

# Acronyms:

## Dictionary Construction & Disambiguation

Kayla Jacobs (Technion), Alon Itai (Technion), Shuly Wintner (University of Haifa). [WG1]

### Abstract

- ❖ Automatically build acronym dictionary
  - Apply to Hebrew
  - Rank multiple expansions by context match
  - Include local acronyms (unaccompanied by expansions)
- ❖ Improve acronym disambiguation
- ❖ Acronym expansions are usually MWEs



"Oh, it's an acronym for 'It Doesn't Stand For Anything.'"

### Why We Care

- ❖ Most acronym expansions are multi-word expressions (MWEs).
- ❖ Acronyms affect NLP applications like search and machine translation.
- ❖ Hand-crafted dictionaries incomplete and require constant updating.

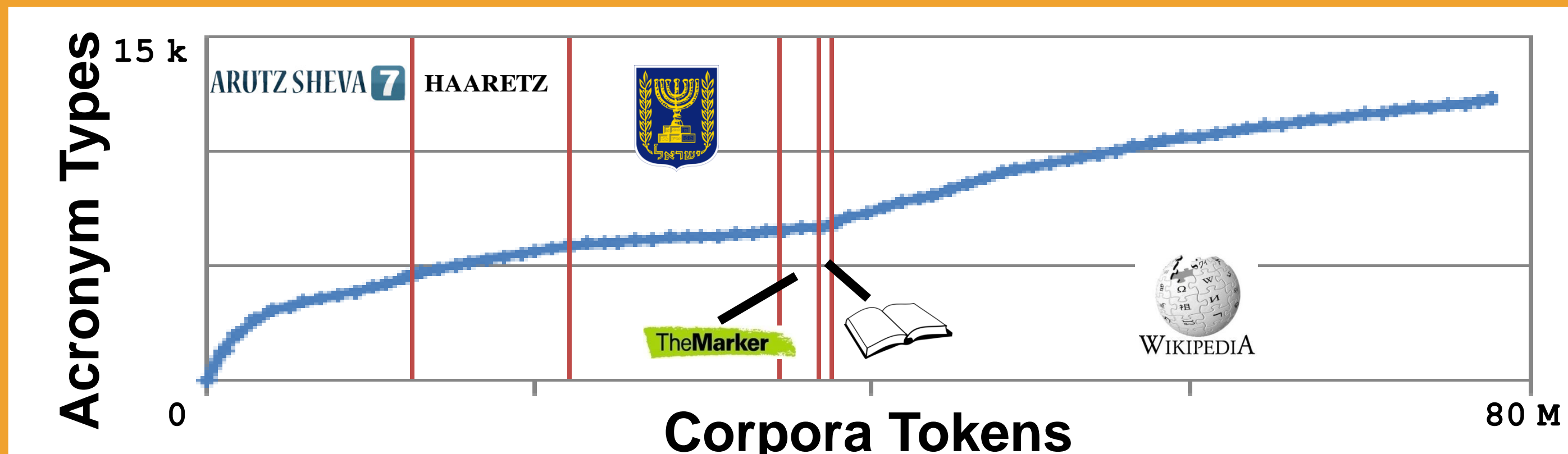
### Previous Work

- ❖ Prior acronym dictionary-building techniques rely on *local acronyms* (acronyms adjacent to their expansions, often in parentheses).
  - "The Central Intelligence Agency (CIA) released its budget."
  - "She works at the Culinary Institute of America (CIA)."
  - "Alumni of the Cleveland Institute of Art support the CIA."
- ❖ Only computational work on Hebrew acronyms: HaCohen-Kerner [04,08,10,13]
  - Disambiguation of Hebrew/Aramaic acronyms in Jewish law domain.
  - Assumes a pre-existing, hand-crafted acronym dictionary.

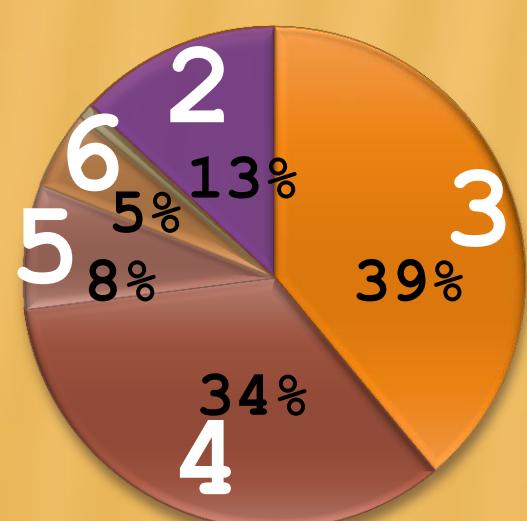
### Hebrew Acronyms

- ❖ In Hebrew corpus, acronyms 1% of word tokens and 3% of types.
- ❖ More common in news and encyclopedia genres than in literature.
- ❖ 2000+ years of frequent usage in Hebrew; ~100 years in English.
- ❖ Challenges from Hebrew's complex morphology and orthography.

