

Hspell

בודק איות עברי חופשי

<http://hspell.ivrix.org.il/>

נדב הראל
דן קניגסברג

מועדון הלינוקס החיפאי
1 פברואר, 2010

מהו Hspell

- בודק איות עברי.
- גם מנתח צורני.
- רישיון חופשי (GPL).
- אוצר מילים נרחב: 23,486 מילים (גרסה 1.1).
- מקפיד על תקן הכתיב חסר הניקוד של האקדמיה ללשון העברית.
- נפוץ: משמש בהפצות לינוקס, אופן-אופיס, פיירפוקס, Gmail, ועוד.
- גרסה 0.1: 12/02. גרסה 1.1: 12/09.

מטרות ההרצאה

- סיכום 7 שנים מאז הגרסה הראשונה של Hspell
 - כיצד הוא פותח
 - כיצד הוא פועל
 - איך שולב בהפצות ובאפליקציות
 - תקן האיות
- **סקירת החלטות שתרמו להצלחת המיזם (=פיתוח מהיר, איכות, ושימוש נרחב).**
- מעבר לבדיקת איות – הווה ועתיד.

רקע היסטורי

- 1983: "תכנה חופשית" (ריצ'רד סטולמן)
- 1991: לינוקס, קרנל חופשי ראשון, משלים מערכת הפעלה חופשית ראשונה.
- המשך שנות ה 90: נוחות, שימושיות, שפות:
 - הפצות לינוקס
 - Unicode
 - UTF8
 - Mosaic, Netscape, Mozilla
 - KDE
 - GNOME
 - Star Office

ראשיתו של hspell

- ינואר 2000: תמיכה עברית מתוכנת בפרויקטים הגדולים (דפדפן, DE, מעבד תמלילים)
- 12 בינואר, 2000 – הצורך בבודק איות עלה לראשונה ב `ivrix-discuss`
- באותו יום: רשימה של 1400 מילים (נטויות) מהתזה שלי.
- 21 בינואר: נטיית ש"ע אוטומטית
- 25 בינואר: דן תורם נטיית פעלים אוטומטית

החלטות ראשונות – מה לא לעשות

- אין טעם להתחנן שחברה מסחרית תשחרר בודק איות או רשימת מילים.
- לא מוסרי להעתיק רשימת מילים.
- אין טעם באיסוף אוטומטי של מילים ממסמכים:
 - לעולם לא יהיה שלם.
 - לא יהיה עקבי.
 - ייאספו גם שגיאות כתיב ומילים לא מילוניות.
- להתמקד בבדיקת איות ללא ניקוד. לא להתפזר עם ניקוד, פירושים ותרגומים – ל"רב מילים" זה דרש מאה שנות אדם!

כיצד Hspell כמעט נכשל

- אחרי פרץ עבודה ראשוני, כמעט התייאשנו:
 - דיונים ברשימת דיוור הראו פחות התלהבות ויותר ניסיונות לשכנענו לשאוב רשימת מילים מבודק איות קיים, או לבקש תרומה של קוד מוכן.
 - הגענו למסקנה ש ispell אינו טוב מספיק לצרכינו.
 - הפיתוח נעצר לשנתיים.
- עד שהפצרות שני אנשים שכנעו אותנו עד כמה המיזם חשוב, וחזרנו לעבוד עליו.
- מסקנה: גם הערכת המשתמשים, וגם עניין אישי - שניהם חשובים לתכנה חופשית!

