰagavantaṃ antarā ca gayaṃ antarā ca bodʰiṃ añña-maggapaṭipannaṃ. Disvāna bʰagavantaṃ etadavoca: "vippasannāni kʰo te āvuso indriyāni. Parisuddʰo cʰavivaṇṇo pariyodāno. Kaṃ si tvaṃ āvuso uddissa pabbajito? Ko vā te sattʰā? Kassa vā tvaṃ dʰammaṃ rovesī?"Ti. |xVin_I_6,_8 Evaṃ vutte bʰagavā apakaṃ ājivakaṃ gātʰāhi ajjʰabʰāsi: -
|s2   "Sabbābʰibʰu sabbavidu'hamasmi
|s3   Sabbesu dʰammesu anupalitto,
|s4   Sabbañajabʰo taṇhakkʰaye vimutto
|s5   Sayaṃ abʰiññāya kamuddiseyyaṃ.
|s6   Na me ācariyo attʰi sadiso me na vijjati,
|s7   Sadevakasmiṃ lokasmiṃ nattʰi me paṭipuggalo.
|s8   Ahaṃ hi arahā loke ahaṃ sattʰā anuttaro,
|s9   Eko'mhi sammāsambuddʰo sītibʰutosmi nibbuto.
|s10  Dhammacakkaṃ pavattetuṃ gaccʰami kāsinaṃ paraṃ,
|s11  Andʰabʰūtasmiṃ lokasmiṃ āhañacʰaṃ2-amatadundubʰi"ti.
|xVin_I_6,_9
|s12  Yatʰā kʰo tvaṃ āvuso paṭijānāsi, arahasi anantajinoti.
|s13  "Mā disā ve jinā honti ye pattā āsavakkʰayaṃ,
|s14  Jinā me pāpakā dʰammā tasmā'haṃ upakā jino"ti.
|c5. |s1 Evaṃ vutte upako ājivako ājivako "huveyyapāvuso"ti3- vatvā sīsaṃ okampetvā ummaggaṃ gahetvā pakkāmi.
|xVin_I_6,_10 |c6.
|s1   Atʰa kʰo bʰagavā anupubbena cārikaṃ caramāno yena bārāṇasī isipatanaṃ migadāyo, yena pañcamaggiyā bʰikkʰu, tenupasaṃkami. Addasaṃsu. Kho pañcavaggiyā bʰikkʰu bʰagavantaṃ durato'va āgaccʰantaṃ. Disvāna aññamaññaṃ santʰapesuṃ: "ayaṃ āvuso samaṇo gotamo āgaccʰati bāhuliko |p9 [PTS Page 009] patʰānavibbʰanto āvatto bāhullāya. So neva abʰivādetabbo. Na paccutʰātabbo. Tassa pattacīvaraṃ paṭiggahetabbaṃ api ca kʰo āsanaṃ tʰapetabbaṃ, sace ākaṅkʰissati, nisīdissati"ti.

Box **2**: Same, after entering of additional references
</td></tr>
</table>

2.1.3. For easy cross-referencing, the greatest amount of information possible should be envisaged, at least with respect to the traditionally used PTS edition (cp. Box 2).

2.2. Example 2: The Udānavarga (Collection of Sanskrit strophes of verses, mostly corresponding to strophes of the Pāli Dhammapāda):

2.2.1. In spite of relatively simple structure (complete collection consisting of 33 Vargas with a maximum of 87 strophes), great divergencies exist between several editions as to the numbering of strophes (e.g.: Uv. 29,24 in the edition by F. Bernhard, Udānavarga, Göttingen 1965 corresponds to Uv. 29,34 in the edition by R. Pischel, Die Turfan-Recensionen des Dhammapada, Berlin 1908 [and 29,23 in the translation of the Tibetan version by W.W. Rockhill, London 1892]).

2.2.2. This may be due to the bad state of preservation of the text, verses and strophes missing in manuscripts, or to secondary additions (cf. Bernhard, o.c., 14).

