

כלי לחיפוש בלשני בעברית

מהו חיפוש בלשני?

- CPQ is flexible and efficient query processor developed by IMS Open Corpus Workbench (CWB).
- It allows to perform linguistically motivated queries, e.g., search for a particular lemma followed by a word (verb?) in past tense.

- כפיר וחיים שלשוילי בהדרכת פרופ' שולי וינטנר מאוניברסיטת חיפה התאימו את CPQ לעברית
- דליה בוז'ן ואלון איתי בנו web interface שמנו דגש הוא על קלות השימוש תוך ויתור על האופציה לחיפוש בוליאני.

כל שאילתה מורכבת מסדרה של עד 4 מילים. כל מילה ניתן לחפש על פי תמנית פני השטח (surface token), הערך הלקסיקלי (הלמה, הערך שמופיע במילון) ולגבי פעלים ניתן לחפש גם על-סמך השורש.

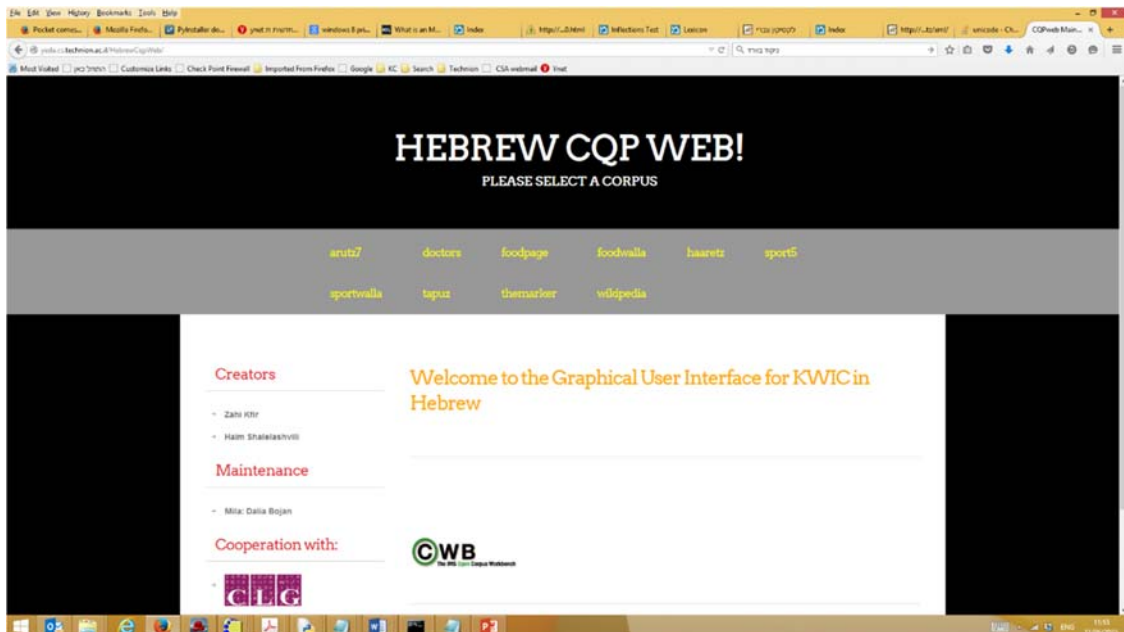
את התכנה ניתן להפעיל בקשורית הבאה <http://yeda.cs.technion.ac.il/HebrewCqpWeb>

בעיות נא לשלוח ל-mila@cs.technion.ac.il

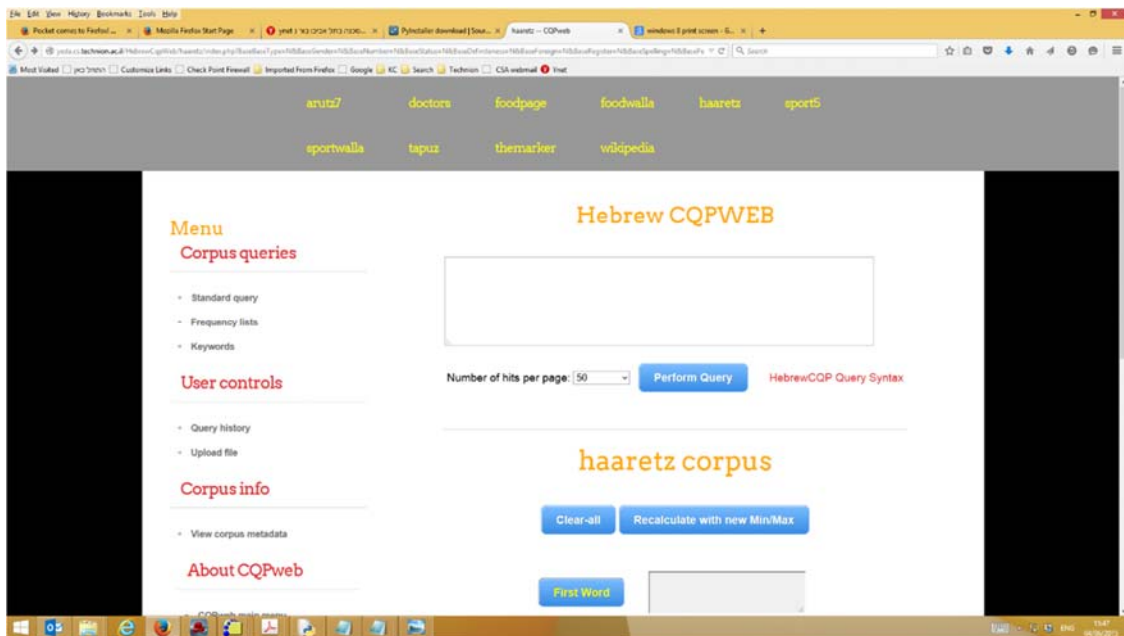
להלן הדגמה:

כניסה לאתר החיפוש דרך אתר מילה

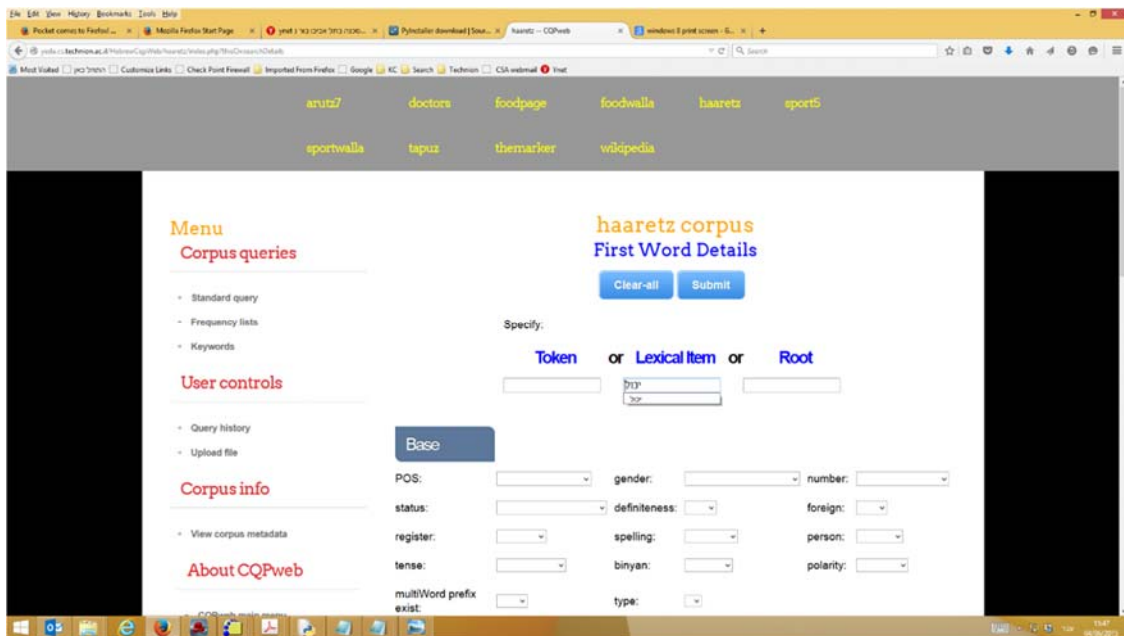
דף הפתיחה לחיפוש



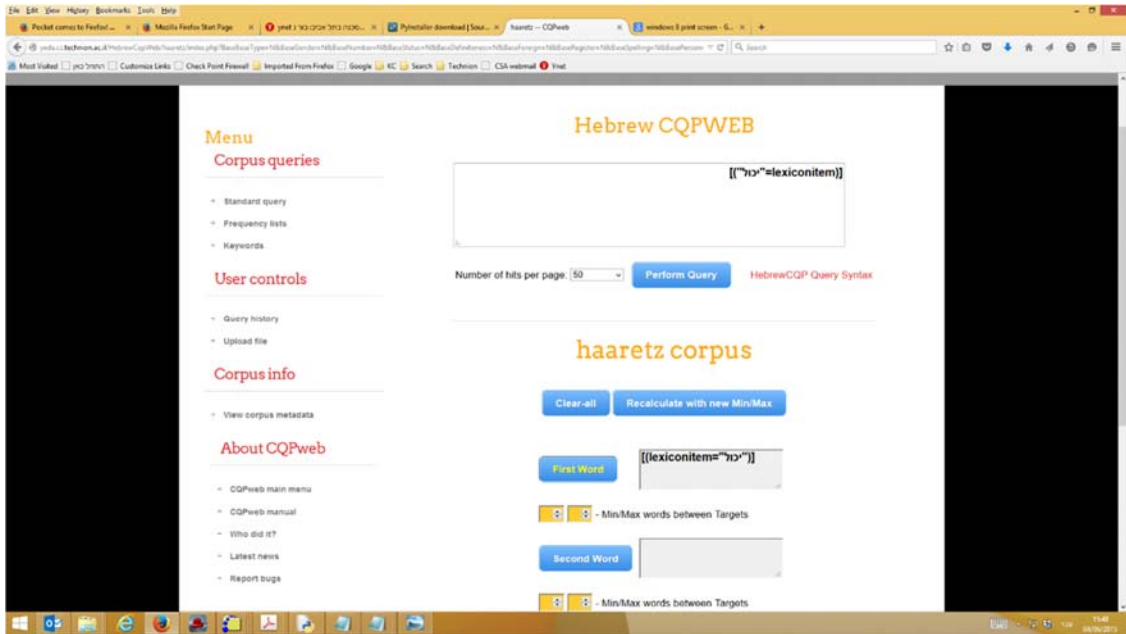
יש כמה קורפוסים. בדוגמא זו נבחר על אתר "הארץ". הקלקה על haaretz (למעלה בצהוב) תוליך אותנו לדף הבא:



נבחר את המילה הראשונה בשאלתה. הקלקת First word תוביל למסך הבא:

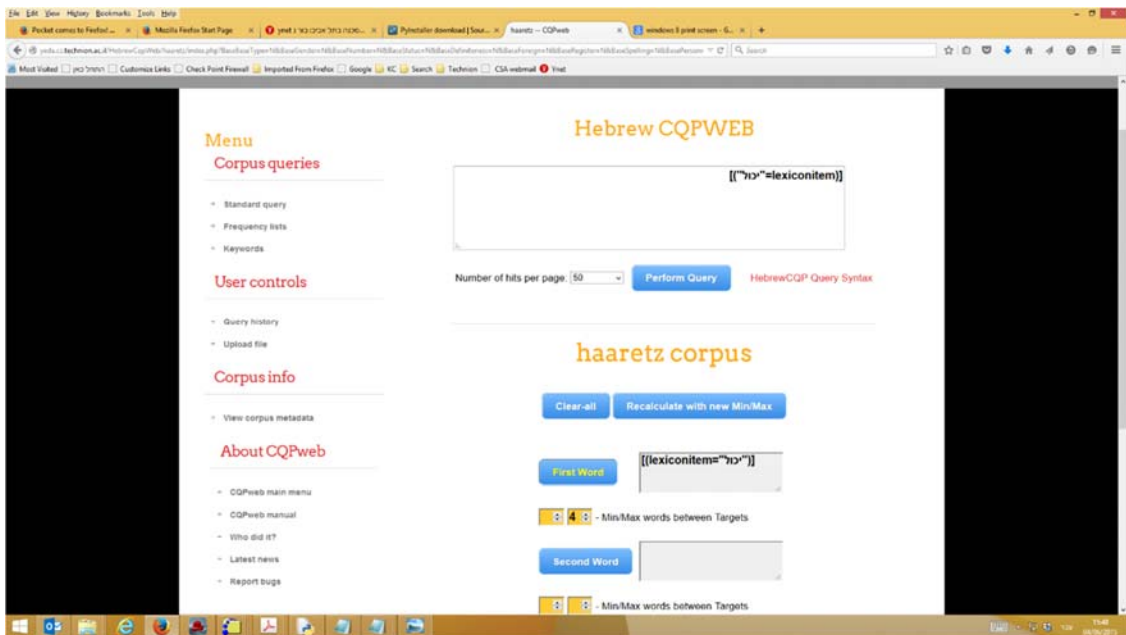


במסך זה סימנו את "יכול" כ-lexical item. לחיצה על Submit תוביל אותנו:

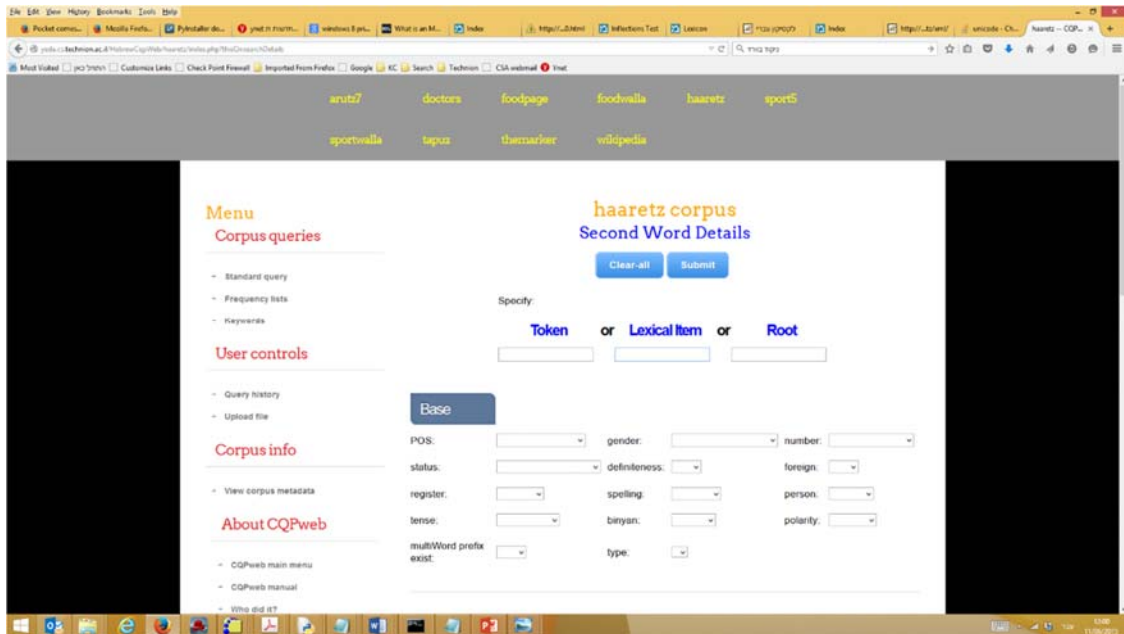


בתיבה העליונה אנו רואים את השאלתה בפורמט של CQP.

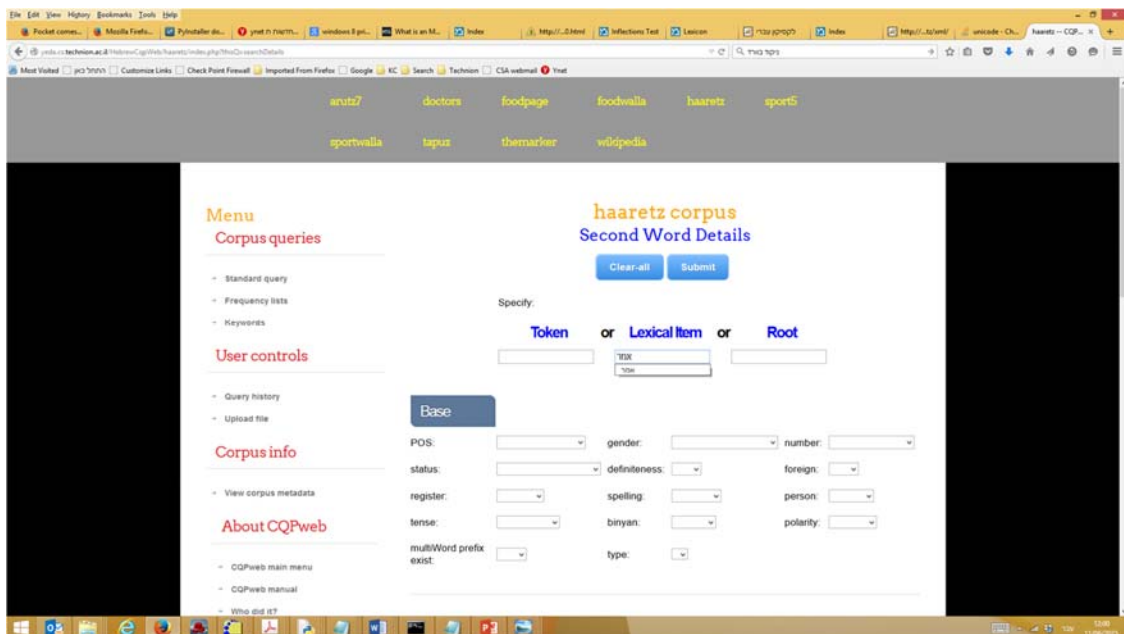
כדי לאפשר למילה השנייה בשאלתה להופיע עד 4 מילים אחרי הראשונה. נציין Min-Max בערך 4. כאן ציינו את max (אם היה ברצוננו שהרווח יהיה לפחות x מילים היינו משנים את המינימום). כיון שלא שינינו את המינימום אין הגבלה, כלומר המילה השנייה יכולה להופיע מיד אחרי הראשונה, אחרי מילה אחת, ... אחרי 4 מילים.



