

UCCA: Hebrew annotation and tokenization guidelines

1. Annotation guidelines

1. Double genitive construction:

- Beit_C o_F [shel_S Dani_A (Beit)_A]_E (בית-ו של דני)
- [Dmut__C o_F]_E [shel_R nachash_C]_C (דמות-ו של נחש)
- Sovlanut_S o_F [shel_R Dani_C]_A azla_D (סובלנות-ו של דני אצל)

2. Static scenes

- [zo_A [ha_F machberet_C [shel_S Yael_A
(machberet)_A]_E]_A (**IMP**)**S**]_H
Zero copula
- Haita_S zo_A [tmunat nachash]_A (היתה זו תמונה נחש)
- Ma_A ze_S [ha_F davar_C [ha_F ze_C]_E]_A? מה זה הדבר (זה)
- 'Et_A ze_S [kli ktiva]_A (עט זה כלי כתיבה)
- [Ein_D/S zo_A kivsa_A]_H אין זו כבשה
- Haita_S [l_R o_C]_A kivsa_A (היתה לו כבשה)

3. Imperatives

Imperatives in Hebrew do not require adding an IMP unit; the imperative itself can be marked as the A since it is conjugated for person.

- [Shlach]_P/A et ha michtav ba doar (שלח את המכתב ב-דוור)
- [Shilchi]_P/A et ha michtav ba doar (שליחי את המכתב ב-דוור)

4. PP Adverbials

We internally annotate PP adverbials as R+C:

- be_R hechlet_C (ב-החלט)
- me_R olam_C (מ-עולם)

- me_R chadash_C (מ-חדש)
- ka_R rauy_C (כ-ראוי)

In rare cases where an expression is completely opaque, it should be left unsegmented at the tokenization stage. If at the annotation stage you encounter a word that was wrongly segmented, you can use the re-tokenize feature to contract it back

5. Reflexive 'atzmi' ('עצמו'):

- Dani_A [rachatz 'atzm o UNA]_P (דני רחץ עצמו-ו)
- [ani_C ['atzm i UNA]_F]_A lo_D haiti_F 'ose_P [davar_C [ka ze]_E]_A (אני עצמי לא היתי עשה דבר כ-זה).

6. Levadi (לבדי) is marked D, internally it is unanalyzable:

- Dani_A 'asa [levad o UNA]_D et ha mesima (דני עשה בלבד-ו את (ה-משימה))

7. Oto (אותו)

- As a direct object it is internally analyzable: Pagashti_P/A [ot_R o_C]_A (פגשתי אותו-ו)
- When conveying sameness it is internally unanalyzable: [Be_R [ot o UNA]_E ha_F sefer_C] (ב-אותו-ה ספר)

8. Annotation of Participants:

In Hebrew, a Participant may appear as:

a. A free word:

- Dani_A halach_P [la_R gan]_A (דני הלך לא-גן)
- Hu_A nimtza sham (הוא נמצא שם)

b. A verb conjugated for person:

- Rainu_P/A seret_A (ראים סרט)
- Shamati_P/A kol_A [ba_R chutz_C]_A (שמעתי קול ב-חווץ)

c. Pronominal suffix:

- Savlanut_S i_A (סבלנות-ו)
- Hadavar_A [lakad et 'eina]_P v_A (ה-דבר לך את עיני-ו)

In some cases, a single scene may contain multiple referents to the same Participant.

We will then need to determine which referent to select as the A :

a. Prefer marking a free word as the A over conjugated verbs and suffixes; any other referent which was not marked as the A should then be marked F.

- **Hu_A** chazar_D ['al bakashat]_P o_F (הוּא חזר על בקשת-וּ)
- **Hu_A** [taman et yad]_{P-} o_F [ba tzalachat]_{-P} (הוּא טמן את יד-וּ ב-צלהת)

b. If there is no free word that can serve as the A, we prefer marking a conjugated verb as the A and then any other referent will be marked F.

- **Nisiti_D/A** [lehotzi mitachat yad]_P i_F tziur_A (ニシティ להוציא מתחת יד-וּ ציר)

c. If the only referent of a Participant is a pronominal suffix then it should be marked as the A:

[Chaschu 'eina]_P v_A (חשכו עיניו-וּ)

9. Remotes:

Generally, we prefer to add free words as Remotes over pronominal suffixes, but if there isn't a free word that can be added, it is OK to add a pronominal suffix.

- [Hora_S/A v_A]_A amru_P [I_R o_C]_A [lehafisk_D le'ashen_P (o)_A]_A (הורין אמרו לו-וּ להפסיק לעשן)

10. Purposive linkage:

Sometimes the infinitive “ל” can express a purposive linkage, but since we do not segment it, it cannot be marked as the Linker. In such cases add an IMP L instead.

[Dani_A halach_P [le_R merkaz_C ha 'ir_C]_A]_H (IMP L) [lirot_P [et_R ha_F hatzaga_P]_A (Dani)_A]_H (דני הלך למרכז העיר לראות את ה-הצגה)

11. Negative polarity items

We will prefer marking the negative word “לא” as the D over any negative polarity item that may also occur in the same scene (e.g. שום דבר, אף אחד, כלום, מואמה)

- **Lo_D raiti_P/A [klum]_A**

12. Ein/lo+ela construction:

a) When it is used as a contrasting negation construction (evokes two scenes) we mark ein/lo as D and ela as L:

- i) **[ein_D [ha yeladim]_A holchim_P [la migrash]_A]_H ela_L [[le beit sefer]_A (holchim)_P (yeladim)_A]_H** (אין ה-ילדים הולכים ל-מגרש אלא ל-בית-ספר)

ii) [Dani **lo_D** ochel_P 'agvania_A]_H **ela_L** [melafefon_A (Dani)_A ochel_P]_H
 (דָנִי **לֹא** אֹכֶל עֲגַבְנִיה **אֶלָא** מַלְפְּפָוָן)

- b) When it is used as an emphasis construction that does not evoke two separate scenes, we mark "ein/lo..ela" as a discontiguous UNA D:

- [ein_{D-} o_{-D} [ela]_A [dvarim tovim]_A lehagid_P]
 (אֵין לְאֶלָא) (דברים טובים להגיד)

13. Interesting examples

- a) Different uses of body parts:

- [ha_F ma'ase]_A orer_D **[be_R lib_C o_C]_A** regashot simcha_P
 (ה-מעשה עורך ב-לִבּ-וּ רגשות שמחה)
- [Ha_F yeled_C]_A **[amar be lib]_P o_F** (ה-ילד אמר ב-לִבּ-וּ)
 (דָנִי הַתָּה אָזֶן לְדִבְרֵי ה-מוֹרָה)
- Dani [hita ozen]_P [le divrei ha more]_A
 (דָנִי נָעַנְעַ רָאשׁ-וּ)
- Dani_A nianeia_P **[rosh_C o_E]_A**
 (דָנִי הַלְּבָבּ אָזֶן)

2. Tokenization guidelines

- 1) The following prefixes are segmented¹:

ב,כ,ל,א	Prepositions	דָנִי הַלְּבָבּ אָזֶן דָנִי נָמַצֵּא בְּ בַּיְתָה' דָנִי הַלְּבָבּ אָשָׁם יְפָה כְּ פָרָח Including where the prefix begins a prepositional phrase (PP) adverbial ² : תְּנָהָג בְּ זְהִירָה
ה	Definite article	הַיְלָד הַלְּבָבּ מָהָרָה. Including when attached to demonstratives ³ :

¹ Note that we do not segment the infinitive “ל”:
לְרֹאָתָה, לְשֻׁתָּועָה

² Generally, prefer to segment prepositional prefixes wherever possible:

דָנִי הַתְּחִילָה מְחֻדֶשׁ, דָנִי בַּיְקָשׁ שֶׁ יַתַּקְשֵׂר אֵלֵי יְהִי בְּ הַקָּדָם, דָנִי בְּ הַחְלָטָה מְעוֹנֵן לְהַשְׁתַּחַף, דָנִי נָמַצֵּא בְּ פְנִים
 דָנִי הַתְּנָהָגָה כְּ רָאוּי

³ An exception would be **הַלְלוּ** which should not be segmented.

		המחשב ה זה עבד האיש ה הוא מוכר
ה	Interrogative	ה רואה אתה את מה ש אני רואה?
ש, ה	Relativizers	הכלב ש דני ראה הוא שחור הילד ה יושב בשורה הראשונה הוא דני
ש, כש, מש, לcas	Subordinating conjunctions	דני ראה ש הכלב מתקרב. כח דני יבוא, נתחיל.
ו	Coordinating conjunction	דני ו דנה
כ	Adverb	דני מכיר כ אלף פרחים

2) Pronominal suffixes are segmented:

a) *Possessive pronominal suffixes appended to nouns:*

ספר: ספר י, ספר ר, ספר ו, ספר ה, ספר נו, ספר כם, ספר כן, ספר מ, ספר ו
 ספרים: ספר י, ספר ר, ספר ו, ספר ה, ספר נו, ספר כם, ספר כן, ספר המ, ספר חן
 שפה: שפת י, שפת ר, שפת ו, שפת ה, שפת נו, שפת כם, שפת כן, שפת מ, שפת ו
 שפות: שפות י, שפות ר, שפות ו, שפות ה, שפות נו, שפות ים, שפות ים, שפות מ, שפות ו

b) *Pronominal suffixes as the objects of prepositions:*

בשביל: בשביל י, בשביל ר, בשביל ו, בשביל ה, בשביל נו, בשביל כם, בשביל כן, בשביל מ, בשביל ו
 לפני: לפני י, לפני ר, לפני ו, לפני ה, לפני נו, לפני כם, לפני כן, לפני המ, לפני חן
 ל: ל י, ל ר, ל ו, ל ה, ל מ, ל נו, ל כם, ל כן, ל המ, ל חן
 מ: ממ י, ממ ר, ממ ו, ממ ה, ממ נו, ממ (מאות נו), ממ כם, ממ כן, ממ המ, ממ חן
 את: את י, את ר, את ו, את ה, את נו, את כם, את כן, את מ, את ו, את חן

c) *Pronominal suffixes as the direct objects of verbs*

הקייף הקיף י, הקיף ר, הקיף ו, הקיף ה, הקיף נו, הקיף כם, הקיף כן, הקיף מ, הקיף ו

d) *Pronominal suffixes appended to the infinitive construct:*

i) *as the subject of the infinitive:*

ב ראות זאת, התחזקת ב' ה הרגשה

ii) as the object of the infinitive:

דני מנסה לראות **ה** כל יומ

3) We segment contracted existential particles from their pronouns

אין **ת**, אין **ר**, אין **ו**, אין **ה**, אין **תּו**, אין **כּו**, אין **כּן**, אין **מּו**, אין **נּו**
יש **ת**, יש **ר**, יש **ו**, יש **ה**, יש **תּו**, יש **כּו**, יש **כּן**

4) We segment contracted reflexives as: "לבד" "עצמ" and "עצמ"

עצמ **ו**, עצם **ר**, עצם **ו**, עצם **ה**, עצם **תּו**, עצם **כּו**, עצם **כּן**, עצם **מּו**, עצם **נּו**
לבד **ו**, לבד **ר**, לבד **ו**, לבד **ה**, לבד **תּו**, לבד **כּו**, לבד **כּן**

5) We segment the following blended forms: זהו, זוהי, מהו, מיהו, איזה

- (a) **זהו** ה בית של דני
- (b) **זו היא** התשובה ה נכונה
- (c) **מהו** ה רעיון ה עומד מאחוריו היוזמה?
- (d) **מי הוא** האיש העומד שם?
- (e) **איזה** עשיר השם ב חלקו

6) Also, note that when conveying obligation, "על" in its contracted forms should be segmented as well:

על **ו** להשתתף
על **ינו** לעשות את הדבר הנכון

7) We do not segment the construct state (e.g. ספרי ילדים, שמלה ה נשף)

8) We do not segment the binyan or subject agreement prefixes/suffixes of verbs (התרכצתי, תלמיד)

9) We do not segment gender/number suffixes on nouns and adjectives: ספרים טובים, תלמידה טובה

10) We do not segment acronyms (e.g. צה"ל).

