

---

# *Text Mining in Hebrew*

---

## *Impact of Morphology Analysis on Topic Analysis and on Search Quality*

Michael Elhadad, Meni Adler, Yoav  
Goldberg, Dudi Gabay and Yael Netzer

---

# Text in Hebrew

- Extract information from text in Hebrew
- Major immediate obstacle
  - Rich morphology
  - Very high number of distinct word forms
  - Very high ambiguity

# Morphological Analysis

## ■ בצלם

- בצָלָם (name of an association)
- בְּצִילָם (while taking a picture)
- בְּצֻלָם (their onion)
- בְּצֵלָם (under their shades)
- בְּצַלְמָם (in a photographer)
- בַּצֵּלָם (in the photographer)
- בְּצֻלָם (in an idol)
- בַּצֵּלָם (in the idol)

---

# How Critical is Morphological Analysis to Text Mining?

- How much does Hebrew morphology affect high-level text analysis tasks?
  - Named Entity Recognition
  - Information Extraction
  - Topic Analysis
  - Information Retrieval

---

# Topic Analysis and Search

## ■ Topic Analysis

- ❑ Unsupervised discovery of topics in text collection
- ❑ Useful to browse a large corpus by theme
- ❑ Difficult to evaluate

## ■ Faceted Search

- ❑ Useful combination of search and browsing
- ❑ Exploratory search (as opposed to fact finding)
- ❑ Enabled by topic analysis

# The Basic Idea

0 עֵשֶׂר קָדָה סָלַע מֵעַה כְּסָף שְׁנֵי יְרוּשָׁלַיִם פְּרִי שׁוֹה ח' ל' אָמַר מֵע"ש ח' מֵשׁ יָצָא הוֹסִיף נָתַן פֶּרֶט דִּינָר זָקָב חֲלָל

1 כּוֹקֵב עֲבֹדָה הַנְּאָה עֲבָד עֹבֵד אָסַר צוּקָה עָשָׂה אָבֵן לָהּ זֹו דָרַךְ בֵּעַל יִשְׁכָּאֵל עוֹלָם נִסְקַל בַּהּ נֶאֱסַר לָקָה בָּנָה

2 כְּתִבָּה בֵּעַל לָהּ אִישׁ נִכְסֵי נָתַן הוֹצִיא מִזֶּזֶן דִּין אוֹתָהּ קָצָה מֵת ל' א' נָטַל יֵשׁ בֵּית עֶקֶר בֵּין אַחַר קָכַר

3 יִהְיֶה עֵשֶׂר קָנָה יָדַע אַחַת הוֹסִיף אֶרְבָּעִים אֶמְצַע כ' ל' שְׁנֵי חֲמֵשׁ

4 שָׁם מֵעַה עוֹף מֵת זֹו נִדְבָּה קָדַשׁ אַחַר מֵעַה שְׁמִים אָמַר קָרַב

5 עַד אָמַר הַעִיד נֶאֱמַן פְּלוּנִי דִין בָּא דָבָר פְּנִים יָדַע קָאָה מֵת פָּה זֹו אַחַר עֲדוּת שְׁנָה קִיָּה אָסַד בָּךְ

אשה איש האשה אשתו אשת נשים האיש אנשי לאשה הנשים  
אנשים לאשתו לאיש והאשה לאשתי ואשה ואיש האנשים  
ואשת באשה ואשתו ואנשי שהאשה אשתי לאנשי ונשים  
באיש באשתו מאיש נשותיהן והאיש כאנשי בנשים לאשת  
מאנשי שאשתו לאנשים נשיו כאיש מאשה והנשים שהנשים  
מאשתו באשת לנשים שהאיש ואנשים

One word איש – about 50 distinct forms in the corpus

---

# Outline

- Objectives
  - Topic Analysis in Hebrew
  - Improved Search
- Topic Analysis with LDA
- Obtaining Precise Morphology in Hebrew
- Combining LDA and Morphological Analysis
- Using Topic Models for Search
- Evaluating Topic Models
- Next Steps

---

# Objectives

- Input:
  - Domain specific text corpus in Hebrew
- Output:
  - Topic model:
    - Discover “topics” discussed in the corpus
    - Recognize topics in unseen text
  - Index text collection by topic
- Task:
  - Search and browse text collection using topics