2.2.3. Arrangement and segmentation of the text do not agree at all with the Pāli Dhammapada so that the segmentation of this text cannot be adopted as it is (all the more since for the Pāli Dhammapada itself, two divergent segmentations are used, viz. one counting only strophes (from 1 to 423), and one dividing the text into 26 vargas with a differing amount of strophes; e.g., Uv. 21 [Tathāgatavarga], 1 corresponds to DP 353 ≈ 24 [Taṅhavagga], 20).

2.2.4. For the sake of cross-referencing, a common segmentation should be envisaged.

3. Problems of establishing a system for cross referencing between related texts

3.1. Example 1: The Gāndhārī Dharmapada (Collection of Prakrit strophes of verses, mostly corresponding to strophes of the Pāli Dhammapāda and the Sanskrit Udānavarga; edition by J. Brough, London 1962):

3.1.1. Although most of the strophes contained in the text are exact equivalents of strophes as present in the Pāli Dhammapāda and/or the Sanskrit Udānavarga, the arrangement is completely different again (e.g., GDP 1 [Brammaṇa], 1ab corresponds to DP 393ab ≈ 26 [Brāhmaṇavagga], 11ab, Uv. 33 [Brāhmaṇavarga], 8ab (Sanskrit version ed. Bernhard) / 33, 11ab (Tibetan version). Cp. Fig. 2 showing the concordance of DP and GDP as present in Brough's edition.

**CONCORDANCE II**

PALI DHAMMAPADA

| | | | | | |
|---|---|---|---|---|---|
| i. *Yamaka* | | 37 | 137a | 74 | |
| 1 | 201 | 38 | 137c | 75 | |
| 2 | 202 | 39 | 137d | | |
| 3 | | 40 | 138b | vi. *Paṇḍita* | |
| 4 | | 41 | 153 | 76 | 231 |
| 5 | | 42 | | 77 | 230 |
| 6 | | 43 | | 78 | |
| 7 | 217 | | | 79 | 224 |
| 8 | 218 | iv. *Puppha* | | 80 | |
| 9 | 192 | 44 | 301 | 81 | 239 |
| 10 | 193 | 45 | 302 | 82 | 225 |
| 11 | 213 | 46 | 300 | 83 | 226 |
| 12 | 214 | 47 | 294 | 84 | 324 |
| 13 | 219 | 48 | (294) | 85 | |
| 14 | 220 | 49 | 292 | 86 | |
| 15 | 205 | 50 | 271 | 87 | |
| 16 | 206 | 51 | 290 | 88 | |
| 17 | 203 | 52 | 291 | 89 | |
| 18 | 204 | 53 | 293 | | |
| 19 | 190 | 54 | 295 | vii. *Arahanta* | |
| 20 | 191 | 55 | 296 | 90 | |
| | | 56 | | 91 | |
| ii. *Appamāda* | | 57 | 297 | 92 | |
| 21 | 115 | 58 | 303 | 93 | |
| 22 | 116 | 59 | 304 | 94 | |
| 23 | (128) | | | 95 | |
| 24 | 112 | v. *Bāla* | | 96 | |
| 25 | 111 | 60 | | 97 | |
| 26 | 117 | 61 | | 98 | |
| 27 | 129, 130 | 62 | | 99 | |
| 28 | 119 | 63 | | | |
| 29 | 118 | 64 | 233 | viii. *Sahassa* | |
| 30 | 120 | 65 | 234 | 100 | 306 |
| 31 | 74 | 66 | | 101 | 308 |
| 32 | 73 | 67 | | 102 | 309 |
| | | 68 | | 103 | 305 |
| iii. *Citta* | | 69 | (283) | 104 | |
| 33 | 136 | 70 | 313 | 105 | |
| 34 | 137b | 71 | | 106 | 310 |
| 35 | | 72 | | 107 | 319, 320 |
| 36 | 138a | 73 | | 108 | 321 |

**Fig. 2**:          Concordance of DP and GDP

3.1.2.      A thorough content-based segmentation for the Gāndhārī Dharmapada is not easy to establish for the same reasons as with the Udānavarga.

