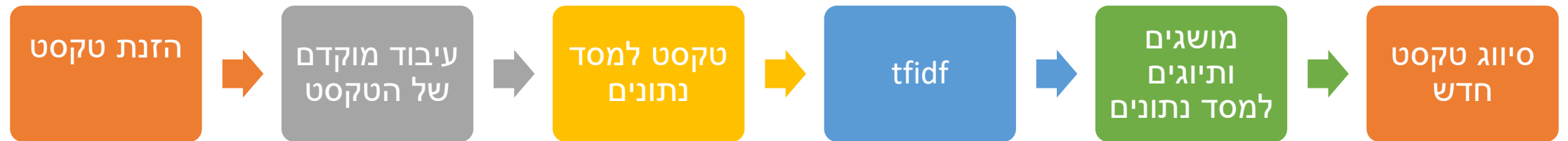


# עבודת גמר: מסוג נושאים

יישום tf-idf

# מבט מלמעלה



# צוותים

- צוות סריקת מידע - 2 או 3 חברים
- צוות עיבוד מקדים - 3 או 4 אנשים
- צוות מסד נתונים - 6 או 7 אנשים
- לשקול חלוקה ל: (1) צוות שאילתות (2) צוות חיבור בפייתון
- צוות גרעין המערכת - 5 או 6
- לשקול חלוקה ל: (1) צוות TFIDF (2) צוות סיווג טקסט חדש

# הזנת טקסט

input: url address, label

output: text (not html) , label

צוות: סריקת מידע

בוחרים מאמר מויקיפדיה  
באנגלית ומעתיקים ידנית את  
הכתובת  
מקבלים גם סיווג כמו  
"כדורסל"