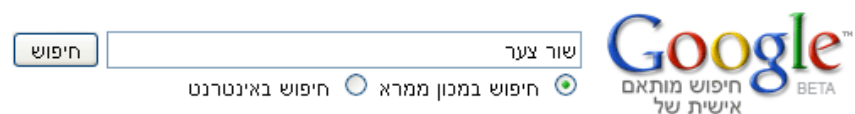
---

# Example: Rambam's Mishne Torah

- Corpus of Mishne Torah
  - Exhaustive code of Halakha
  - Written by Maimonides 1170-1180
  - 14 books, 85 sections, 1000 chapters, 15K articles, 350K words.
- Creative compilation of laws from multiple sources:
  - Torah, Talmud (Bavli and Yerushalmi), Tosefta, halakhic midrashim (sifra and sifre), Geonim.
  - Synthetic hierarchical organization

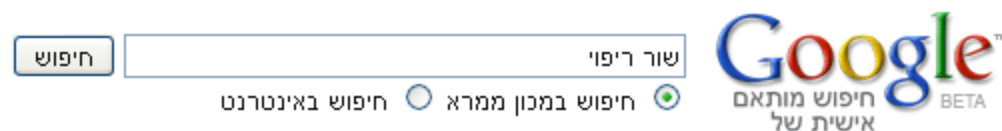
# Problems with Existing Search

## ■ Morphology



[משנה תורה - ספר נזקים - הלכות חובל ומזיק פרק ד](#)

[יב] מי שחציו עבד וחציו בן חורין, שביישו אדם, או ציערו, או שנגחחו **שור**, ... אבל **צער** ובושת וריפוי הרי הוא שלה; וכן נזק שאינו פוחת אותה מכספה, הרי הוא שלה. ...  
[www.mechon-mamre.org/i/b404.htm](http://www.mechon-mamre.org/i/b404.htm)



A single “ו” and the word is not found...

החיפוש שלך - **שור ריפוי** – לא תאם אף מסמך

# Problems: Explore complex topics

- “שור” refers to many complex halakhic topics:
  - Damages (שור נוגח)
  - Kosher meat (שחיטה)
  - Sacrifices (קרבנות)
  - Shabat (שבת)
  - Calendar (מזל שור)
- Queries must be disambiguated
  - שור+שבת?

# Exploratory/Faceted Search

- How to deal with ambiguous query terms?
  - Propose refinements according to contexts
  - “*Do you mean: damages, meat, shabat...*”
  - Propose facets for query refinement
- Where do the topics (facets) come from?
- How do we disambiguate the query terms?
- Given a disambiguated topic, how do we refine the query?

---

# Outline

- Objectives
- Topic Analysis with LDA
- Obtaining Precise Morphology in Hebrew
- Combining LDA and Morphological Analysis
- Using Topic Models for Search
- Evaluating Topic Models
- Next Steps

---

# Discovering Topic Models: LDA

- Latent Dirichlet Allocation
  - Blei and Jordan 1993
- Discover (unsupervised) topic structures in a document collection
- Topics are modeled as distributions of words
- Probabilistic generative model of text

# שור Topics for

- 65 שלם הויק בעל נזק שור בהמה קשות חגב נזוק חץ פטר בור מועד נפל הקה אדם את לו חבירו דין 0.0217290799815
- 83 הביא מנחה מן כלי אומר עלה נתן שמו לוג עשרון שר לבונה מין ק'טז יצא שם חץ נדר מזבח גב 0.00679501698754
- 72 בהמה חמור בה עלה פכה צריך אנה אדם מחשבה עגלה כשר חשב חי קשר הכניס שור מלאכה אין קלב בעל 0.00636205899364
- 73 קאה מול ירח עשר יהיה קאנה מעלה שני מן ר'אש צפון ידע ראשון קשה גרע דרום ליל עולם חלק ממנו 0.00604686318972
- 54 בשר מן חי קלב בהמה עוף דם אבר וית דג אכל טמא טהור ביצה בה אפר מין עצם תוכה שחט 0.002
- 35 מום בכור בהמה בעל תמוכה זו אומר כ'הן קדש בה בו אנה תמים עשה נולד קבוע בית הקדוש מזבח ח'ל 0.00176574455562
- 92 גר אש הדליק אור גב תנור עץ כלי קדרה בשל חגכה- גחלת תבשיל הניח הוסיף תוך עלה צריך חשכה חם 0.00129449838188
- 66 פתוח סתום כי אל דבר משה אומר 'שתים בן יי איש את שלש עשה ויהי כולן יום הוא שש 0.00125680770842
- 53 שחט פסח שחיטה אכל שני את סכין שם עשה ראשון שר אחר תבוכה עזקה דחה דם ישחוט עליו ח'ל ארבעה 0.00125588697017
- 82 נפל כ'ל מאה ח'ל סאה עצם אפר תרומה תוך התערב יין מין אפור כהם אי נתן קלא עלה דבר שש 0.00098167539267
- 81 קדש מעל הקדוש נהנה דם הקדש בדק בית פדה פרט נתן מזבח יצא בהמה ח'מש יוכל דבר אותה נפל הוסיף 0.000669792364367
- 96 שלם גבב תשלום קרן ארבעה חמש דם דין גבב כפל נתן בעל חגב את ח'מש פטר מן מכר קנס ממון 0.000576701268743
- 49 מים יד קגל חצץ גוף ר'אש עין אדם נטל ידיים עלה : נתן על בשר עבל צריך מבילה גב שער 0.000497636227917
- 84 ( ) א'כל חוץ אכל : נכנס דו בג א אל יד י נתן הז ח עני 0.000467508181393
- 64 קדש איש אומר לי לה זו קדש קדושין את מקדש נתן לו דבר אנה פרט ספק צריך בת נמצא מנה 0.000322684737012