3.1.3.      As the text can hardly be treated without permanent comparison of its Pāli and Sanskrit equivalents (it is the only larger Gāndhārī Prakrit text preserved at all), cross-referencing is especially important. Reliable segmentation of all three texts is presupposed.

3.2. Example 2: Tocharian (A) text relating the tale of the conversion of Upaka (Upage) by Buddha (obviously fragment from the Udānālaṃkāra, i.e., an Udānavarga commentary): THT 850 sq. / Toch. A (edition Sieg-Siegling), No. 217 sq.

3.2.1.      "Internal" segmentation as used in the printed edition: lines of manuscript pages (e.g., 217a 5) vs. content-based segmentation suggesting itself from verse structure which is indicated in the manuscript by numbers. Problem: The fragmentary status does not enable us to establish higher units (chapter 21, corresponding to the "Tathāgatavarga" of the Udānavarga as indicated in F. Bernhard's edition [p. 278]?)

No. 217 = T Ⅲ Š 79. 15

Der Länge nach ziemlich vollständiges, teilweise aber stark beschädigtes Blatt. Nach Herstellung der Photographie ist noch der sehr zerstörte Rest der linken Seite gefunden worden. Vgl. Tafel 29.

Vorderseite

1 (nicht erhalten)           **217 a**

2 //// kñ[ä]ññä [ta] – – – [p·] sne y·· y· wāryāñc · [śś ·] – – – ptā ////

3 //// ·· e skākā wärpāt͚ p(tā)ṉ̃ḵat͚ ārkiśoṣṣis krant͚ ṃarkaṃpal͚ āksis(s)i – –·– – – – – –
     [p]ū̱ḵ knānmāṃ ṯṃaṣ͚ bram poñcäṃ wältsa –

4 – – – sn · – · k(ā)ckeyo ✵ p̱aklyoṣas wrasañⁱ͚ pūk kācke parsācⁱ͚ pū̱ḵ͚ knānmāṃ – k · ṉ̃ḵat͚
     p̱arko parnont͚ māgat ṣiṃ ypeyaṃ ✵ wärpā

5 – ks[ï]ssi 3 krañcäṃ ṃarkampal͚ māryu praṣtaṃ okñäṣ ñāktas napenas sam̱͚ oṅkraci ✵ 8
     p̱alskāt͚ pū̱ḵ͚ knānmāṃ ke maltw āksisam͚ lyāklyäṃ k̂upā

3.2.2.      Manifold necessity of cross-referencing both with poetic texts (Udānavarga, Dhammapada, Gāndhārī Dharmapada) and prose texts (Sanskrit Catuṣpariṣatsūtra, Lalitavistara, Mahāvastu; Pāli Vinayapiṭaka and many more); cp. E. Sieg / W. Siegling, Festschrift M. Winternitz, Leipzig 1933, 167 sqq. or, vice versa, F. Bernhard's edition of the Udānavarga (l.c.) or E. Waldschmidt's edition of the Catuṣpariṣatsūtra (Vol. I, Berlin 1952). E.g., 217b 6 sq. (≈ UA 21, 12) corresponds to Uv. 21,1, DP 353, CPS 10,5, MV III, 326, 5-8, Vin. I, 6,8 and others.

4.       Cross-referencing and electronic retrieval

4.1.      Task: Simple indication of parallels as in many printed editions (cp. Fig. 4 showing Brough's edition of the GDP with parallels from Uv. and DP indicated and Fig. 5 showing cross references indicated in Bernhard's edition of Uv.) should be overcome by immediate automatic access to parallel texts when treated electronically.