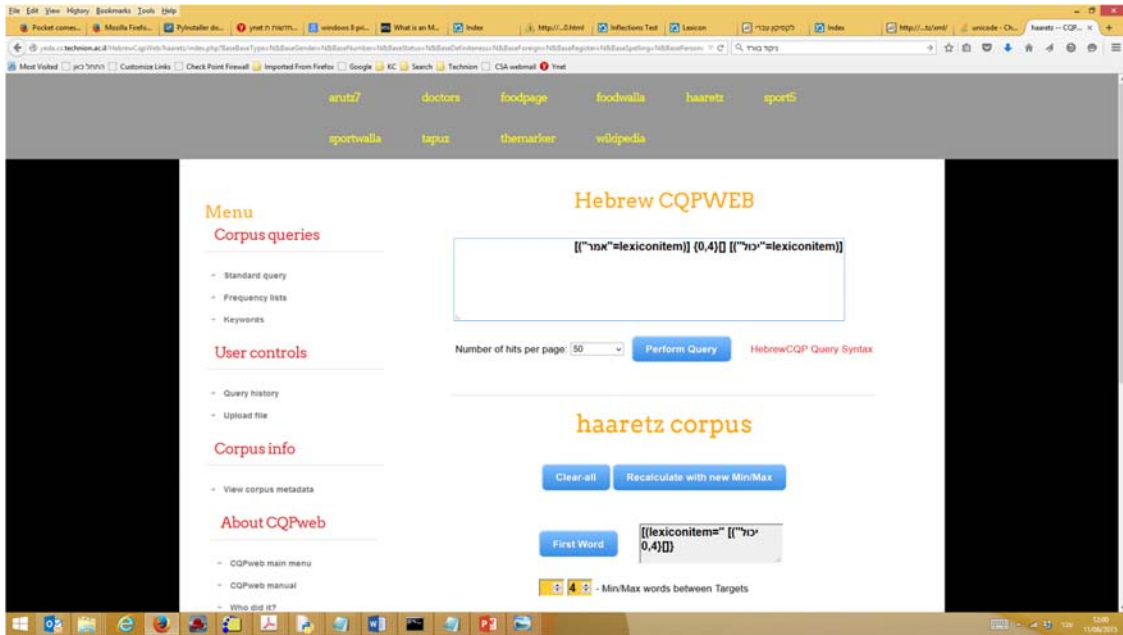
נעבור עתה למילה השנייה (לחיצה על Second Word)



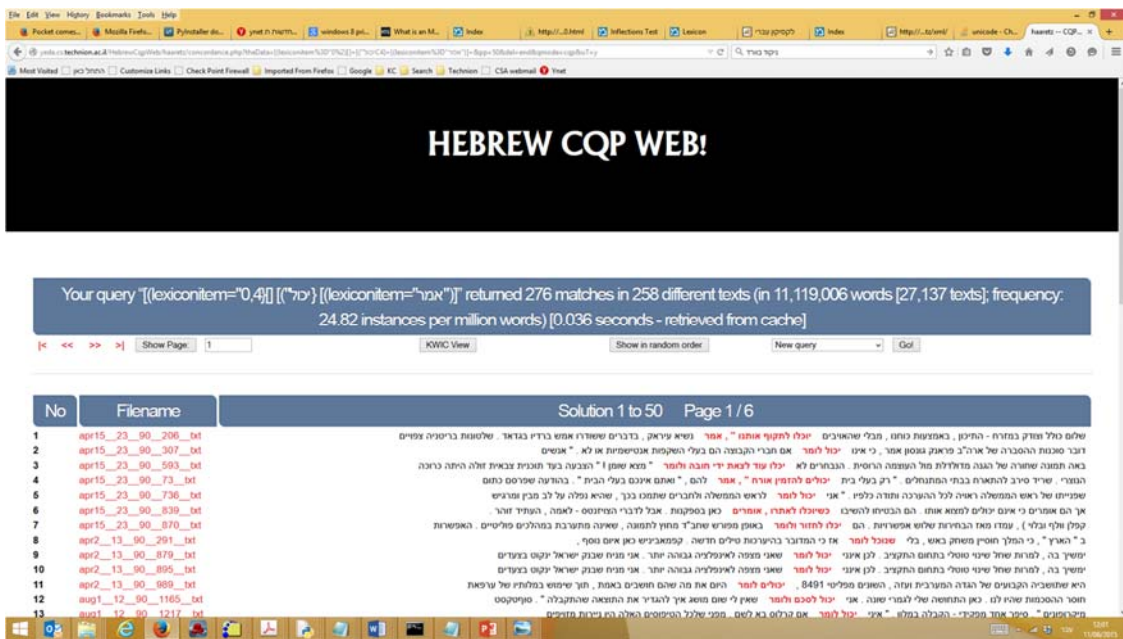
כאן נבחר את הערך הלקסיקלי "אמר"