השתמשו בחבילות  
"requests" ו "html2text"  
תחלצו רק טקסט בלבד!  
שישמר במשתנה



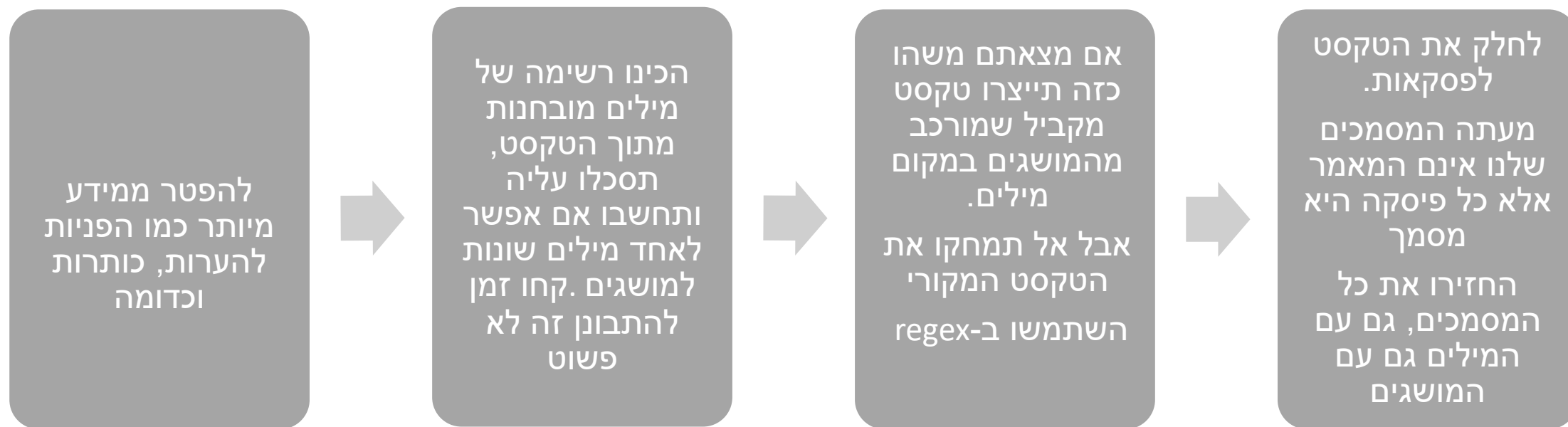
המודול מחזיר צמד של  
טקסט והסיווג שלו

# עיבוד מוקדם של טקסט

input: raw text

output: list of documents

צוות: עיבוד מקדים

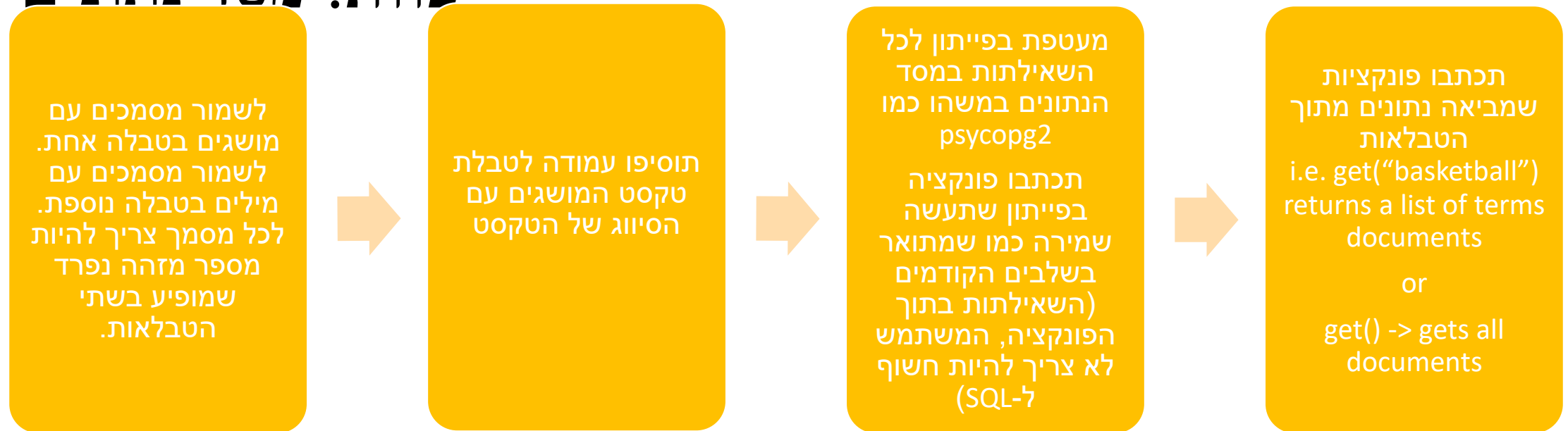


# טקסט למסד נתונים

input: list of documents (terms & words) + labels

output: get method in python that selects documents by label

## צוות: מסד נתונים



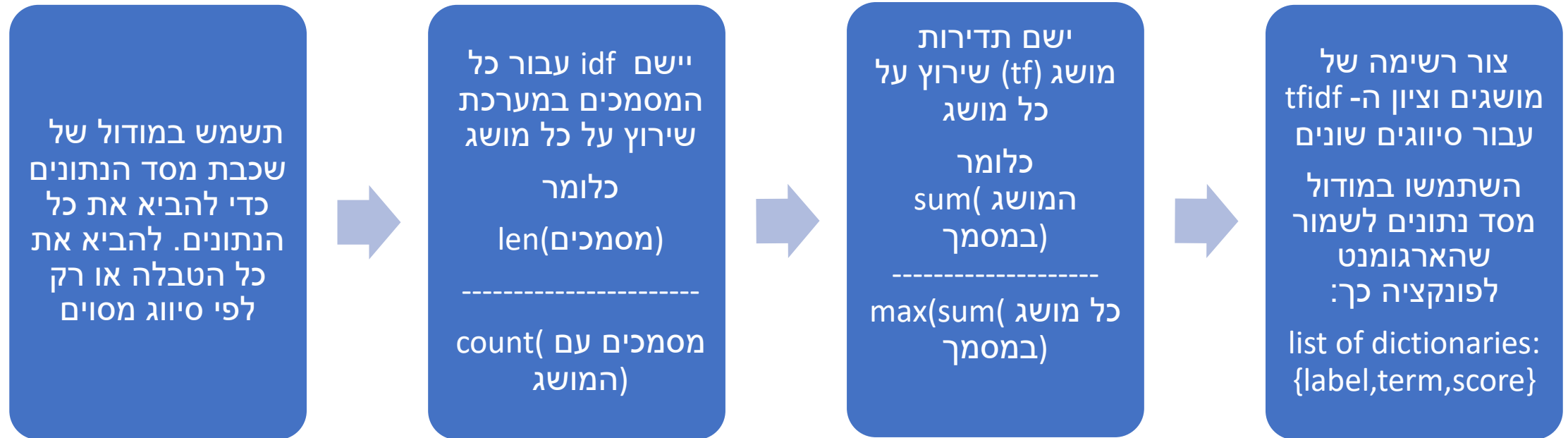
בניית המידע המזוין: מסמך מושגים הוא מסמך שיש בו הרבה טקסט ולא רק רשימת מושגים, רק הוא מחליף מילים במושגים מסמך מילים זה אותו דבר אבל עם המילים המקוריות