**Fig. 4:** Indication of cross references in GDP edition



**Fig. 5**: Indication of cross references in Uv. edition

## 4.2. Simple solution: editing of related text passages side by side within one file (cp. Fig. 6 showing GDP 1, 2 with its equivalent, DP 394, arranged interlinearily)

### 4.2.1. Shortcoming: Loss of readability

### 4.2.2. Problem not to be covered easily by retrieval software: Separate treatment of several languages within one text (here: Gāndhārī Prakrit vs. Pāli)



**Fig. 6**:      GDP 1,2 contrasted with DP 394

## 4.3. Sophisticated solution: "Synchronizing of texts" (Wordcruncher solution; cf. http://www.wordcruncher.com)

### 4.3.1. Requirement: Common segmentation of the texts to be synchronized

#### 4.3.1.1. This is easy in, e.g., Bible tradition where the same segmentation of texts has been used traditionally (e.g. Mt. 6,9 [cp. Fig. 7 showing synchronous arrangement of the Armenian and Greek New Testament] or 2.Chr. 13,12) with but a few exceptions (esp. in OT: Jer.)

**Fig. 7**:　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　Mt. 6, 9 in Armenian and Greek

4.3.1.2.　　Synchronizing is not easy where no such common segmentation exists as in the case of the Dhammapāda equivalents, not to mention the prose texts.

4.3.2.　Actual solution: Inserting common references as highest order segmentation units in the texts. Cp. Fig. 8 showing GDP 1,2 synchronized with DP 394 and Uv. 33,6; Fig. 9 showing Toch. A 217b 6 [21,12] synchronized [in two arrangements] with DP 353 and Uv. 21,1; Fig. 10 showing GDP 1,6 synchronized with the Pāli prose text Saṃyuttanikāya I, 167; Fig. 11 showing Toch. A 217b 5 [21,12] synchronized [in two arrangements] with CPS 10,2.

4.4.　Task for the future: Extending the reference system to further texts of the Pāli Canon as well as to other branches of Buddhist tradition (Chinese, Tibetan, etc.).

**Book #1 -- Dharmapada (Gandhari-Prakrit)**
File  Search  View  Options  Window  Help
DP 394: DhP 1,

2

ki di jaḍa'i drumed$^b$a

{kiṃ te jaṭābi dummed$^b$a}

ki di ayina-śaḍi'a

{kiṃ te ajina-sātiyā}

adara gahana kitva

{abb$^b$antaraṃ te gahanaṃ}

bahire parimajasi. (2)

{bāhiraṃ parimajjasi.}

{O_3}

**Book #2 -- Dhammapada (Pali)**
File  Search  View  Options  Window  Help
DP 394: DhP 26, 394

394

kiṃ te jaṭābi dummed$^b$a kiṃ te ajinasāṭiyā
abb$^b$antaraṃ te gahanaṃ bāhiraṃ
parimajjasi \ 26\ 394

395

paṃsukūlad$^b$araṃ jantuṃ kisaṃ
d$^b$amanisant$^b$ataṃ
ekaṃ vanasmiṃ j$^b$āyantaṃ tam ahaṃ brūmi