נלחץ על submit ונראה את השאילתה בפורמט CPQ

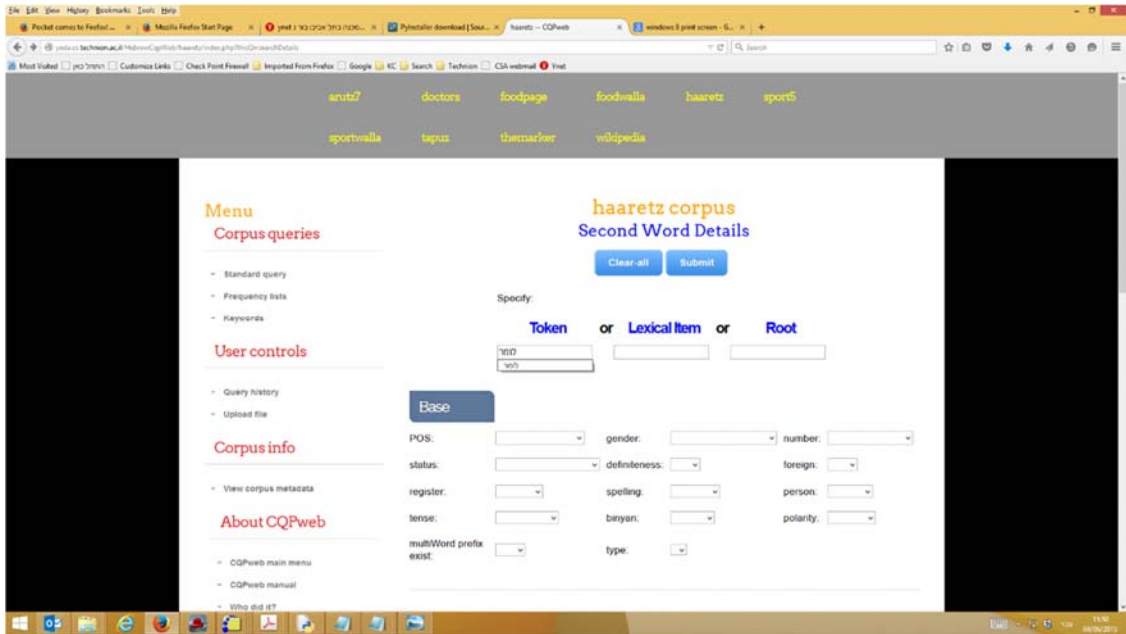


לחיצה על perform query תיתן לנו את תוצאות השאלתה:

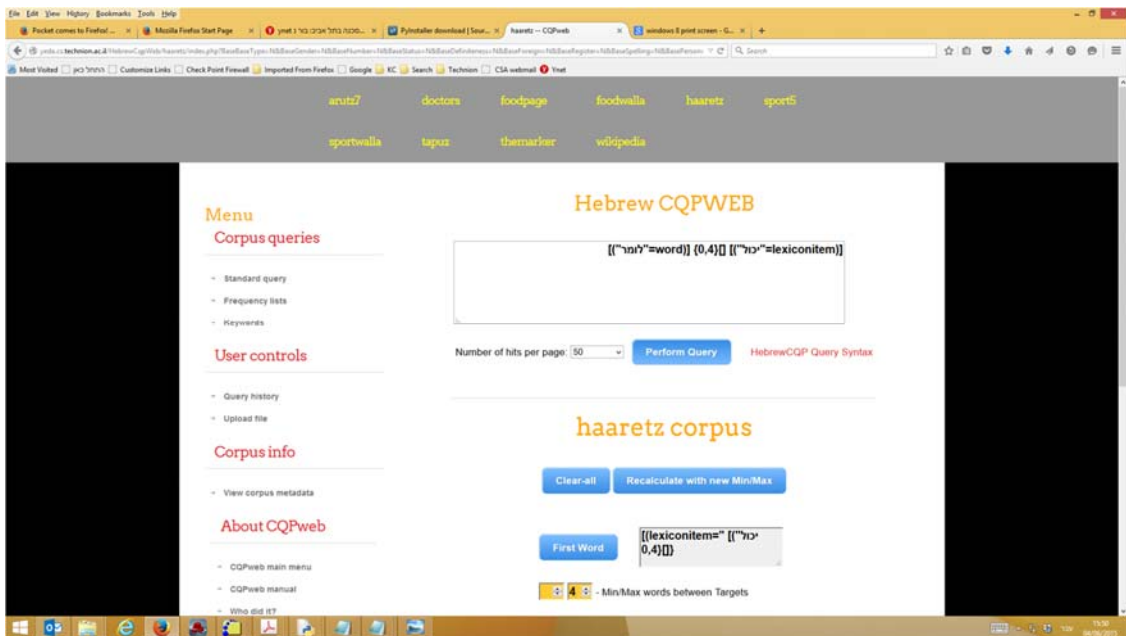


אנו רואים שכל אחת משתי המילים ("יכול" ו"אמר") יכולות להופיע בכל אחת מהנטייות שלהן.

עתה נרצה להגביל את המילה השנייה רק לצורה "לומר". לשם כך נחזור לתפריט המילה השנייה, נגקה אותה (clear) ונרשום "לומר" בתיבה Token.



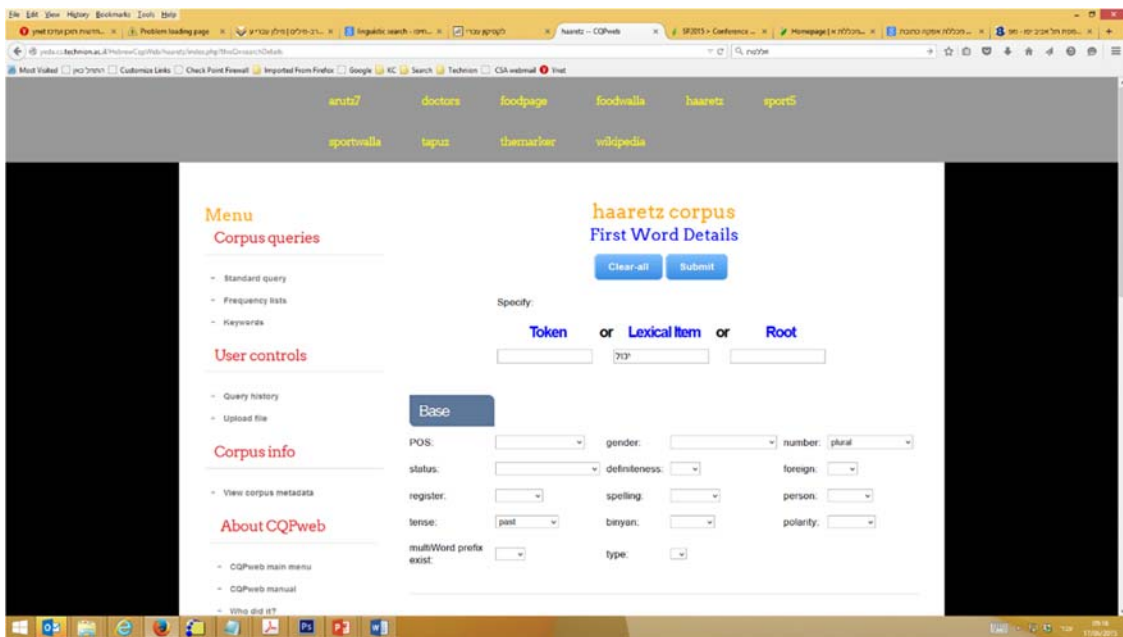
הקלקה על submit תוביל אותנו למסך הבא שבו רואים שהשאלתה מתייחסת למילה השנייה בתוך word ולא Lexical Item.