# tfidf

input: all the documents

output: list of dict {terms, labels & scores}

צוות: גרעין המערכת



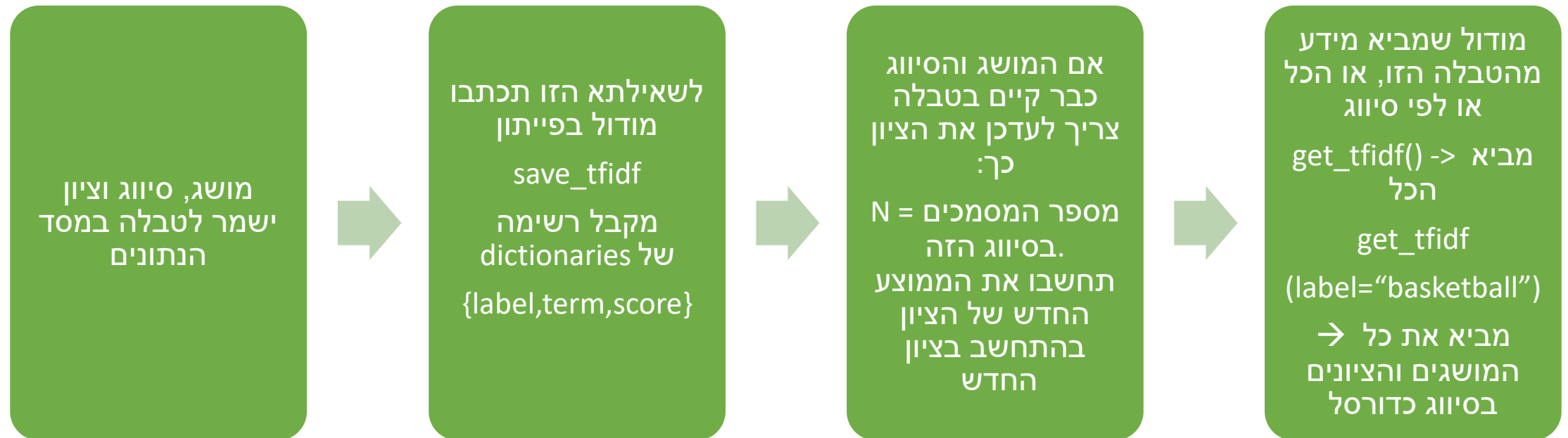
אתם לא צריכים לדאוג למודולים של שכבת מסד הנתונים ההנחה היא שאתם מקבלים אותם מצוות אחר

# טבלאות מבד נתונים ל tfidf

input: list of {term, label, tfidf\_score}

output: save\_tfidf([{t,l,score}]), get\_tfidf(\*\*label)

צוות: מסד נתונים



the list of dict is none of your concern they will be implemented by core team

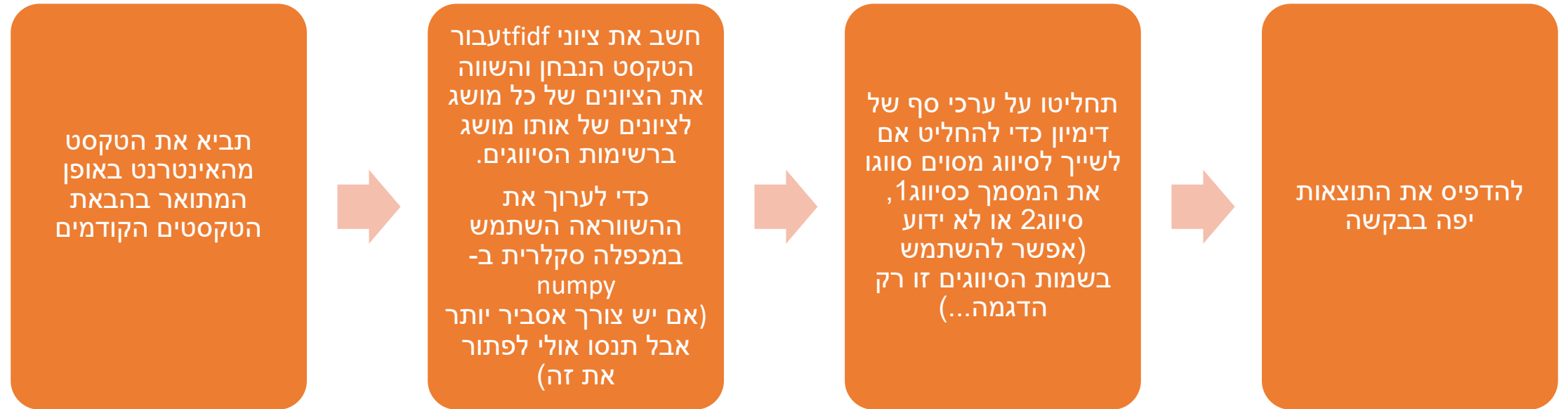


# סיווג של טקסט חדש

input: url address

output: assumed label

צוות: סריקת מידע וגרעין המערכת



# אינטגרציה

- כל חלק שתיארנו כאן ישמר כחבילה (כדאי אולי לשמור את כל שכבת המידע כחבילה אחת עם כמה מודולים)
- חברי צוות האינטגרציה הם ראשי הצוותים של שאר הצוותים ובראשם ראש צוות אינטגרציה (שגם הוא יכול להיות חלק מצוות אחר). תעשו ישיבה מקדימה לפני שמתחילים לפתח כדי להבין טוב מה כל צוות צריך להוציא בסוף ולראות שכולם מבינים אותו דבר, ואחרי זה תתכנסו עוד פעם אחת באמצע העבודה לדווח זה לזה מה קורה וכשמסיימים לבנות את המרכיבים תתחילו לתפור.
- הרעיון הוא להפטר מכל הזיופים במרכיבים השונים וליצור חיבורים אמיתיים לזרימת הנתונים

בהצלחה