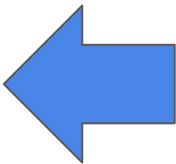


עיבוד שפה טבעית ש9:

תחביר

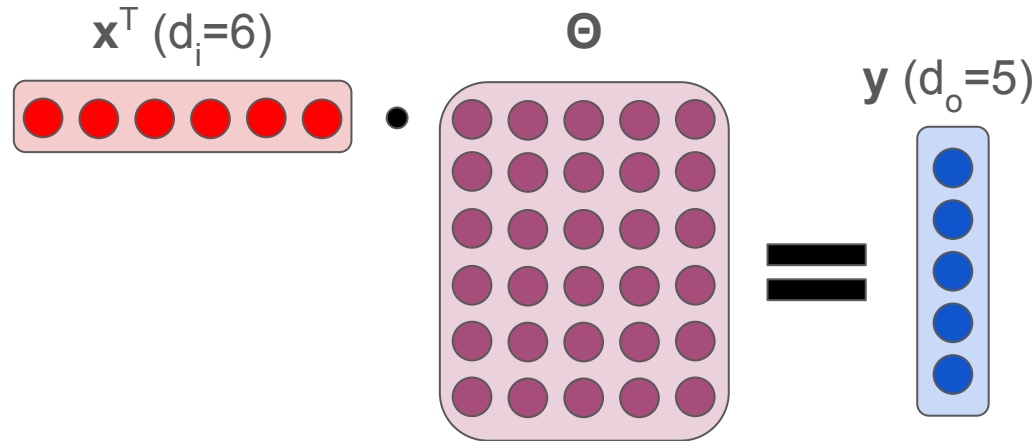
(SLP 17-18, E 9-11.1)

פתרון הבוחן



השלמה - כפל מטריצות כפעולה ברשת נוירונים

- ברשת FFN רגילה: שכבה מקבלת וקטור x ממימד d_i ומוציאה וקטור y ממימד d_o



- ולכן Θ הוא מטריצה:
 $y = x^T \Theta$

השלמה - כפל מטריצות כפעולה ברשת נוירונים

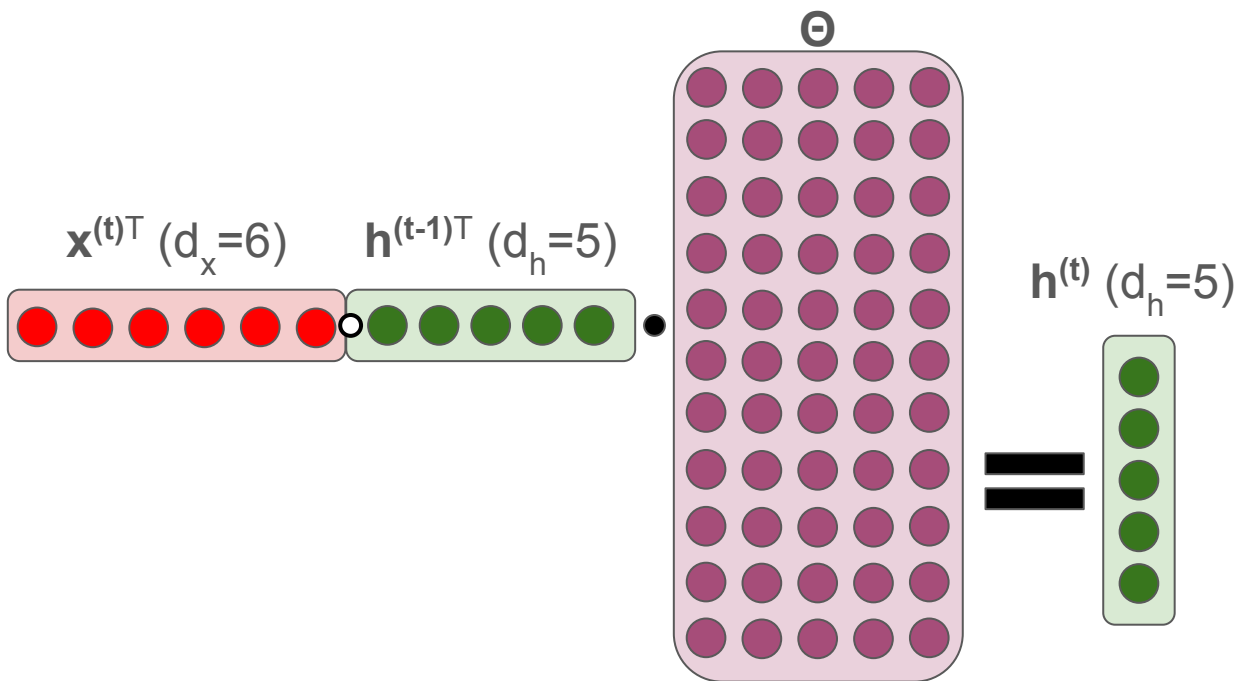
- ברשת נשנית: שכבה מקבלת וקטור x ממימד d_x