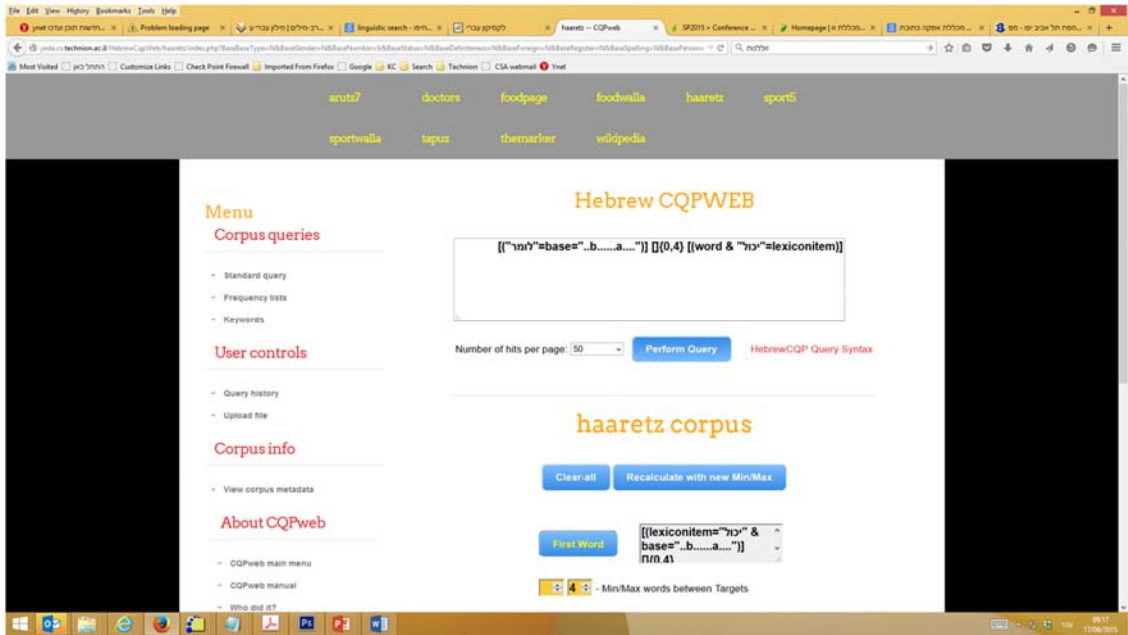
לחיצה על perform query תראה לנו את תוצאות החיפוש. בעוד שהמילה הראשונה מופיעה בנטיית (יכול, יכולים) כל מופעי המילה השנייה הם "לומר".



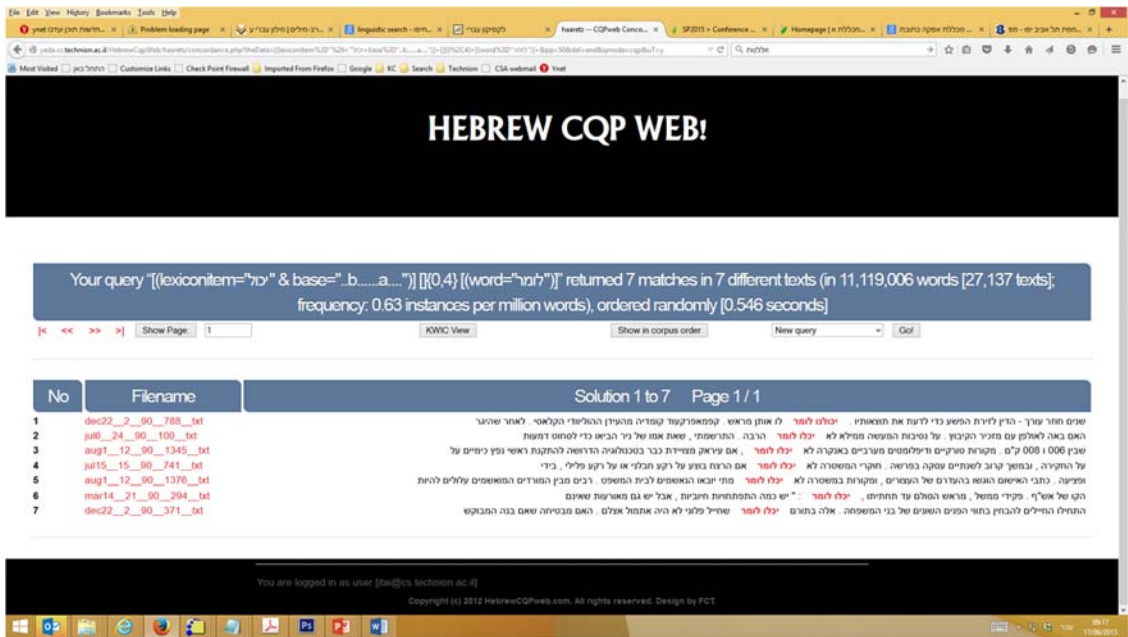
לבסוף, נגביל את הערכים של המילה הראשונה להטיות של עבר ורבים. לשם כך נחזור לתפריט של המילה הראשונה ונבחר ו-tense = past ו-number = plural.