התלבטות ראשונה

ניתוח בזמן ריצה (גישה "אנליטית")	רשימת מילים "טפשה" ("גישה סינטטית")		
איסוף מילות בסיס ורמזי נטייה	איסוף מילות בסיס ורמזי נטייה		שלב פיתוח
-	הטיה אוטומטית של מילים אלו (לא כולל תחיליות)		שלב קומפילציה
ניסיון לנחש את מילת הבסיס של המילה הנתונה, ואז בדיקה האם הנטייה הנ"ל באמת חוקית	בדיקת כל מילה ברשימת המילים החוקיות, אולי בתוספת תחיליות מש"ה וכל"ב בהתאם לסוג המילה		שלב ריצה (בדיקת איות)
צריכת זיכרון ודיסק נמוכה, בקלות	אפשרות שילוב בבודקי איות מבוססי-רשימות קיימים. קוד פשוט יותר.		יתרונות

התלבטות ראשונה (המשך)

- החלטנו: **הגישה הסינתטית.**
- יצירת רשימות מילים:
- איסוף ידני של מילות בסיס ורמזי נטייה
- הטיה אוטומטית בעזרת תכניות perl
- בדיקת איות:
- Ispell לא תמך כראוי בהוספת תחיליות
- גרסה 0.1: בודק איות בפרל, אטי זולל זיכרון.
- מגרסה 0.6 (אוגוסט 2003): בודק איות יעיל ב C.
- מגרסה 0.8 (יוני 2004): **בדיקת איות גם עם aspell.**

התלבטות ראשונה (המשך)

- בדיקת איות היום:

- פקודה hspell, וספרייה, יעילות ובמסורת Unix .

- בכל זאת, כמעט אף אחד לא משתמש בהן:

- יצרנו רשימות מילים עבור hunspell, myspell, aspell.

- מילון עברי כזה נכלל ברוב הפצות הלינוקס. מנגנון

- myspell קיים גם ב Mac OS X Snow Leopard.

- פלג-אין בודק איות עברי באופן-אופיס ופיירפוקס משתמש

- במילון עברי דרך hunspell או myspell.

- כך גם (כנראה) עושה gmail.

- **גישת רשימת המילים: בחירה נכונה וחשובה**

השוואת ביצועים (בדיקת איות בלבד)

hunspell -d he_IL (double-affix-compressed since Hspell 1.1)	aspell -l he	hspell	
3457 K	8458 K	153 K (2.7 bits per word!)	שטח דיסק (מילון, התאמת תחליות)
0.15 sec	0.01 sec (he.rws mapped)	0.03 sec	זמן טעינת מילון
87,268	436,000	3,800,000	מילים לשנייה (בדיקה לבד, ללא תיקון)
10552 K	780 K + 8396 K	5744 K	צריכת זיכרון pmap -d

דחיסת המילון

חתול חתולי חתולים חתולך כלב כלבי כלבים כלבך

461,000 מילים, 3500 KB.

• דחיסת gzip מוצאת קטעים משותפים. 979 KB.

חתול סי 0ם 2ך כלב סי 0ם 2ך

• מידול זה כבר מקטין ל 1017 KB.

• דחיסת gzip מקטינה ל 110 KB !

תקן האיות

- אנו בודקים איות במלל ללא ניקוד.
- אך מהו "כתיב מלא" נכון?
- נבחר תקן: **הכתיב חסר הניקוד** של האקדמיה ללשון העברית.
- התקן במלואו מופיע כנספח ברוב המילונים, תקציר באתר האקדמיה.
- רק לאחרונה החל להופיע הכתיב חסר הניקוד בערכים במילונים. האם מוזכר בבתי ספר?
- מסמך שלנו "סוגיות בכתיב חסר הניקוד".

הכתיב חסר הניקוד - דוגמאות

- הוספת וי"ו לציון קיבוץ:
חולצה.
- הוספת וי"ו לציון חולם, קמץ קטן וחטף-קמץ:
חודש, חודשי, חודשים.
- אך רק אם היה חולם במילת הבסיס:
חכמה, תכנית, צהריים.
- ובפרט לא אם כלל לא הייתה תנועת ס:
פועל-פועלו, שמר-לשמרו.

הכתיב חסר הניקוד - דוגמאות

- מוסיפים יו"ד לציון חיריק שאין אחריו שווא נח:
סיפור, שנייה, יישום, אבל שמחה, עברית,
וכן **דיוק**.
- במילים לועזיות, היו"ד כבר בכתיב המנוקד:
היסטוריה.
- הוספת יו"ד בנטיות תלויה בצורה היסודית:
זיכרון-זיכרונות, לב-לבי.
אבל עיפרון-עפרונות.
- הוספת יו"ד לציון צירי במקרים מסוימים: **תיבה,**
חירש, אבל לרוב לא: **שרפה, ממד.**

הכתיב חסר הניקוד - דוגמאות

- הכפלת וי"ו עיצורית:
עוול, הוועד, עכשווי, אבל ועד, קו, מצוות.
- הכפלת יו"ד עיצורית:
בניין, הייתה, עליי, קריית-, אבל קריה, מצוין, ילד, הילד, בנאי, בית.

רשימת המילים ובנייתה

- מה יש ברשימת המילים?
- **שמות עצם.** הטיה אוטומטית ליחיד, רבים, נפרד, נסמך, כינויים.
דוגמה: **כלב:**

כלבים כלבי- כלביי
כלבינו כלביך כלבייך
כלביכם כלביכן כלביו
כלביה כלביהן כלביהם

כלב כלב- כלבי כלבנו
כלבך כלבך כלבכם
כלבכן כלבו כלבה
כלבן כלבם

רשימת המילים ובנייתה (המשך)

- **תארים.** הטיה אוטומטית ליחיד, רבים, נסמך, זכר, נקבה.
דוגמה: **ירוק:**

ירוקה ירוקת- ירוקות
ירוקות-

ירוק ירוק- ירוקים ירוקי-

רשימת המילים ובנייתה (המשך)

- **פעלים.** הטיה אוטומטית לבניינים אפשריים, זמנים, שמות פעולה, וכו'.
דוגמה: **שמר**, בניין קל (רשימה מקוצרת):

לשמור	שמור	אשמור	שומר	שמרתי	לשמור
שמירה	שמרי	תשמור	שומרת	שמרת	שמרתיו וכו
	שמרו	תשמרי	שומרים	שמרת	
	שמורנה	ישמור	שומרות	שמר	
		תשמור	שמור	שמרה	
		נשמור	שמורה	שמרנו	
		תשמרו	שמורת-	שמרתם	
		תשמורנה	שמורים	שמרתן	
		ישמרו	שמורי-	שמרו	
			שמורות		

רשימת המילים ובנייתה (המשך)

- **מילים נוספות**, ללא הטיה אוטומטית:
 - מיליות ונטיותיהן (של, שלי, ...)
 - תוארי הפועל (היכן, שלשום, הרבה, ...)
 - שמות גוף (אני, אלה, מיהו, ...)
 - שמות פרטיים של אנשים, מקומות, חודשים, אותיות וכד'.
 - מספרים (שניים, שניהם, שני)

רשימת המילים ובנייתה (המשך)

- גרסה 1.1 כללה נטיות של:

- 12445 שמות עצם

- 3583 שמות תואר

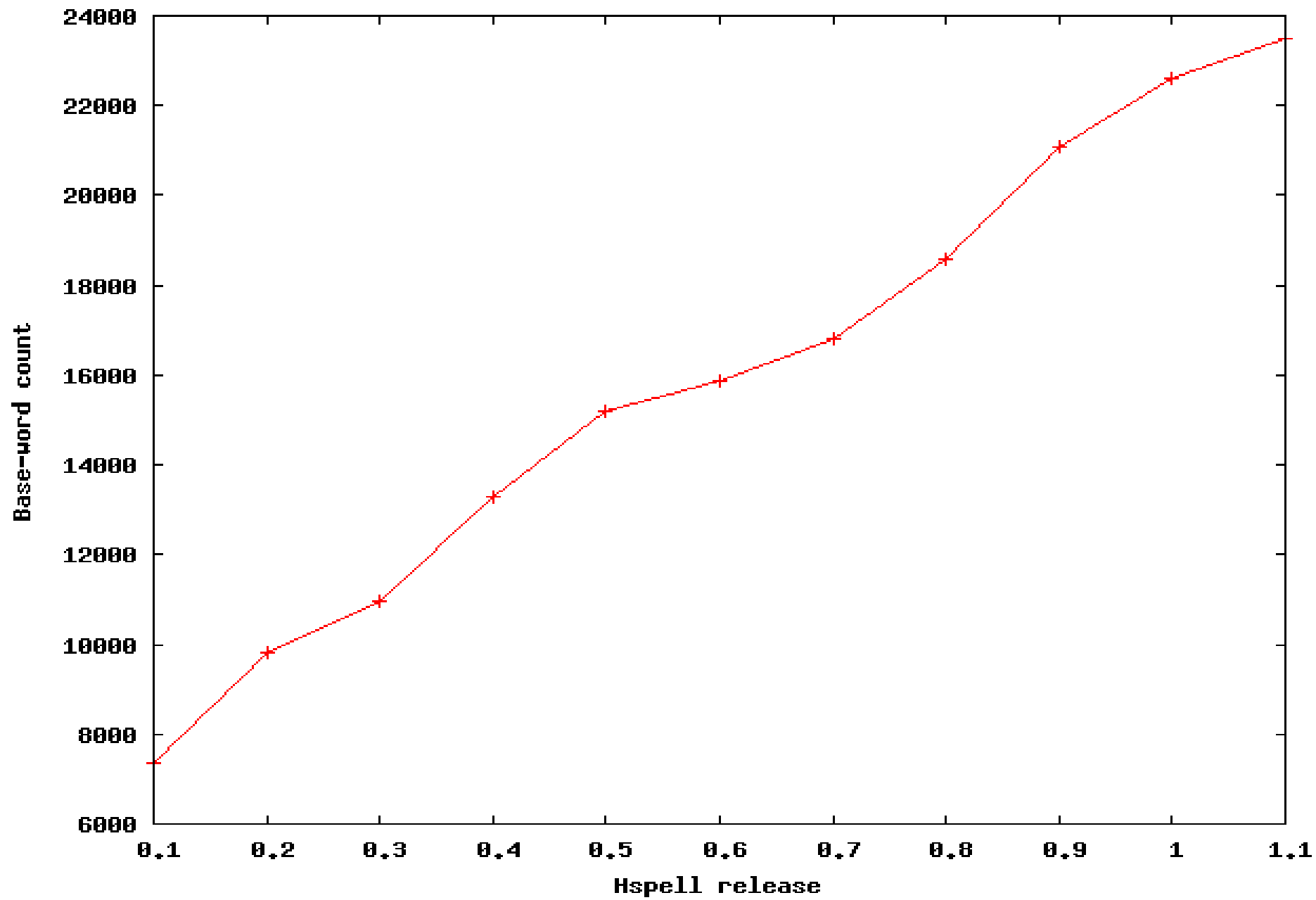
- 5242 פעלים

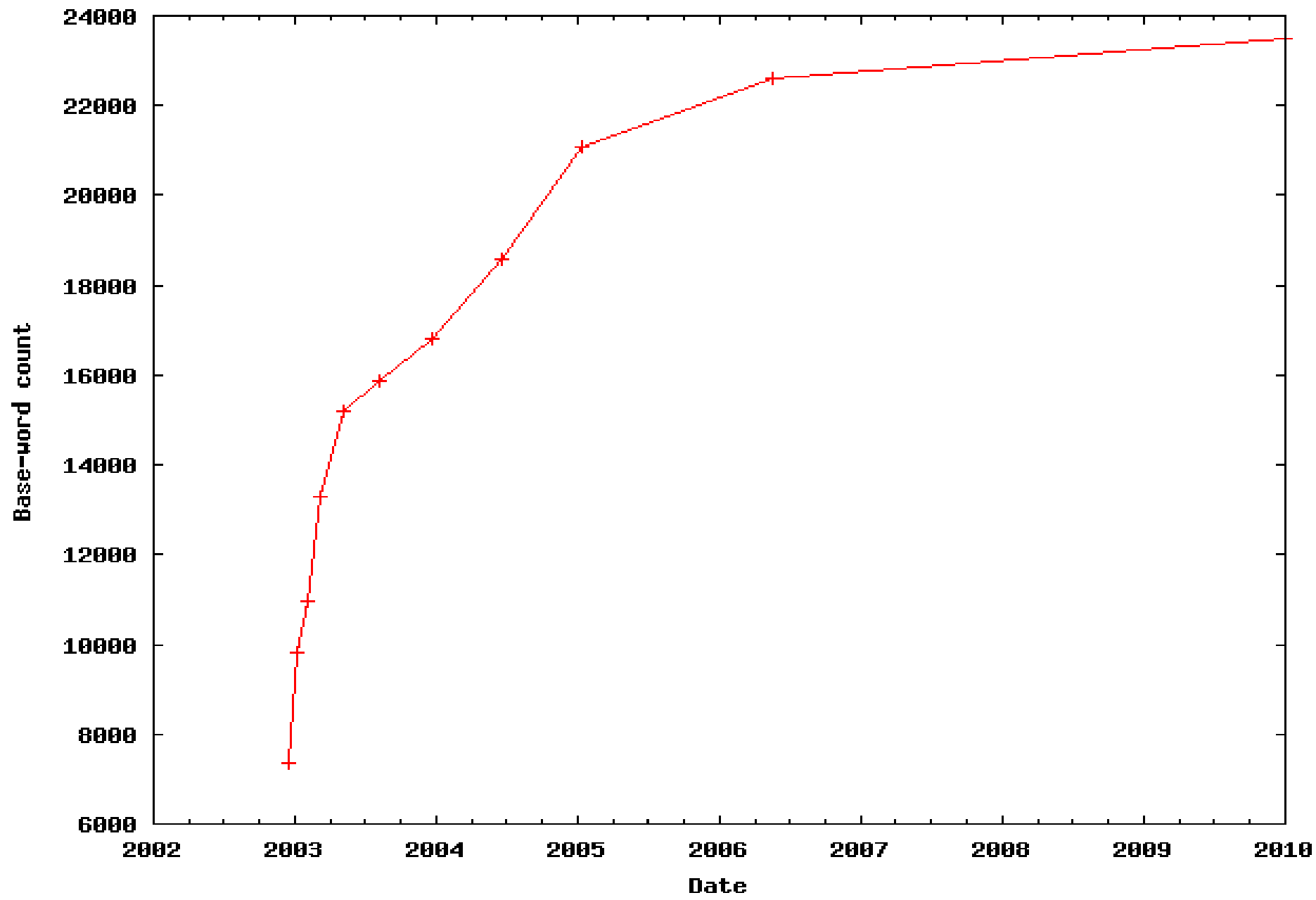
- וכן 2216 מילים נוספות

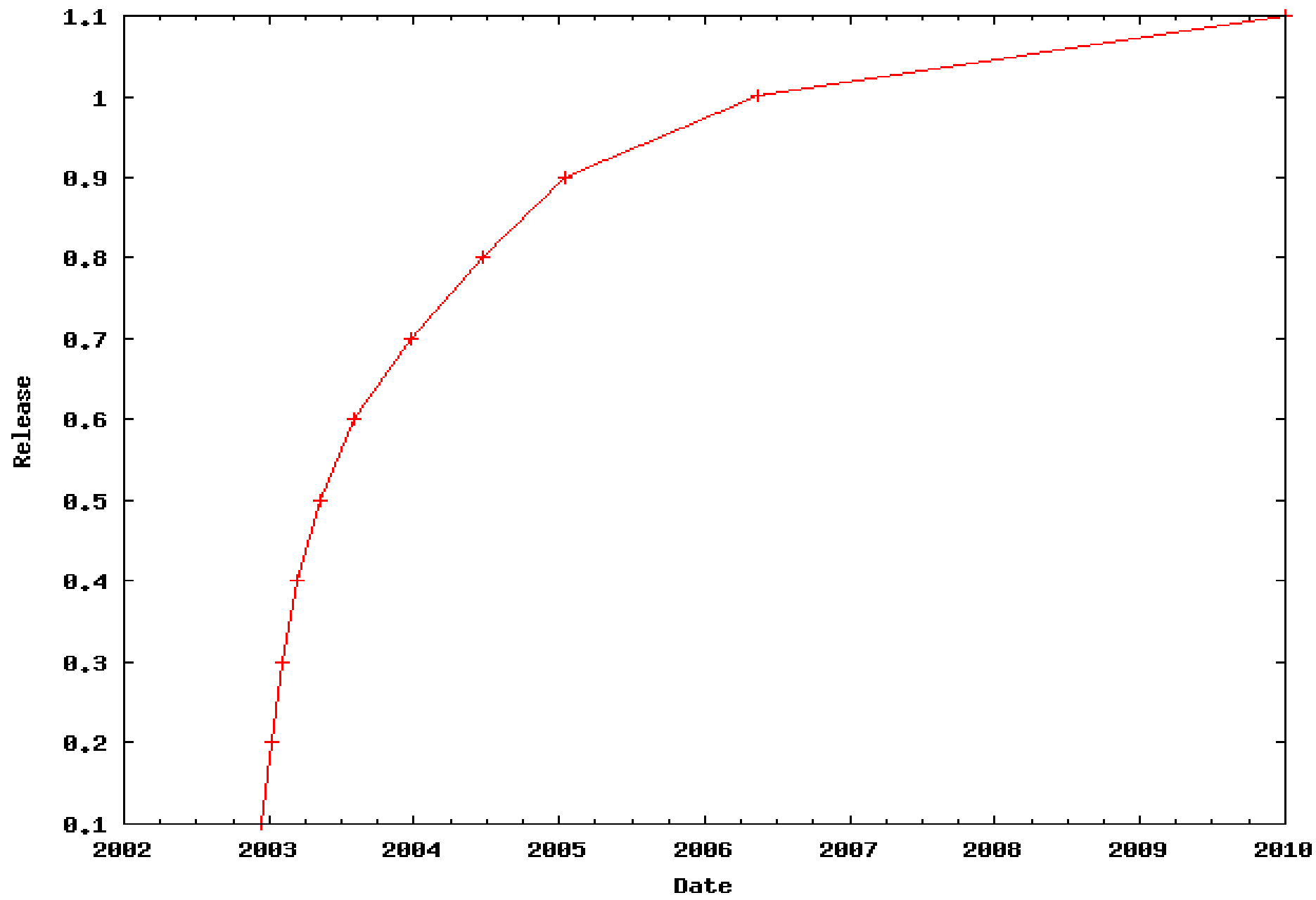
- ובסך הכול:

- 23,486 מילות בסיס

- 461,326 נטיות ומילים







הטיית שמות עצם

- התכנית wolig.pl מקבלת רשימת שמות-עצם ומטה כל אחד מהם, בעזרת רמזים שנתן המשתמש.
- ב"לוח השמות השלם" של ד"ר שאול ברקלי מופיעות כ- 350 צורות נטייה לשמות עצם מנוקדים.

הטיית שמות עצם (המשך)

- למזלנו רוב הצורות מתבררות כזהות ללא ניקוד.
למשל

כלב-כלבי דגל-דגלי

- שתי אלו מוטות באותה צורה ע"י שורת הקלט:
כלב ע

- והפלט הוא

כלב כלב- כלבי כלבנו כלבך כלבך כלבכם
כלבכך כלבו כלבה כלבן כלבם כלבים כלבי-
כלביי כלבינו כלביך כלבייך כלביכם כלביכך
כלביו כלביה כלביהן כלביהם

הטיית שמות עצם (המשך)

- לעתים חייבים רמזים לנטיות הנכונות, למשל לגבי צורת הריבוי.