ווקטור חבוי h ממימד d_h

ומוציאה וקטור h ממימד d_h

- Θ נשאר מטרצה:

$$h^{(t)} = \Theta((x \circ h^{(t-1)})^T)$$

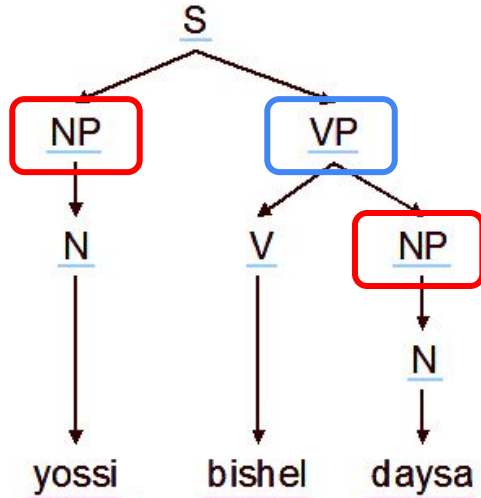
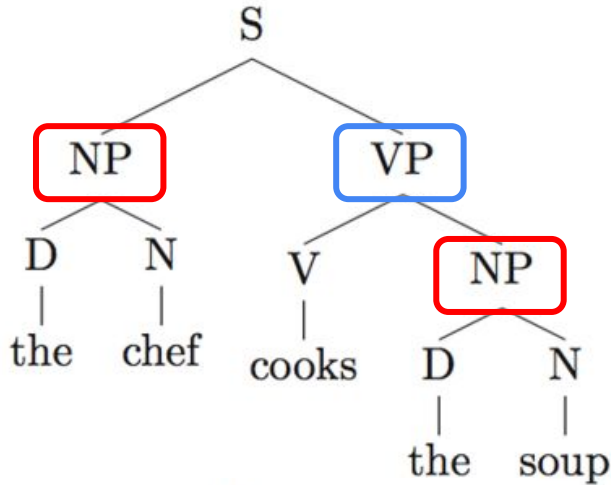


יחידת יסוד תחבירית: צירוף (phrase)

- הילדה אכלה את התפוח
 - [הילדה] [אכלה את התפוח]
 - [הילדה] [אכלה את [התפוח]]
 - [מרב] אכלה את התפוח
 - [היא] אכלה את התפוח
 - הילדה [נסעה]
 - הילדה אכלה את [זה]
- צירוף שמני (NP - noun phrase)
- צירוף פעלי (VP - verb phrase)
- ?

לכל ביטוי יש ראש תחבירי לפי הסוג שלו ולפי חלק הדיבר המתאים

תחביר ומבנים



- בדקדוק מרכיבי (constituency grammar) ניתן לחלק משפטים לצירופים (phrases) סמוכים ומקוננים
- הצירופים נקראים על-שם ה"ראש", שהוא המרכיב המגדיר אותם לשונית
- אתגר במציאת מבנה המשפט (לאנשים ולא לגוריתמים) - דו-משמעות
 - לקסיקלית - מילה יכולה לבטא כמה משמעויות, או כמה חלקי דיבר (חולצה מטיילת בואדי)
 - תחבירית (מבנית) - אותן מילים, אותם חלקי דיבר, פרשנות מבנה שונה (\Leftarrow)

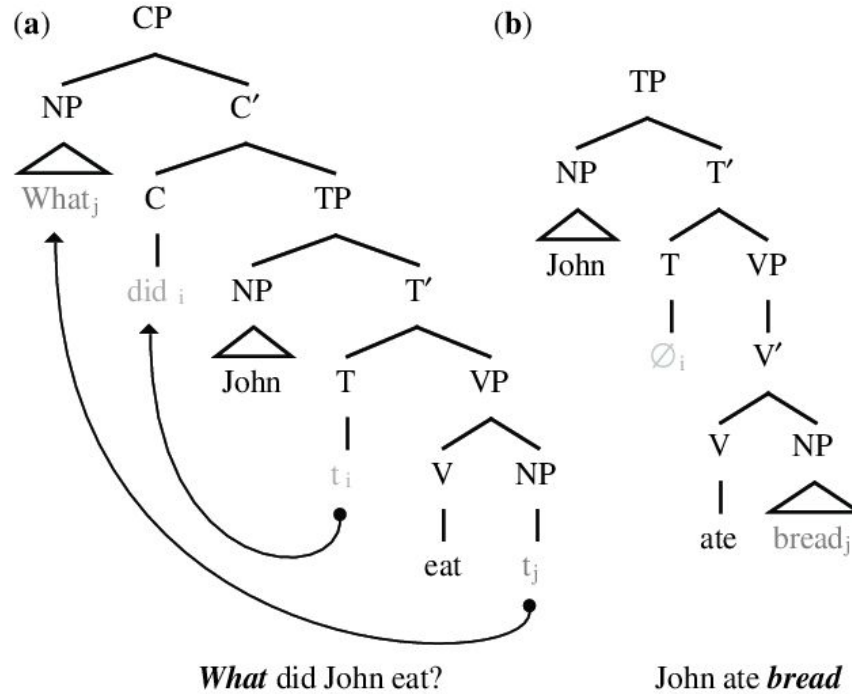
דוגמאות לדו-משמעות מבנית

- הממשלה מבטלת את התו הירוק בבריכות, במסעדות ובאירועים זועמים.
- כסא אופניים לילד שמתחבר לסבל
- מכתב שכתבה מורה לתיאטרון

תחביר כהסבר: קורפרנס (coreference), אנאפורה (anaphora)

- (האינדקסים j, i מציינים שמדובר באותו אדם)
- שפרה_i מדברת אליה_j.
- * שפרה_i מדברת אליה_j.
- שפרה_i מדברת אל עצמה_i.
- * היא_i מדברת לשפרה_j.
- שפרה_i אוהבת שהיא_j לפעמים שותקת.
- שפרה_i אוהבת שהיא_j לפעמים שותקת.
- שפרה_i ביקשה מאגאתה_j לעזור לה_i.
- * שפרה_i ביקשה מאגאתה_j לעזור לעצמה_i.
- שפרה_i ביקשה מאגאתה_j לעזור לעצמה_j.
- (מעוניינים בהסברים בלשניים? גגלו Binding theory)

Movement “תנועה”



תחביר כהסבר (II): מגבלות על תנועה

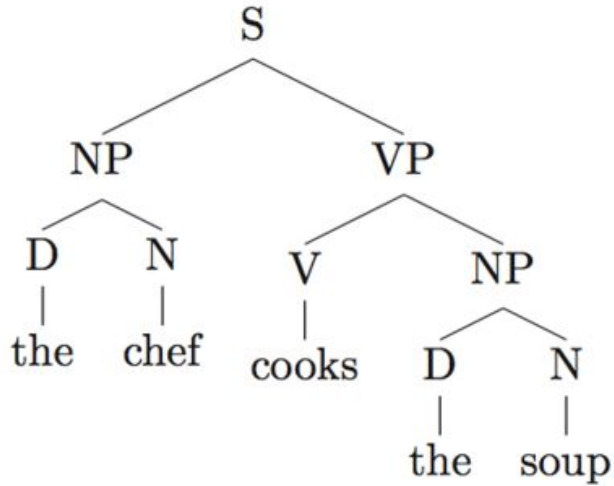
- יוסי בישל דיסה.
 - מי בישל דיסה?
 - מה יוסי בישל?
 - את מי יוסי פגש?
- רינה טענה שיוסי בישל דיסה.
 - מי רינה טענה שבישל דיסה?
 - מה רינה טענה שיוסי בישל?
- רינה קבעה את הקביעה שיוסי בישל דיסה.
 - * מה רינה קבעה את הקביעה שיוסי בישל?
- Rina bought Yossi's book
 - Whose book did Rina buy?
 - * Whose did Rina buy book?
- (מעוניינותים בהסברים בלשניים? גגלו Ross's Islands)

דקדוק חסר-הקשר (CFG - Context Free Grammar)

- אמצעי "מעלה-מטה" (top-down) להתחיל מצומת השורש ולאט לאט **לגזור** עוד צמתים עד שמגיעים למשפט
 - גנרטיבי מאוד
- כלל המשפטים החוקיים בשפה נקבעים ע"י **כללי הגזירה**, שיכולים להיות רקורסיביים
- פורמלית:

צירופים ותגים	N a set of non-terminal symbols (or variables)
מילים	Σ a set of terminal symbols (disjoint from N)
כללי גזירה	R a set of rules or productions, each of the form $A \rightarrow \beta$, where A is a non-terminal, β is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$
שורש	S a designated start symbol and a member of N

כללי גזירה



- $S \rightarrow NP VP$
- $NP \rightarrow D N$
- $VP \rightarrow V NP$
- $D \rightarrow \text{the}$
- $N \rightarrow \text{chef} \mid \text{soup}$
- $V \rightarrow \text{cooks}$

• (בדקדוק הזה יש רקורסיה?)

• צורת חומסקי נורמלית (**CNF**): כל כלל כולל מימין או (1) שני נונטרמינלים (nonterminals), או (2) מילה אחת.

○ (האם הדקדוק לעיל הוא CNF?)

עוצמת CFG [דילגנו בשיעור]

- יותר כללי גזירה - יותר עצים אפשריים
- אבל, כמה כללי גזירה צריך?
- כל צורת ביטוי אפשרי בשפה דורשת עוד כלל

Grammar	Lexicon
$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid the \mid a$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid money$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid prefer$
$NP \rightarrow Pronoun$	$Pronoun \rightarrow I \mid she \mid me$
$NP \rightarrow Proper-Noun$	$Proper-Noun \rightarrow Houston \mid NWA$
$NP \rightarrow Det Nominal$	$Aux \rightarrow does$
$Nominal \rightarrow Noun$	$Preposition \rightarrow from \mid to \mid on \mid near \mid through$
$Nominal \rightarrow Nominal Noun$	
$Nominal \rightarrow Nominal PP$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	
$VP \rightarrow Verb NP PP$	
$VP \rightarrow Verb PP$	
$VP \rightarrow VP PP$	
$PP \rightarrow Preposition NP$	

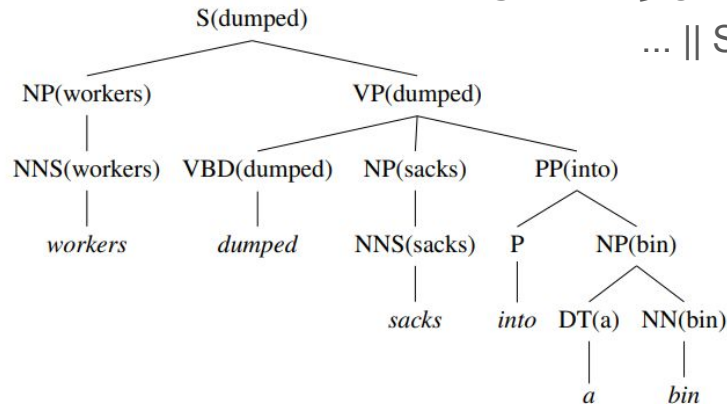
Figure 13.1 The \mathcal{L}_1 miniature English grammar and lexicon.

לקסיקליזציה (אייזנסטיין 10.5.2)

- דקדוק מועשר בידע מילוני

- מסגרת סבטגוריזציה subcategorization frame - "מה המילה מקבלת כמשלים"

- עמדה: נושא NP
- כתבה: נושא NP (+ מושא ישיר NP)
- נתנה: נושא NP + מושא ישיר NP (+ מושא עקיף PP-ל)
- רצתה: נושא NP + מושא ישיר NP || נושא NP + מושא פסוקי משועבד SBAR
- אמרה: נושא NP + מושא עקיף NP-ל + מושא פסוקי לא נטוי S || ...



- תת-אפיון של הצירופים

- NP-SBJ - צירוף בעמדת נושא, NP-OBJ בעמדת מושא, וכו'

- סימון ראש הצירוף על-גבי העץ

אוניברסליות?

- כבר באנגלית, הרבה תופעות צריכות הסברים בנוסח "תנועה" או "אלמנטים ריקים"
 - What did the man want?
 - I brought the notebook to class and the laptop to the office
- בשפות עם סדר מילים יותר חופשי, קשה לקבץ הכל בצירופים או לתרץ את המבנים
- בשפות שבהן כללי התאם חשובים, לפעמים קשה להסיק עליהם מתוך המבנים

מבנה תלויות (Dependency structure) (אייזנסטיין 11.1)

- נתאר ישירות קשר בין מילים במשפט ונבנה היררכיה שבה רק המילים הן צמתים

the chef cooks the soup

מבנה תלויות (Dependency structure) (אייזנסטיין 11.1)

- נתאר ישירות קשר בין מילים במשפט ונבנה היררכיה שבה רק המילים הן צמתים

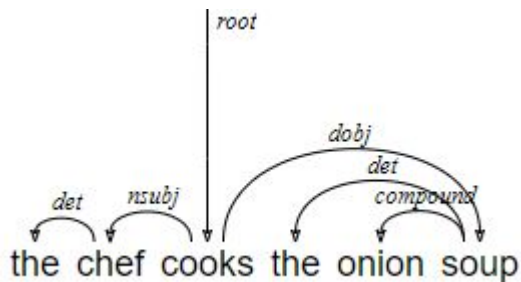
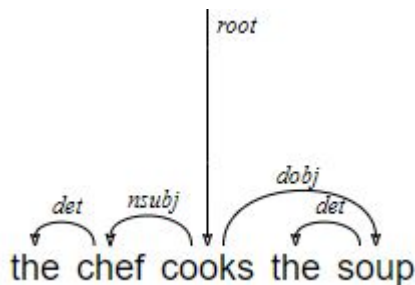
- בפועל, נדבר על עץ תלויות (dependency tree)

- קשת מוגדרת ע"י:

- המקור שלה, או ראשה **head**

- היעד שלה, או תלויה **dependent**

- הסוג שלה **type / label / relation**



סוגי קשתות בעץ תלויות

- קשתות בין נושא (predicate) לארגומנט ליבה (core argument)

- נושא שמני nsubj
- נושא פסוקי csubj
- מושא ישיר dobj
- מושא עקיף iobj
- פסוקית משלימה ccomp, xcomp

- קשתות בין נושא למרכיבי עזר (non-core dependent)

- נספחים שמניים obl
- פסוקית אופן advcl
- פועל עזר aux, cop

סכימות התיוג
מתעדכנות עם
השנים

כשלא מוגבלים לעצים, כאן ייתכנו קשתות נכנסות כפולות

סוגי קשתות בעץ תלויות

סכימות התיוג
מתעדכנות עם
השנים

- קשתות בין נשוא (predicate) לארגומנט ליבה (core argument)
- קשתות בין נשוא למרכיבי עזר (non-core dependent)
- קשתות בתוך צירופים שמניים
 - תווית יידוע determiner
 - סמיכות nmod
 - שם תואר amod
 - מילת יחס או יחסה case

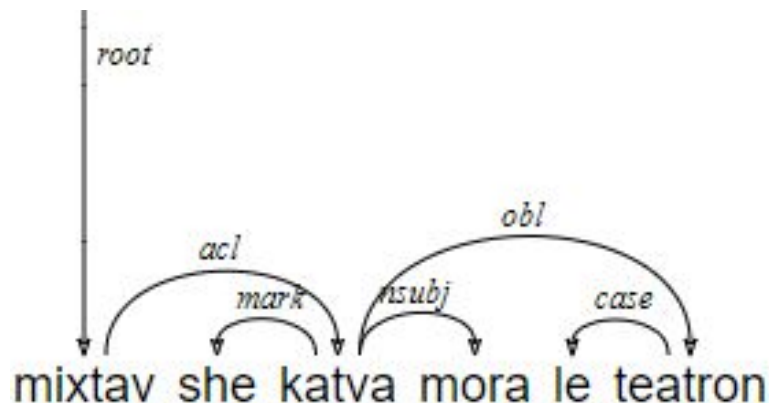
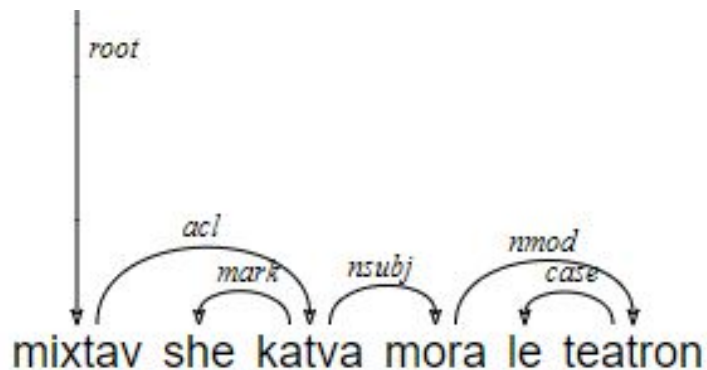
כשלא מוגבלים לעצים, כאן ייתכנו קשתות נכנסות כפולות

● אחרות

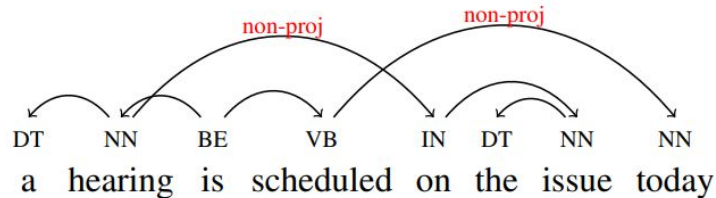
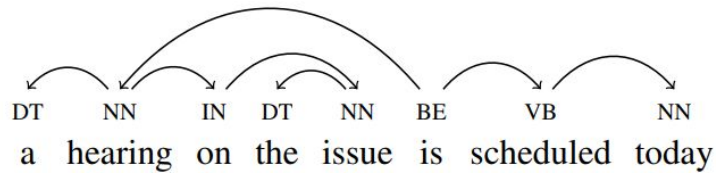
- חיבור conj, cc, list
- "מיוחדים" goeswith

דו-משמעות בעץ תלויות

- מכתב שכתבה מורה לתיאטרון



היטליות (Projectivity)



- התכונה לפיה קשתות העץ לא חוצות זו את זו

- התייחסות יחידה לסדר המילים בדקדוק תלויות
- האם יש התייחסות כזו בדקדוק מרכיבי?

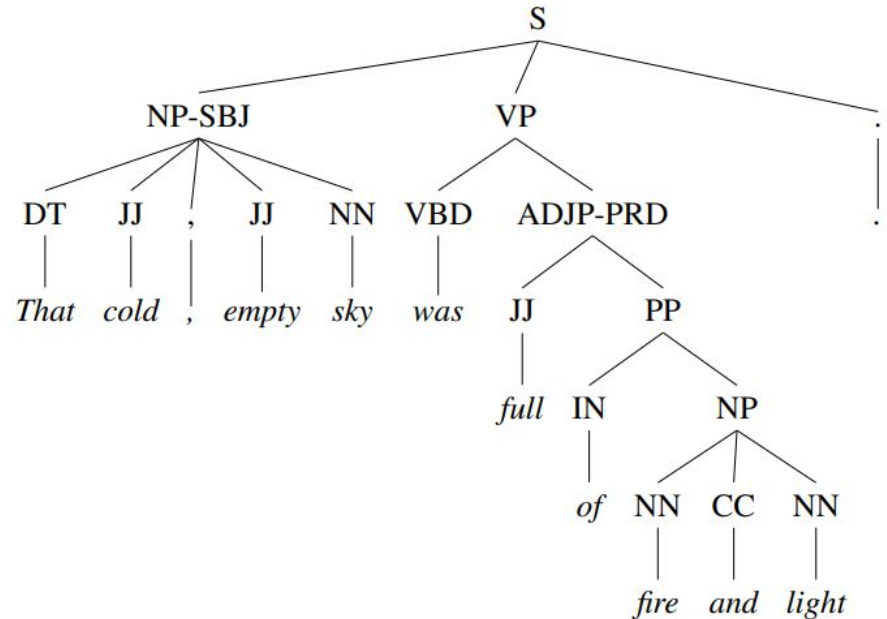
- האם יש בשפה טבעית משפטים שהעץ שלהם אינו היטלי?

- הגדרה פורמלית: קשת מ-h ל-d היא היטלית ⇔ כל המילים שבין h ו-d ברצף המילים הן צאצאיות של h בגרף; עץ הוא היטלי ⇔ כל הקשתות בו היטליות

משאבים (מרכיבי)

- אוסף העצים של פן (Penn Treebank, PTB)

```
((S
  (NP-SBJ (DT That)
    (JJ cold) (, ,)
    (JJ empty) (NN sky) )
  (VP (VBD was)
    (ADJP-PRD (JJ full)
      (PP (IN of)
        (NP (NN fire)
          (CC and)
          (NN light) ))))
  (. .) ))
```



משאבים (תלויות)

● תלויות אוניברסליות (UD) Universal Dependencies

Abaza	1	<1K	☐	Northwest Caucasian	Hindi	2	375K	☐W	IE, Indic	Odia	1	<1K	☐☐	IE, Indic
Afrikaans	1	49K	☐☐	IE, Germanic	Hindi English	1	26K	☐	Code switching	Old Church Slavonic	1	57K	☐	IE, Slavic
Akkadian	2	25K	☐☐	Tupian, Tupari	Hittite	1	1K	☐	IE, Anatolian	Old East Slavic	3	184K	☐☐	IE, Slavic
Akuntsu	1	<1K	☐☐	IE, Albanian	Hungarian	1	42K	☐☐☐☐W	Uralic, Ugric	Old French	1	199K	☐☐☐☐	IE, Romance
Albanian	1	<1K	W	Afro-Asiatic, Semitic	Icelandic	3	1,162K	☐☐☐☐W	IE, Germanic	Old Irish	2	21K	☐☐☐☐	IE, Celtic
Amharic	2	10K	☐☐☐☐☐	IE, Greek	Indonesian	3	169K	☐☐☐☐W	Austronesian, Malayo-Sumbawan	Old Japanese	1	<1K	☐	Japanese
Ancient Greek	2	416K	☐☐☐☐	IE, Armenian	Irish	2	131K	☐☐☐☐W	IE, Celtic	Old Turkish	1	<1K	☐	Turkic, Northeastern
Apurina	1	<1K	☐☐	Afro-Asiatic, Semitic	Italian	7	818K	☐☐☐☐☐W	IE, Romance	Persian	2	654K	☐☐☐☐☐☐☐☐☐	IE, Iranian
Arabic	3	1,042K	☐☐W	Mande	Japanese	8	2,849K	☐☐☐☐☐☐W	Japanese	Polish	3	499K	☐☐☐☐☐	IE, Slavic
Armenian	2	55K	☐☐☐☐☐☐☐	Basque	Kanuri	1	2K	☐☐	Austronesian, Javanese	Portuguese	4	570K	☐☐☐☐☐	IE, Romance
Assyrian	1	<1K	☐☐	Afro-Asiatic, Semitic	Kabyle	1	47K	☐☐	Tupian, Tupi-Guarani	Prakrit	1	<1K	☐☐	IE, Indic
Bambara	1	13K	☐☐	Afro-Asiatic, Cushitic	Karelian	1	2K	☐☐	Afro-Asiatic, Berber	Romanian	4	937K	☐☐☐☐☐☐☐W	IE, Romance
Basque	1	121K	☐	IE, Slavic	Karo	1	2K	☐☐☐	IE, Indic	Russian	4	1,832K	☐☐☐☐☐☐W	IE, Slavic
Beja	1	1K	☐	IE, Indic	Kazakh	1	10K	☐☐☐	Uralic, Finnic	Sanskrit	2	28K	☐☐	IE, Indic
Belarusian	1	305K	☐☐☐☐☐W	IE, Celtic	Khunsari	1	<1K	☐☐☐	Tupian, Purubora-Ramarama	Scottish Gaelic	1	84K	☐☐☐☐	IE, Celtic
Bengali	2	<1K	☐☐☐W	IE, Slavic	Kiche	1	10K	☐☐☐☐☐W	Turkic, Northwestern	Serbian	1	97K	☐☐	IE, Slavic
Bhojpuri	2	6K	☐☐	Mongolic	Komi Permyak	1	<1K	☐	Mayan	Sindhi	1	6K	☐☐	IE, Indic
Breton	1	10K	☐☐☐☐☐W	Sino-Tibetan	Komi Zyrjan	2	10K	☐☐	Uralic, Permian	Skolt Sami	1	2K	☐☐☐	Uralic, Sami
Bulgarian	1	156K	☐☐☐	IE, Romance	Korean	5	446K	☐☐☐☐☐☐W	IE, Celtic	Slovak	1	106K	☐☐	IE, Slavic
Buryat	1	10K	☐☐☐☐☐	Sino-Tibetan	Kurmanji	1	10K	☐☐☐	IE, Iranian	Slovenian	2	170K	☐☐☐☐	IE, Slavic