# Topics for a Document

## :topics

65 [שלם הויק בעל נזק שור בהטת לשות תנב נזק חץ פטר בור מוגד נפל תה אדם את לו חבירו דין 0.0822931114193](#)

56 [תנב פטר לשות הוציא גד הניח מים רב קמד חוך פנים שני חפץ זרק דרך חוץ הר פבר חבירו קקר 0.0153750259713](#)

40 [אטה תה עפח ארבע עשר ארבעה הוא יש כ' תל יתר חוך בין עשה פחות מקום מן זו שלש מחצה בו 0.00681000097286](#)

50 [+ ש נ' " prefix=כ עשה עבר שתה לה שם '+ דבר מזה דבר ה' ' רע אשר נתן פרט 0.00546919976972](#)

## content

החופר בור ברשות הרבים ונפל לתוכו שור חמור ומת אפילו היה הבור מלא צמר וכיוצא בהן בעל הבור חייב לשלם נזק שנאמר בעל הבור ישלם ואחד שור וחמור שאר מיני בהמה ויהי ועוף נאמר שור חמור בהווה אחד החופר בור ברשות הרבים החופר ברשותו ופתח לרשות הרבים פתח לרשות חבירו שחפר ופתח ברשותו רשותו הפקיר חייב בנזקי הפקיר רשותו הפקיר שברשותו הקדישו פטור שנאמר בעל הבור ישלם מי שיש לו בעלים וזה הפקר ברשות חפר מפני שחפר ברשותו אחד החופר הבור מאליו בהמה חיה הואיל והוא חייב עשה חייב בנזקיו ואחד החופר הלוקח שנתן לו במתנה שנאמר בעל הבור ישלם מי שיש לו בעלים מכל מקום אחד החופר מקום שהיה מכוסה שנאמר + שמות כ"א prefix="ל" ג + וכי יפתח איש בור כי איש בור כראוי אע"פ מתוכו ונפל לתוכו שור ומת פטור שנאמר הא פטור בדבר שיכול לעמוד בפני שוורים ואינו יכול לעמוד בפני גמלים והלכו עליו גמלים והלכו עליו שוורים ונפלו בו הגמלים מצויין באותו מקום פטור מפני שזה אונס יבואו שם גמלים אפילו חייב מתוכו ונפלו בו שוורים אע"פ מצויין שם תמיד והרי הוא פושע הואיל ומחמת נפלו בו פטור כיוצא בזה המוצא בור וחזר בעל הבור חייב וזה האחרון פטור תמנו בעפר וחזר והוציא את העפר האחרון חייב שכיון בעפר מעשה הראשון בור שני שותפין ועבר עליו הראשון והשני הראשון חייב לשני לשני מוטו נפטר הראשון ונתחייב השני הראשון ובא השני ומצאו מגולה השני חייב ועד אימתי יהיה השני לבדו חייב שידע הראשון מגולה וכדי פועלים וכל שימות בו בתוך זמן השני לבדו חייב בו וכל שימות בו אחר זמן כזה שניהן חייבין לשלם שהרי שניהן בו המוסר לשומר השומר חייב בנזקיו לחרש שוטה וקטן אע"פ שהיה מכוסה הבעלים חייבין עשוי ואלו בהן דעת חבירו ובא בעל הדלי ונטל בעל הבור חייב אחד החופר בור מערה חריץ ולמה נאמר בור שיהיה בו כדי להמית וכמה כדי להמית עומק עשרה טפחים היה פחות מעשרה ונפל לתוכו שור שאר בהמה חיה ועוף ומת פטור חייב בעל נזק שלם היה עומק הבור תשעה ומהן טפח אחד מים חייב מים חשוב שני טפחים ביבשה היה שמונה ומהן שני טפחים מים שהיה שבעה ומהן שלשה טפחים מים ונפל לתוכו שור וכיוצא בו ומת מחייבין אותו לשלם תפש הניזק מוציאין מידו שהדברים האלו ספק יש בהן החופר בור עמוק עשרה טפחים ובא אחר לעשרים ובא שלישי לשלשים כולן חייבין חפר הראשון פחות מעשרה אפילו טפח ובא האחרון לעשרה בין שחפר בו טפח שהגביה בנין שפתו טפח האחרון חייב סתם טפח שהוסיף טפח שבנה ספק כבר מעשה הראשון עדיין חפר הראשון בור עמוק ובא האחרון ונפל לתוכו שור ומת מחמת מת האחרון פטור שהרי מיעט מחמת מת האחרון חייב שהרי הוא הקריב בור נפל השור מאותו הצד האחרון האחרון חייב שהרי הקריב בור אע"פ שמת מן ההבל מן הצד שחפר הראשון נפל הראשון חייב שזה האחרון מיעט בור עליו התורה אפילו מתה הבהמה ואין צריך לומר מתה לפיכך היה עומק הבור לו הבל בו הבהמה ומתה פטור היה יתר רחבו יש לו הבל מתה בו הבהמה חייב אע"פ שלא עשה תל גבוה ברשות הרבים בו הבהמה ומתה היה גבוה עשרה טפחים חייב לשלם היה פחות מעשרה פטור מיתת הבהמה הם בלבד חייב לשלם נזק שלם ואפילו גבוה שהוא בכל שהוא דבר מצוי וידוע ואין המיתה בכל שהוא מצויה והרי הוא כמו אונס אינו חייב מיתת הבהמה