הקלקה על submit תוביל ל-

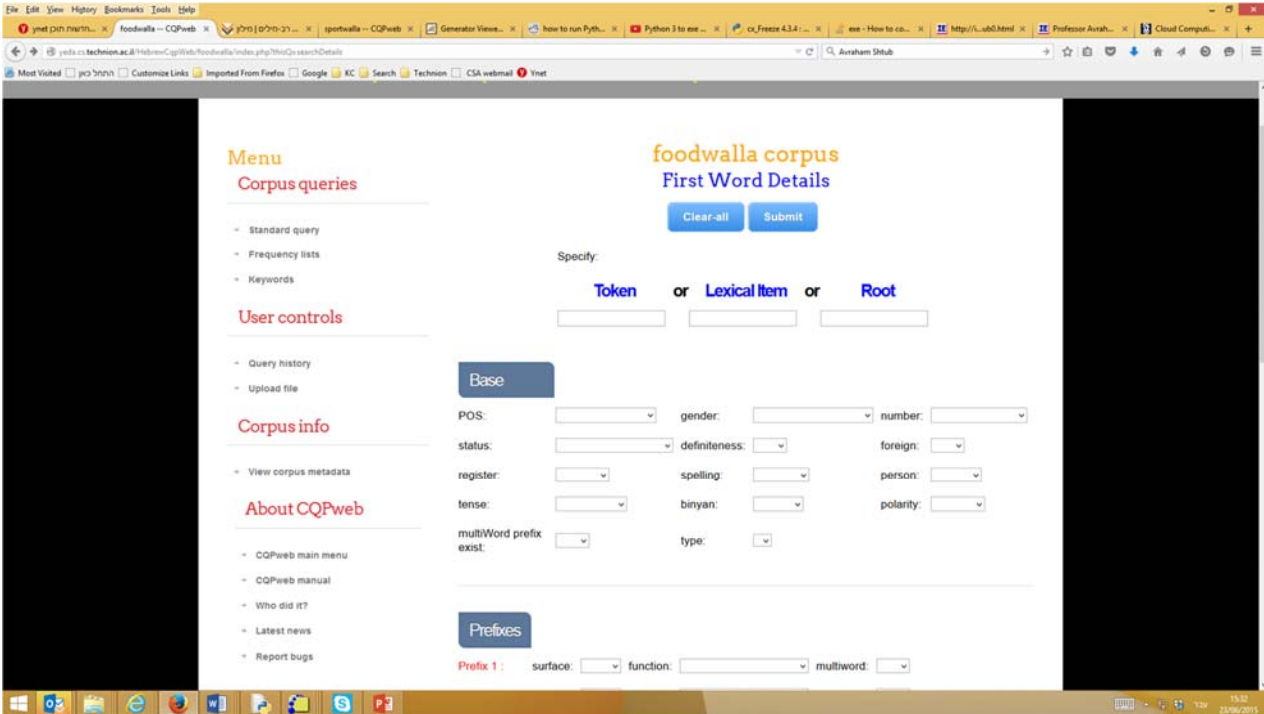


כאן אנו רואים שלמילה הראשונה נוספה דרישת base. (זה לא נראה טוב בגלל שמאל ימין).
הקלקה על perform query מבצע את השאילתה המתוקנת ומראה שבקורפוס הארץ יש 7 משפטים שממלאים את תנאי השאילתה. בצד שמאל מופיעים שמות הקבצים שבהם נמצא כל מופע.

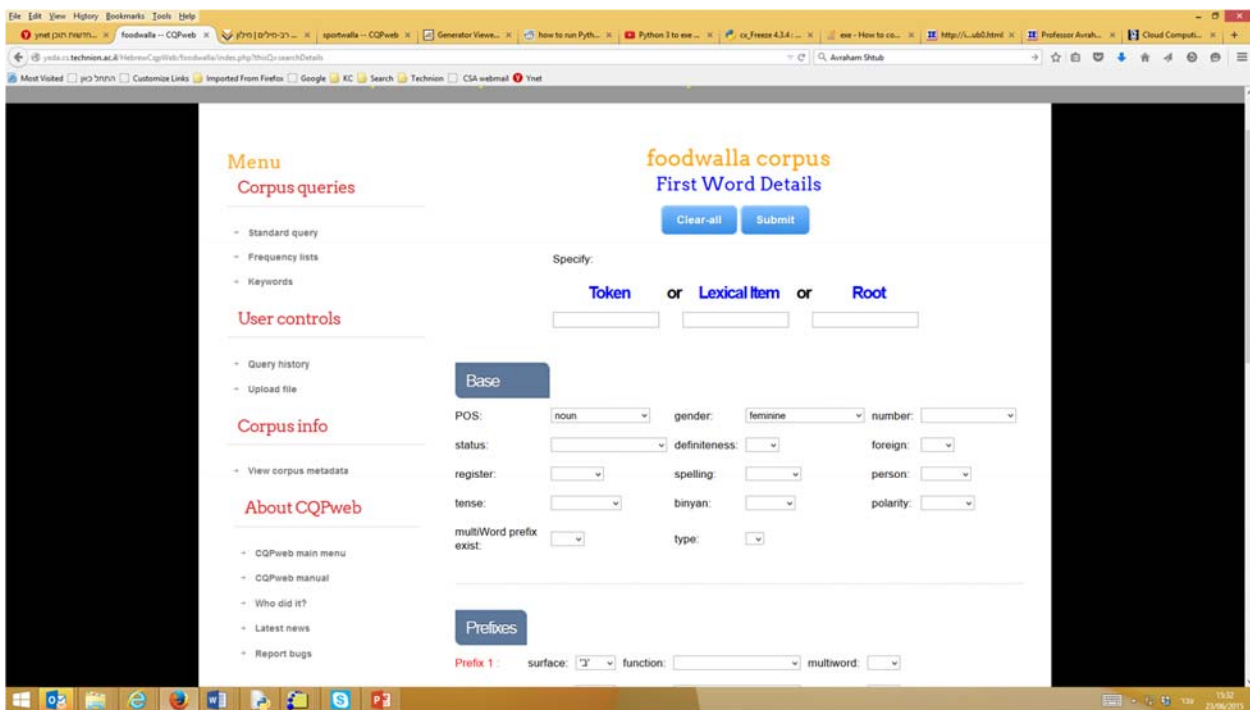


דוגמא 2:

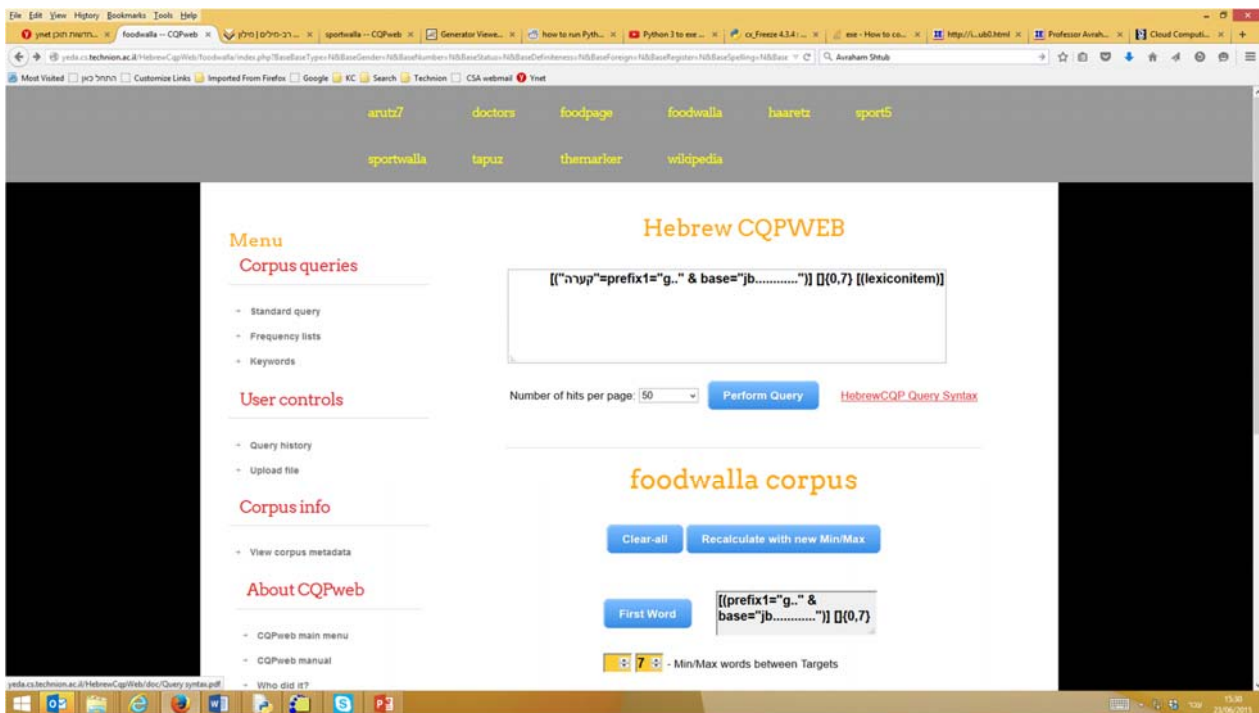
בבחר קורפוס אחר: foodwalla



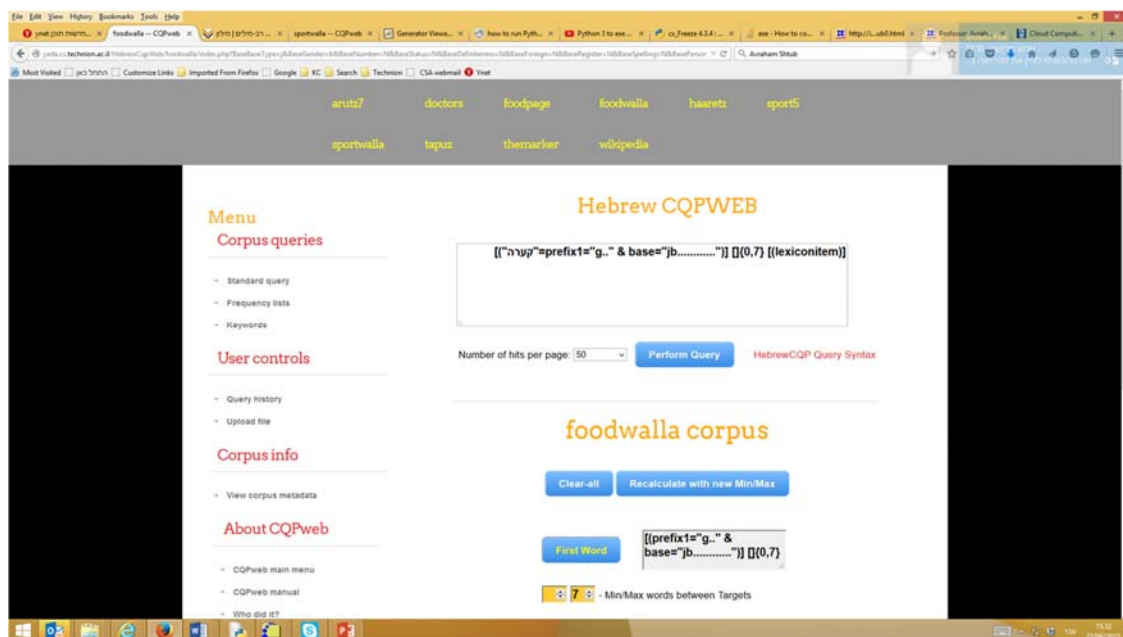
בקורפוס זה, המילה הראשונה תוגבל לשם עצם בנקבה, שמתחילה במילית השימוש "ב":



המילה השנייה תבוא אחרי רווח של עד 7 מילים, ותוגבל לערך הלקסיקלי "קערה"



בצע submit



לאחריה perform query ונקבל:

