

## האילמה ילוקוטורפ רוזחא: Task # 1428 - החותפ תסנכ

<b>Status:</b>	New	<b>Priority:</b>	Normal
<b>Author:</b>	בירק םדא	<b>Category:</b>	מינותג תריצק
<b>Created:</b>	10/26/2010	<b>Assignee:</b>	Ori Hoch
<b>Updated:</b>	01/21/2013	<b>Due Date::</b>	
<b>הנושאר הלטמכ מיאתמ:</b> No			
<b>Description:</b> הדעוו לכ םע םויה השעגש המל המודב			

### History

12/24/2012 07:29 pm - Ori Hoch

- Assignee set to Ori Hoch

- Category set to תריצק מינותג

- הנושאר הלטמכ מיאתמ set to No

need to get the protocols from this url:

[http://www.knesset.gov.il/plenum/heb/plenum\\_search.aspx](http://www.knesset.gov.il/plenum/heb/plenum_search.aspx)

and then parse the doc files

12/29/2012 05:35 pm - Ofir Carny

- File parse.py added

Attached script contains python retrieval code, openoffice based parsing and HTML conversion, and some post conversion HTML parsing and Hebrew unicode handling

01/07/2013 10:08 pm - Ori Hoch

- File 03581012.doc.awdb.xml added

- File 04085512.doc.awdb.xml added

- File 18\_ptm\_219371.doc.awdb.xml added

attached: sample xml files produced using antiword

wiki page for planning the parsing process of these xml files:

[[<https://github.com/astupidog/Open-Knesset/wiki/parse-plenum-protocols>]]

01/13/2013 03:38 pm - Ori Hoch

- File plenum\_htmls.zip added

parsing process is finished and working but there are probably many bugs..

attached are html files that were parsed, hopefully someone will look over them and provide some bug reports

the plenum\_htmls.zip file should be extracted somewhere, then you should open the index.html file using a web browser (tested with chrome but should

work on other browsers as well)

#### 01/21/2013 08:45 pm - Ori Hoch

רשמאל תודוקג:

1. קר הלה. "האילמ תדעו"כ וא "הדעו"כ הלא נמסיש הדעו לכל "type" הדש פסוול רשאכ תודעו לש תואלבטו מילדום מלאב שמתשנ מינונה דסמב מילדום תניחבמ. תחא "האילמ תדעו".
2. הנוש parsing עצבת זאו "האילמ תדעו"ב רבודמ מא קודבתש כ create\_protocol\_parts לש היצקנופל יאנת תפסוה.
3. 'וכו "רבוד אלל טסקט", "תרתוכ", "רבוד": גוס לש הדש פסוול שאכ protocolParts ודבועי מיקלחה לכ.
4. תסנכה רלאמ לוקוטורפה תא בוש דירוהל ילב re-parse עצבל לכונש כ לוקוטורפה לש יקנה טסקטה תא רומשל שמשל protocol\_text הדשה.

#### Files

parse.py	7.8 kB	12/29/2012	Ofir Carny
18_ptm_219371.doc.awdb.xml	459.9 kB	01/07/2013	Ori Hoch
03581012.doc.awdb.xml	414.8 kB	01/07/2013	Ori Hoch
04085512.doc.awdb.xml	249.7 kB	01/07/2013	Ori Hoch
plenum_htmls.zip	2.9 MB	01/13/2013	Ori Hoch