- השוו: קוף-קופים עוף-עופות

- אן: שירות-שירותים חירות-חירויות

- לשם כך חייבים לספק רמזים:

עוף, עות

שירות, עים

- והפלט למשל: עוף עוף- עופי עופנו עופך עופך עופכם עופכן עופו עופה עופן עופם **עופות** עופות- עופותיי עופותינו עופותיך עופותיך עופותיכם עופותיכן עופותיו עופותיה עופותיהן עופותיהם

הטיית שמות עצם (המשך)

• צורות ריבוי נוספות:

• משנה ע, יות

• קצבה ע, אות

• גרב ע, יים

• עשן ע, יחיד

• בת ע, רבים=בנות

הטיית שמות עצם (המשך)

• ישנם גם שמות-עצם עם כמה צורות ריבוי חוקיות:

- חודש ע,ים,יים
- שעה ע,ות,יים
- קבר ע,ות,ים
- איש ע,ים,רבים=אנשים
- שפה ע,ות,יים,רבים=שפתות

הטיית שמות עצם (המשך)

- התכנית wolig.pl מנסה לנחש את הריבוי הנכון לפי סיומת צורת היחיד:

<u>יחיד</u>	<u>רבים</u>
מלך	מלכים
מלכה	מלכות
כותרת	כותרות
כמות	כמויות
אחריות	-אין-

- כך כ 90% משמות העצם אינם דורשים רמז ריבוי

הטיית שמות עצם (המשך)

- כאמור, למזלנו לא מעניין אותנו אם תנועה כמו סגול, פתח או צירי משתנה בין הנטיות השונות:
 - כלב – כלבים, כלבי
 - ספר – ספרים
 - מטוס – מטוסים
- גם במקרה חיריק, כללי הכתיב חסר הניקוד לרוב לא ידרשו יו"ד:
 - ספר – ספרי
 - עז – עזים
- יוצא דופן המקרה בו ביחיד חיריק, בנטיות נעלם
 - עיפרון – עפרוני, עפרונות. עיפרון ע,ות,אבד_י

הטיית שמות עצם (המשך)

- למרבה הצער, דרושים עוד מספר רמזים.
דוגמאות:

- חנית – חניתות (לא חניות): חנית ע, שמורת

- רובה – רובי (לא רובתי): רובה ע,ים,סגולה

- צומת – צמתים, צומתי: צומת ע,ים,אבדו

- אח, אחיו (לא אחו): אח ע,מיוחד_אח

- שן, שיניים (לא שניים): אח ע,מיוחד_שן

- מגורים, מגורי-, מגוריי...: מגור ע,אין_יחיד

- גברת ע,רבים=גברות,נסמכים=גבירות

הטיית שמות עצם (המשך)

- טיפול אוטומטי בכללי הכתיב חסר הניקוד, תוך שימוש פנימי באותיות עזר:
 - קריה ע: קרעה קריה קרעת- קריית-
 - עירייה ע: עיריעה עירייה עיריעות עירות
 - קו ע: קש קו קשים קווים
 - שואק ע, אבד_ו: שואק שוק ששקים שווקים שואקיכם שוקיכם

בודק איות חופשי

- דרגות חופש אפשריות לבודק איות:
 - קנייני ומסחרי (קוד נעול, תמורת תשלום).
 - חינמי אך קוד נעול (freeware).
 - קוד פתוח, רשימת המילים הסופית נתונה כקובץ.
מבחינת RMS, זה חופשי. לדעתנו לא מספיק:
איך מפתח חדש מוסיף או משנה מילים?
 - קוד פתוח, נתונה רשימת מילות בסיס וקוד נטייה פתוח.
איך יודעים מהו תקן האיות המשותף לכל המילים?
 - קוד פתוח, מילות בסיס וקוד נטייה פתוח, ובנוסף תיעוד מפורט של תקן האיות הממומש. **זו בחירתנו.**

בודק איות חופשי

- כמובן שקוד או תוכן מועתק אינו חופשי, גם אם מדביקים עליו GPL:
- המילון פותח ללא העתקת רשימות מילים ממוצרים אחרים, מילונים, וכד'.