11) Note that because of an automatic tokenization process that occurs before you receive the passage, you may encounter a word that was unnecessarily tokenized for containing a geresh (e.g. ג'ונגל may appear ג'ונגָל). In such cases please delete the unnecessary space (should be simply ג'ונגל).

12) You might also see that due to automatic tokenization spaces occur between a makaf and the words it connects (בית - ספר). Those spaces are correct, so please do not delete them.

13) We segment prepositions in fixed expressions like ב. בקשה

- 14) We do not segment derivational morphemes like עירא**ק**, מהי**רות**
- 15) We do not segment morphemes in borrowed words like אונט**יבש****מי**.
- 16) We do not segment complex and compound prepositions (that are combinations of a preposition+noun or preposition+preposition). For example, the preposition באמצ**יעת** should not be segmented.
- 17) We do not segment compound conjunctions like מאחר ש-, מפני ש-
- 18) We do not segment compound questions words such as למה and לאן, כמה, איפה (note that both למה and לאן should not be segmented)

Prepositions in the Hebrew Treebank (UD) and their Frequency

Use this list to identify prepositions that you are uncertain of. If you find in the text a preposition that does not appear on this list, please mention this in the comment field in the follow-up spreadsheet.

ב	7928
של	4830
ל	4438
את	1998
מ	1697
על	1503
כ	573
עם	457
בין	281
מן	213
עד	211
כדי	185
לפנ י	153
לאחר	148
אל	147
כמו	141
נגד	124
אחר י	114
לא	71
לפי	56
לעומ ת	50
אחר	47
באמצ יעת	46
למרות	44
בגל	44
בידי	43
מפני	39

38 בעקבות	38
38 במשן	38
37 תור	37
36 בל'י	36
35 ליד	35
35 לגב'י	35
32 במקום	32
30 לבין	30
29 בפנ'י	29
29 אצל	29
28 מול	28
28 בשל	28
27 מתוֹר	27
27 בתוֹר	27
23 לקראת	23
23 לעבר	23
22 למען	22
22 כלפי	22
22 בעוד	22
21 תחת	21
21 מעל	21
20 תמורה	20
20 בקרוב	20
19 עברו	19
18 בשביל	18
18 בגין	18
17 סביב	17
16 מאחוריו	16
16 بعد	16
15 מצד	15
15 בלי	15
14 מבין	14
12 לנוכח	12
11 מתחת	11
11 מעבר	11
11 לאורך	11
11 בטרם	11
11 בזכות	11
10 לשם	10
10 ליד'	10
10 החל	10
9 מבלתי	9
9 מאות	9
9 לצד	9
9 כש	9
9 בעבר	9
8 נוכח	8
8 כרגע	8

7	הודות
6	ע"י
6	בלבד
6	כנגד
6	דרך
5	מבعد
5	לאור
5	באוזני
4	לזכות
4	חו"מ
4	בעבור
4	בגדר
3	מש
3	לרגל
3	לפנות
3	לכבוד
3	בהתחשב
3	בתור
3	בצד
3	במו
3	בלעדיו
2	פרט ל
2	סמור ל
2	נסוף ל
2	משך
2	מחוצה ל
2	لتוך
2	למן
2	לכדי
2	כאל
2	בתוככי
2	בלית
1	קדם ל
1	משל (כמו "כאילו")
1	לצורך
1	למעין
1	לכעין
1	לדידי
1	כען
1	זולת
1	במסגרת
1	בגנות
1	אודות

ביטויי יחס נוספים: על מנת, מחוץ, בפי