**Book #3 -- Samyuttanikaya**
File  Search  View  Options  Window  Help
: , , : (n. ): ,

TITUS

Pālī Canon

Saṃyuttanikāya

edited by the
Maharagama Bhikkhu Training Centre

**Book #4 -- Udanavarga**
File  Search  View  Options  Window  Help
DP 394: Uv. 33,

DP_394

(kiṃ te jaṭāb$^b$ir durbudd$^b$e)
(kim cāpy ajinaśaṭib$^b$iḥ) /
(abb$^b$yantaraṃ te kaluṣam)
(bāhyakaṃ parimārjasi)[1] // 6A
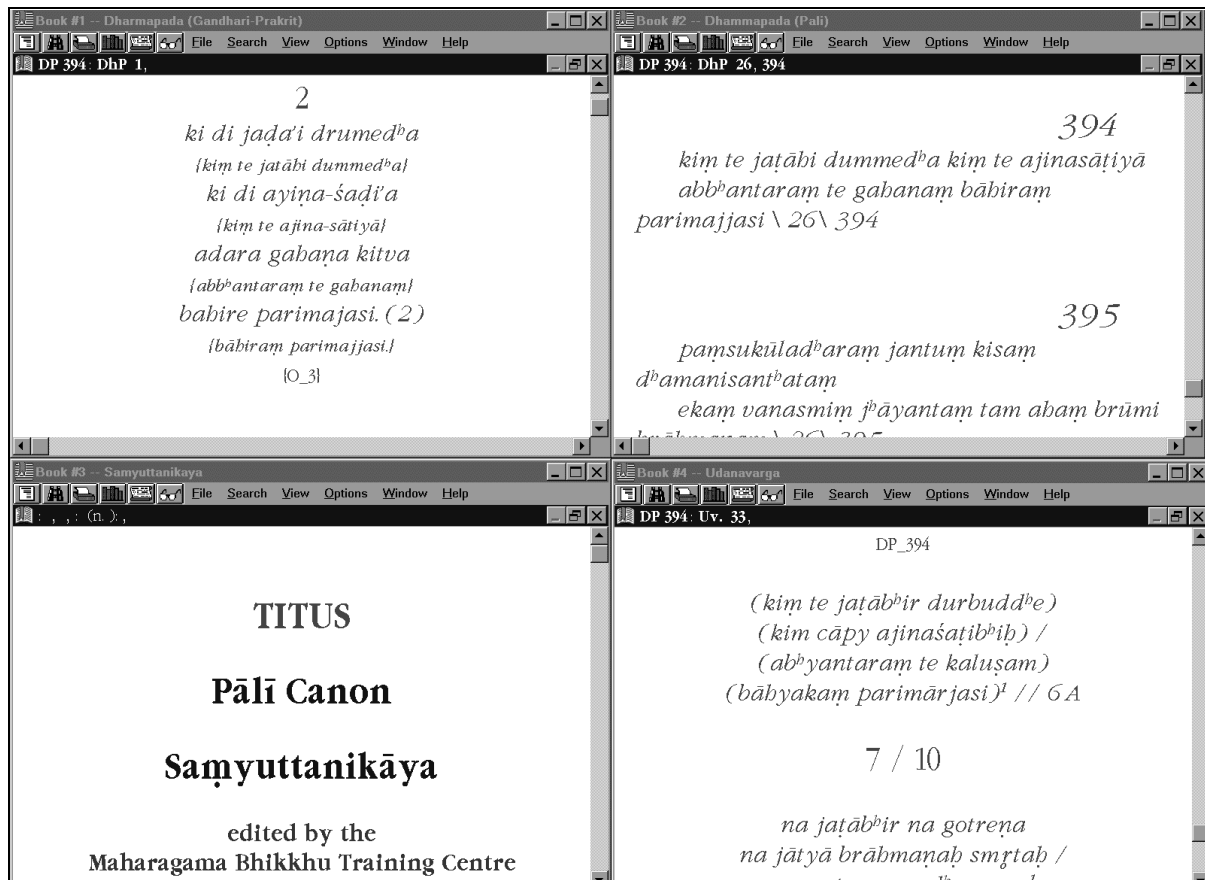
7 / 10

na jaṭāb$^b$ir na gotreṇa
na jātyā brāhmaṇaḥ smṛtaḥ /

**Fig. 8:**                        GDP 1, 2 synchronized with DP and Uv. 33,6

**Book #1 -- Tocharian A 217-218**
File  Search  View  Options  Window  Help
DP 353 THT 850, 217b, 6: 21,12

sabbab$^b$ib$^b$u sabbavudū 'ham asmi   sabbesu d$^b$ammesu
anūpalitto
evaṃ vutte b$^b$agavā[28] Upakaṃ ājīvikaṃ gāt$^b$ābi aj$^b$ab$^b$āsi[29] :
sabbāb$^b$ib$^b$ū[30] sabbavidu 'ham asmi[31]  sabbesu[32] d$^b$ammesu[34]
anupalitto[33]
ptāñKAT KAssī[28] we [29]pūK' ly[utā]r[30] . . -
- - - - nasaM' pūkyo[32] sne loplune[33]