Cantonese	1	13K	☐☐	Chukotko-Kamchatkan	Latin	5	977K	☐☐☐☐☐☐	Mayan	Soi	1	<1K	☐☐	IE, Iranian
Catalan	1	553K	☐	Sino-Tibetan	Latvian	1	265K	☐☐☐☐☐	IE, Latin	South Levantine Arabic	1	<1K	☐☐	Afro-Asiatic, Semitic
Chinese	5	285K	☐☐☐☐W	Afro-Asiatic, Egyptian	Laz	1	2K	☐☐	IE, Baltic	Spanish	3	1,022K	☐☐☐☐☐W	IE, Romance
Chukchi	1	6K	☐☐	IE, Slavic	Ligurian	1	6K	☐☐☐☐☐☐W	Kartvelian	Swedish	3	206K	☐☐☐☐☐W	IE, Germanic
Classical Chinese	1	283K	☐☐	IE, Germanic	Lithuanian	2	75K	☐☐☐☐	IE, Romance	Swedish Sign Language	1	1K	☐	Sign Language
Coptic	1	52K	☐☐☐	Uralic, Finnic	Lluli	1	1K	☐☐☐☐	IE, Baltic	Swiss German	1	1K	☐☐☐☐W	IE, Germanic
Croatian	1	199K	☐☐☐☐	IE, Germanic	Low Saxon	1	2K	☐☐	IE, Celtic	Tagalog	2	1K	☐☐☐	Austronesian, Central Philippine
Czech	6	3,428K	☐☐☐☐☐☐W	IE, Romance	Magahi	2	7K	☐☐☐☐	IE, Indic	Tamil	2	12K	☐☐	Dravidian, Southern
Danish	2	100K	☐☐☐☐	IE, Germanic	Malay	1	<1K	☐☐	Tupian, Tupari	Tatar	1	1K	☐☐	Turkic, Northwestern
Dutch	2	306K	☐☐☐	IE, Germanic	Maltese	1	44K	☐☐☐☐	Afro-Asiatic, Semitic	Telugu	1	6K	☐☐	Dravidian, South Central
English	11	1,943K	☐☐☐☐☐☐☐☐☐☐☐W	IE, Germanic	Manx	1	20K	☐☐☐☐☐☐W	IE, Celtic	Thai	1	22K	☐☐☐W	Tai-Kadai
Erzya	1	17K	☐☐☐☐☐	Uralic, Mordvin	Marathi	1	3K	☐☐☐	IE, Indic	Tupinamba	1	1K	☐☐	Tupian, Tupi-Guarani
Estonian	2	511K	☐☐☐☐☐☐	Uralic, Finnic	Mbaia Guaraní	2	13K	☐☐	Tupian, Tupi-Guarani	Turkish	9	733K	☐☐☐☐W	Turkic, Southwestern
Farose	2	50K	☐☐☐☐	IE, Germanic	Middle Irish	2	<1K	☐☐☐☐	Creole	Turkish German	1	37K	☐	Code switching
Finnish	4	397K	☐☐☐☐☐☐☐W	IE, Romance	Moksha	1	3K	☐☐	Uralic, Mordvin	Ukrainian	1	122K	☐☐☐☐☐☐☐W	IE, Slavic
French	9	1,195K	☐☐☐☐☐☐☐☐W	IE, Germanic	Mundurucu	1	<1K	☐☐	Tupian, Mundurucu	Upper Sorbian	1	11K	☐☐	IE, Slavic
Frisian	1	51K	☐☐☐☐☐	Code switching	Najia	1	140K	☐	Creole	Urdu	1	138K	☐☐	IE, Indic
Frisian Dutch	1	3K	☐	IE, Romance	Nayini	1	<1K	☐☐	Uralic, Permian	Uyghur	1	40K	☐	Turkic, Southeastern
Galician	2	164K	☐☐☐☐☐	IE, Germanic	Neapolitan	1	<1K	☐☐	IE, Romance	Vietnamese	1	42K	☐	Astro-Asiatic, Viet-Muong
German	4	3,810K	☐☐☐☐☐☐W	Tupian, Tupi-Guarani	North Sami	1	26K	☐☐	Uralic, Sami	Warlpiri	1	<1K	☐☐	Pama-Nyungan
Gothic	1	55K	☐☐☐	Afro-Asiatic, Semitic	Norwegian	3	666K	☐☐☐☐	IE, Germanic	Welsh	1	41K	☐☐☐☐W	IE, Celtic
Greek	1	63K	☐☐☐W	IE, Slavic						Western Armenian	1	92K	☐☐☐☐	IE, Armenian
Guajajara	1	2K	☐☐	IE, Indic						Wolof	1	44K	☐☐W	Niger-Congo, Northern Atlantic
Hebrew	1	161K	☐☐	IE, Germanic						Xibe	1	15K	☐☐	Tungusic
										Yakut	1	<1K	☐☐	Turkic, Northeastern
										Yoruba	1	8K	☐☐W	Niger-Congo, Defoid
										Yupik	1	2K	☐☐	Eskimo-Aleut

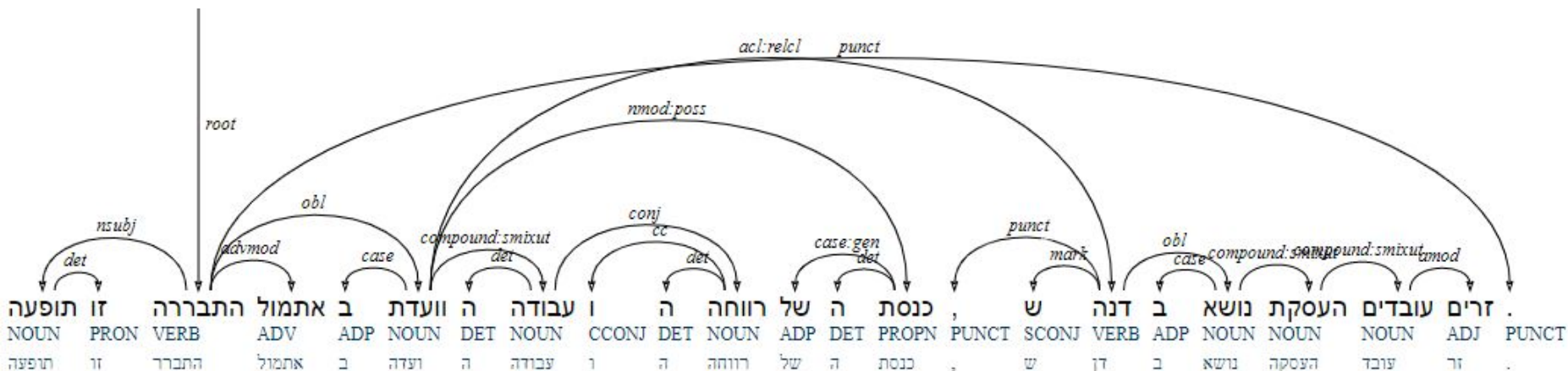
פורמט CoNLL-U

```
# sent_id = 2
```

text = תופעה זו התבררה אתמול בוועדת העבודה והרווחה של הכנסת, שדנה בנושא העסקת צובדים זרים

[illegible]

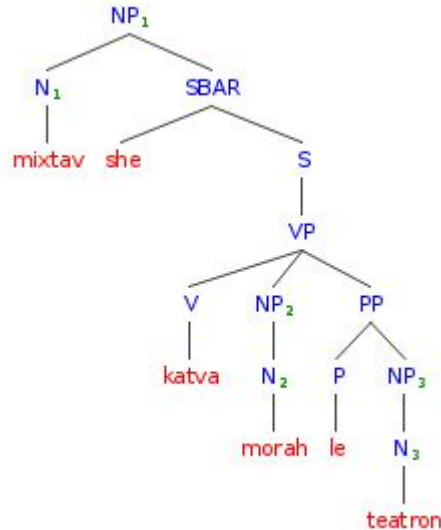
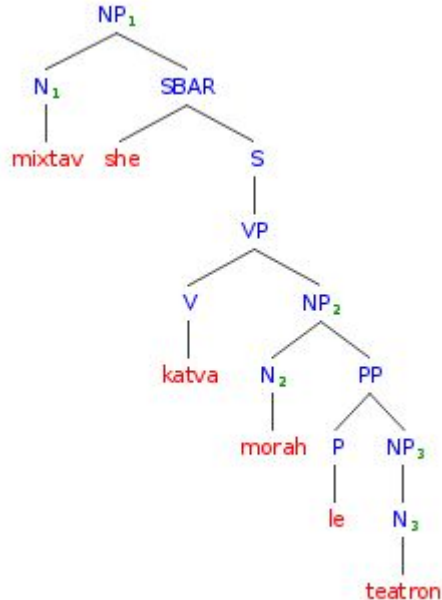
פורמט CoNLL-U



<https://arborator.ilpqa.fr/q.cgi>

מטריקות (מרכיבי) [דילגנו בשיעור]

- פרסבל PARSEVAL: לוקחים את כל המרכיבים בעץ ה"נכון" (gold) ובודקים כמה המערכת זיהתה נכון לגמרי (recall), לוקחים את מרכיבי העץ החזוי ובודקים כמה מהם נכונים (precision), מדווחים F1



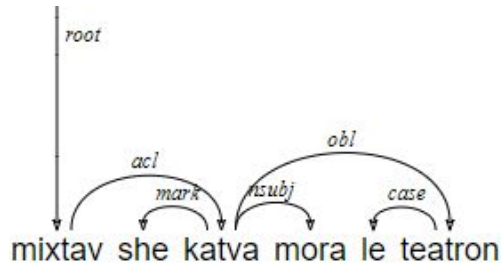
- חציית סוגריים (cross-bracketing): יחס המרכיבים הנכונים בלי "הפרעה במבנה": סוגר [בלתי סגור בתוך]

מטריקות (תלויות)

- ציונים פר-מילה (ראשה או תלויה?)

- עם חשיבות לסוג הקשת LAS - Labeled Attachment Score
 - בלי חשיבות לסוג הקשת UAS - Unlabeled Attachment Score
 - רק סוג הקשת, בלי חשיבות לראשה LS - Label Score
- איזו מין בעיה זו?

- אינטואיציה: כמה תקבל מערכת על משפט שפורסר נכון חוץ משגיאה בתליית צירוף יחס?



- עוד אפשרות: דיוק פר סוג קשת (לאורך הדאטאסט)

- מטריקות?

ניתוח תחבירי Syntactic Parsing

- מחר. בינתיים:
- כמה עצי מרכיבים אפשריים יש למשפט בן n מילים?
- כמה עצי CNF אפשריים יש למשפט בן n מילים?
- כמה עצי תלויות אפשריים יש למשפט בן n מילים?
- כמה עצי תלויות היטליים אפשריים יש למשפט בן n מילים?