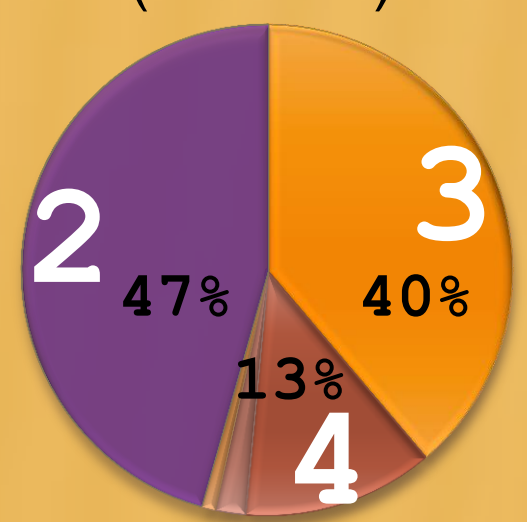
A never-ending story for unique acronyms:  
new acronyms continue to be found as more text is read



#### Acronym Lengths (# letters)



#### Expansion Lengths (# words)



| Letters | %   | Formation Rule  | Example  |
|---------|-----|-----------------|--|
| 2       | 98% | ■□□ ■□□         | ש"ק = שקל חדש (shekel new, "New Israeli Shekel")     |
| 3       | 48% | ■□□ ■□□ ■□□     | כ"א = אבל (but if thus, "unless")                    |
|         | 18% | ■□□ ■□□         | ב"ח = בית חולים (house-of sick-people, "hospital")   |
|         | 18% | ■□□ ■□□         | ת"ת = תתן (take and+give, "negotiation")             |
| 4       | 21% | ■□□ ■□□ ■□□ ■□□ | א"ע"כ = אף על פי כן (yet on as thus, "nevertheless") |
|         | 18% | ■□□ ■□□         | דוא"ל = דואר אלקטרוני (mail electronic, "e-mail")    |
|         | 13% | ■□□ ■□□         | ש"ש = שבת (exits-of Sabbath, "Saturday night")       |

### Building a Dictionary

#### 1 Identify acronyms

- ❖ Easy in Hebrew: unambiguous orthographic marking (internal " mark).
  - יו"ר = יושב ראש (sitter head, "chairperson")
- ❖ Difficult in English: capitalization and punctuation vary widely:
  - M.S. / MS / M.Sc. / MSc / MSC = Master of Science
  - au = atomic unit

#### 2 Identify potential expansions

- ❖ Collect corpus  $n$ -grams ( $2 \leq n \leq 5$ ). Public relations is easy.
- ❖ Discard  $n$ -grams that are infrequent or end with a preposition or quantifier.

| $n$ | $n$ -grams  | Freq.   |
|-----|---|---|
| 2   | public relations  | 1092  |
|     | relations <span style="color: red;">is</span> <span style="color: red;">X</span>        | 152   |
|     | is easy   | 5224  |
| 3   | public relations <span style="color: red;">is</span> <span style="color: red;">X</span> | 102   |
|     | relations is easy   | 23  |
| 4   | public relations is easy  | <span style="color: red;">1</span> <span style="color: red;">X</span> |

#### 3 Pair acronyms and expansions

- ❖ For each  $n$ -gram, generate all possible frequent acronyms via common formation rules.
- ❖ Tag with contextual info from LDA topic model.

public relations

| Rule    | Acronym | Freq.   |
|---------|---------|---|
| ■□□ ■□□ | PR      | 5293  |
| ■□□ ■□□ | PRE     | <span style="color: red;">2</span> <span style="color: red;">X</span> |
| ■□□ ■□□ | PURE    | 53  |

#### 4 Train classifier

- ❖ Train SVM to recognize matches:
  - ➕ Pairs from gold dictionary
  - ➖ Gold dictionary acronym paired with non-gold  $n$ -gram
- ❖ Linguistically-motivated classification features:
  - $n$ -gram PMI
  - acronym and  $n$ -gram document frequencies
  - formation rule acronym and  $n$ -gram lengths
  - LDA topic similarity score

| Match-Recognition Approach                                   | Precision   | Recall      | F-score     |
|--|-------------|-------------|-------------|
| <b>Baseline</b>  |             |             |             |
| Guess acronym's most-frequent $n$ -gram is correct expansion | 55 %        | 3 %         | 5 %         |
| <b>Our classifier</b>  | <b>82 %</b> | <b>81 %</b> | <b>82 %</b> |

### Acronym Disambiguation

- ❖ Extrinsically evaluated dictionary on acronym disambiguation task.
- ❖ Given 200 acronyms and their contexts, how many of the *correct* expansions are in the top  $r$  dictionary results for the acronyms?

| Dictionary  | $r = 1$     | $r = 2$     | $r = 3$     | $r = \infty$ |
|---|-------------|-------------|-------------|--------------|
| <b>Baseline #1:</b><br>Dictionary of local parenthetical acronyms |             |             |             | 52 %         |
| <b>Baseline #2:</b><br>Gold dictionary                            | 66 %        | 78 %        | 81 %        | 83 %         |
| <b>Our dictionary</b>   | <b>72 %</b> | <b>79 %</b> | <b>77 %</b> | <b>85 %</b>  |

| Error Rate Reduction           | $r = 1$ | $r = 2$ | $r = 3$ | $r = \infty$ |
|--------------------------------|---------|---------|---------|--------------|
| Our Dictionary vs. Baseline #1 |         |         |         | 69 %         |
| Our Dictionary vs. Baseline #2 | 18 %    | 8 %     | 14 %    | 14 %         |

### Future Work

- ❖ Exploit for identifying multi-word expressions (MWEs).
- ❖ Apply to other languages
  - Hebrew advantages: Easy acronym identification, very widespread acronym use.
  - Hebrew disadvantages: Complex morphology/orthography, poor NLP resources.
- ❖ Additional applications: search and machine translation.