בודק איות חופשי

- החלטנו לשחרר את Hspell תחת ה GPL. האם זו הייתה החלטה נכונה?
- שירותים רשת יכולים להשתמש ב Hspell כי הם לא מפיצים תכנה. לדוגמה, gmail.
- האם ואיך מגן ה GPL עד המילון? אנחנו יודעים על לפחות מקרה אחד בו המילון שלנו (אך לא הקוד) הועתק בניגוד לתנאי ה GPL.
- תכנות שאינן GPL כמו openoffice לא יכלו להכניס את hspell. שימוש ב aspell עקף את הבעיה.
- באותו אופן, יכולה מערכת הפעלה קניינית להכיל בודק איות עברי מבוסס hspell.

מעבר לבדיקת איות – הווה ועתיד

- עד עתה דיברנו על בדיקת איות.
- מהן היכולות הלשוניות הנוספות שקיימות ב hspell כבר היום?
?
- מהן היכולות הלשוניות החסרות ושאולי נצטרך בעשור הבא?
?

מנתח צורני (מורפולוגי)

- התכנית שמייצרת את נטיות "כלב" יכולה לזכור בדיוק מדוע יצרה כל נטייה:

מין שמות עצם קיים למרות
שלא דרוש לבדיקת איות



כלב ע,ז,יחיד
כלב ע,ז,יחיד,סמיכות
כלבי ע,ז,יחיד,של/אני
...
כלבים ע,ז,רבים
....

- כך בהינתן מילה בטקסט (נטויה, אולי עם תחיליות), נוכל לדעת את הצורה או הצורות האפשריות לנתח אותה.

מנתח צורני (מורפולוגי):

- דוגמאות מ -al -hspell:

משטרה	הרכבת	כלבים
<p>משטרה: משטרה (ע,נ,יחיד) משטר (ע,ז,יחיד,של/היא) מ+שטרה: שטר (ע,ז,יחיד,של/היא)</p>	<p>הרכבת: הרכיב (פ,נ,2,יחיד,עבר) הרכיב (פ,ז,2,יחיד,עבר) הרכבה (ע,נ,יחיד,סמיכות) ה+רכבת: רכבת (ע,נ,יחיד) ה+רכבת: (ה"א השאלה) רכב (פ,נ,2,יחיד,עבר) רכב (פ,ז,2,יחיד,עבר) רכבת (ע,נ,יחיד)</p>	<p>כלבים: כלב (ע, ז, רבים)</p>

מה עוד חסר?

- יכולות חדשות שלדעתי נצטרך בעשור הקרוב בעולם התכנה החופשית: (סדר חשיבות יורד)
 - תכנת הקראה (TTS) של טקסט עברי לא מנוקד.
 - בודק איות שמזהה גם מילים נכונות בהקשר לא נכון. לדוגמה, "זהו כתיב לא תיקני".
 - תכנת חיפוש טקסט עברי מדויקת, שמאפשרת למצוא את המילה שהתכוונת אליה (למשל רכבת, לא רכבת).
 - (עדיפות נמוכה) תכנת ניקוד אוטומטי
 - (עדיפות נמוכה מאוד) בודק איות לכתיב מנוקד.

מה עוד חסר? (המשך)

- איך לעשות זאת? לדעתי:
- תכנה לניתוח דקדוקי בהינתן ניתוח מורפולוגי.
- שיפור דיוק על-ידי לקסיקון עם מחלקות סמנטיות (semantic frames). לדוגמה: "שתיתי חלב".
- הוספת ניקוד – או לפחות תנועות – ללקסיקון ולתכנת הנטייה של hspell.
- הכרת וריאנטים נפוצים (לוא דווקא נכונים) למנתח המורפולוגי של hspell.
- תכנת הקראה של טקסט עברי עם תנועות.
- לחלק מהיכולות לא דרושים כל השלבים.

חוסרים קטנים

- שיפור פיצול מילים ב openoffice (ראשי תיבות)
- כיווץ עץ המילים בזיכרון: מ 5 MB לכיוון 150 KB.
- יאפשר גם mapping במקום קריאה, ו embedded
- שיפור אלגוריתם הצעת התיקונים של hspell.
- שיפור הצעות התיקונים ב hunspell.
- עבודה על יעילות hunspell.
- ייעול המנתח הצורני (כרגע צורך 17 MB).
- שילוב המנתח הצורני ב hunspell.