**Book #2 -- Dhammapada (Pali)**
File  Search  View  Options  Window  Help
DP 353: DhP 24, 353

353

sabbāb$^b$ib$^b$ū sabbavudū 'ham asmi sabbesu
d$^b$ammesu anūpalitto
sabbañjaho taṇhakkhaye vimutto sayaṃ
ab$^b$iññāya kam uddiseyyaṃ \ 24\ 353

**Book #4 -- Tocharisch A**
File  Search  View  Options  Window  Help
DP 353: THT 850 217b,

DP_353     Uv_21,1
_ _ _ _ _ _ _ _ Ṣ P Arklune _ .k. ṣi tñi
kāruṇiK\ : ptāñKAT KAssī we pūK\
ly[utā]r.. _ _ _ _ _ nasaM\ pūkyo sne
loplune wraM Aṃny[o] (:)
* * * * * * (*) Ṣ^PA-rklu-+e * +k+ ṣi tñi kā-ru-ṇi-K\,
: ptā-ñKA-T^KA-ssi we pū-K\ ly+-(*) r+ * * * * *

**Book #5 -- Udanavarga**
File  Search  View  Options  Window  Help
DP 353: Uv. 21,

DP_353

sarvāb$^b$ib$^b$ūḥ sarvavid eva cāsmi
sarvaiś ca d$^b$armaiḥ satataṃ na liptaḥ /
sarvaṃjahaḥ sarvab$^b$ayād vimuktaḥ
svayaṃ hy ab$^b$ijñāya kam uddiśeyam // 1

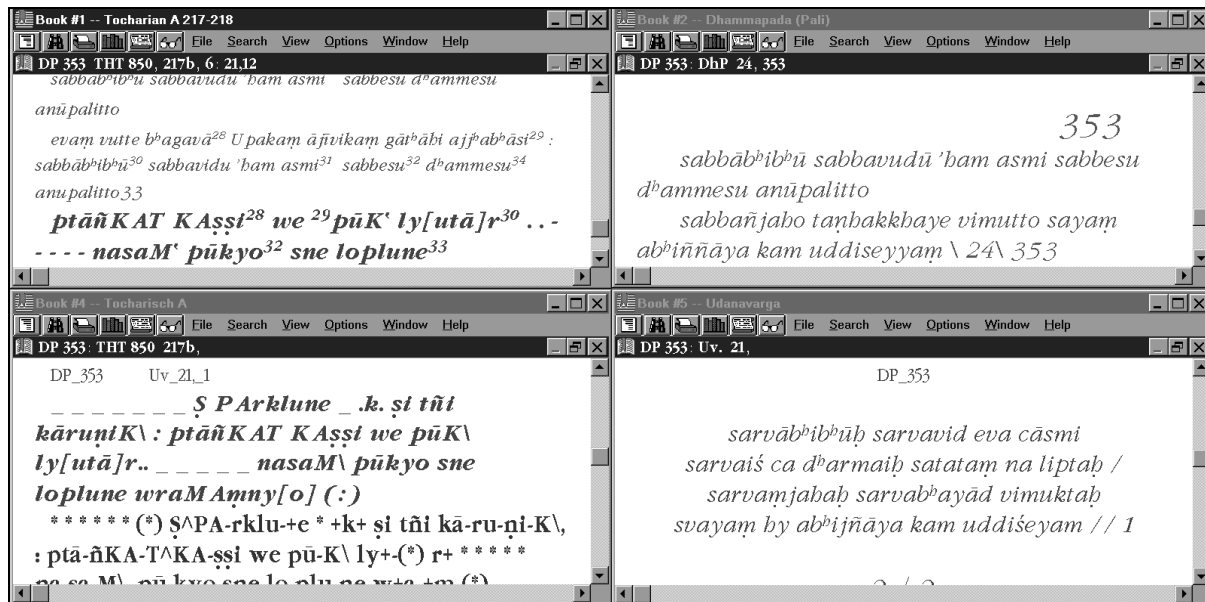**Fig. 9:**                        Toch. A 217b 6 [21,12] synchronized with DP 353 and Uv. 21,1

**Fig. 10**:                    GDP 1,6 with Saṃ. I, 167



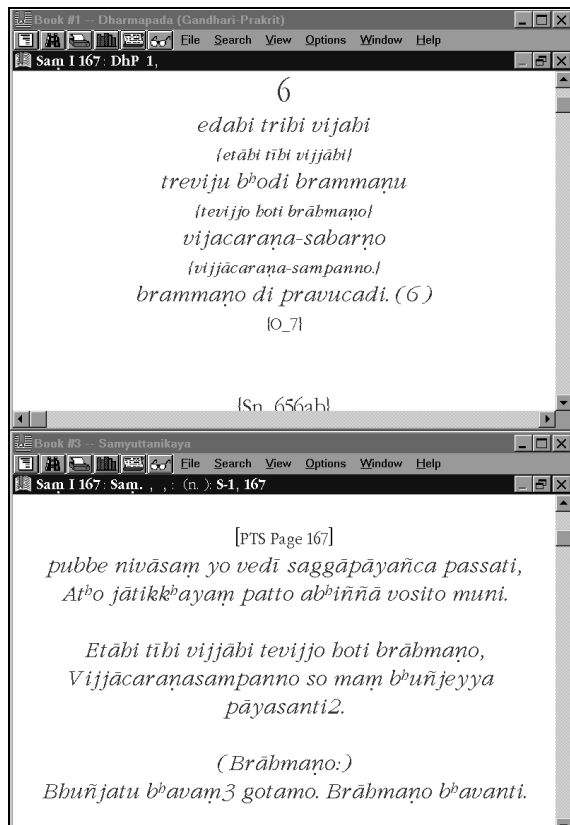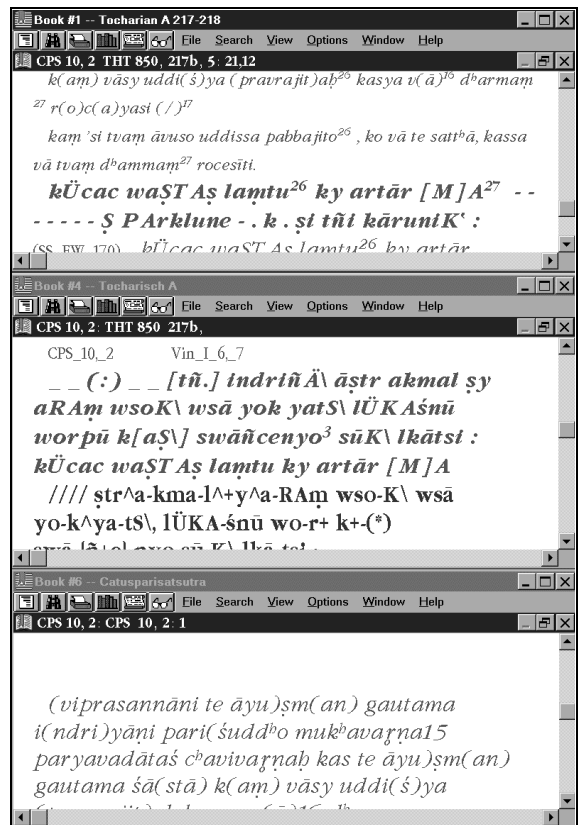**Fig. 11**:                    Toch. A 217b 5 with CPS 10,2