---

# The LDA Model

- Observation: documents are composed of words.
- Intuition: documents exhibit multiple topics
- Generative probabilistic model:
  - Each document is a mixture of topics
  - Each word is drawn from the topics active in the document

# Structure of the LDA Model

## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

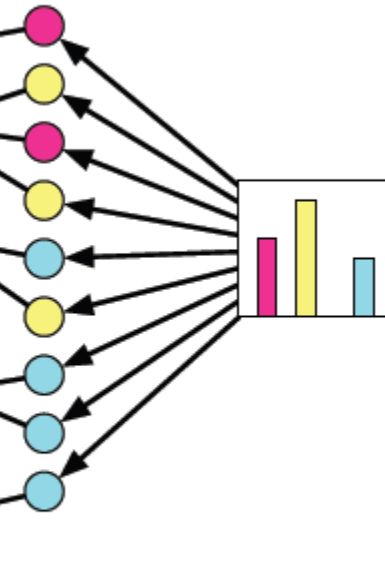
\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

“are not all that far apart,” especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. “It may be a way of organizing any newly sequenced genome,” explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.



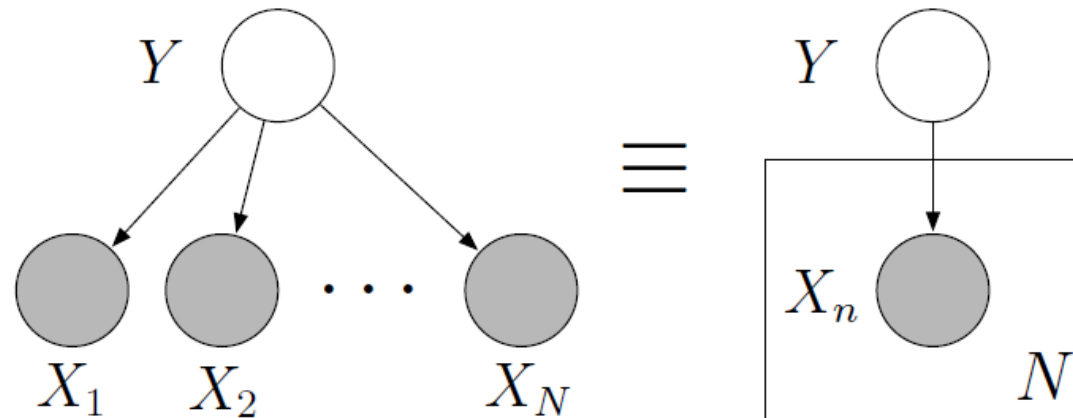
From (Blei 2008)

---

# Learning an LDA Model from Observations

- Observation: documents and words
- Objective: infer an underlying topic structure
  - What are the topics?
  - How are the documents divided according to those topics?

# Graphical Models

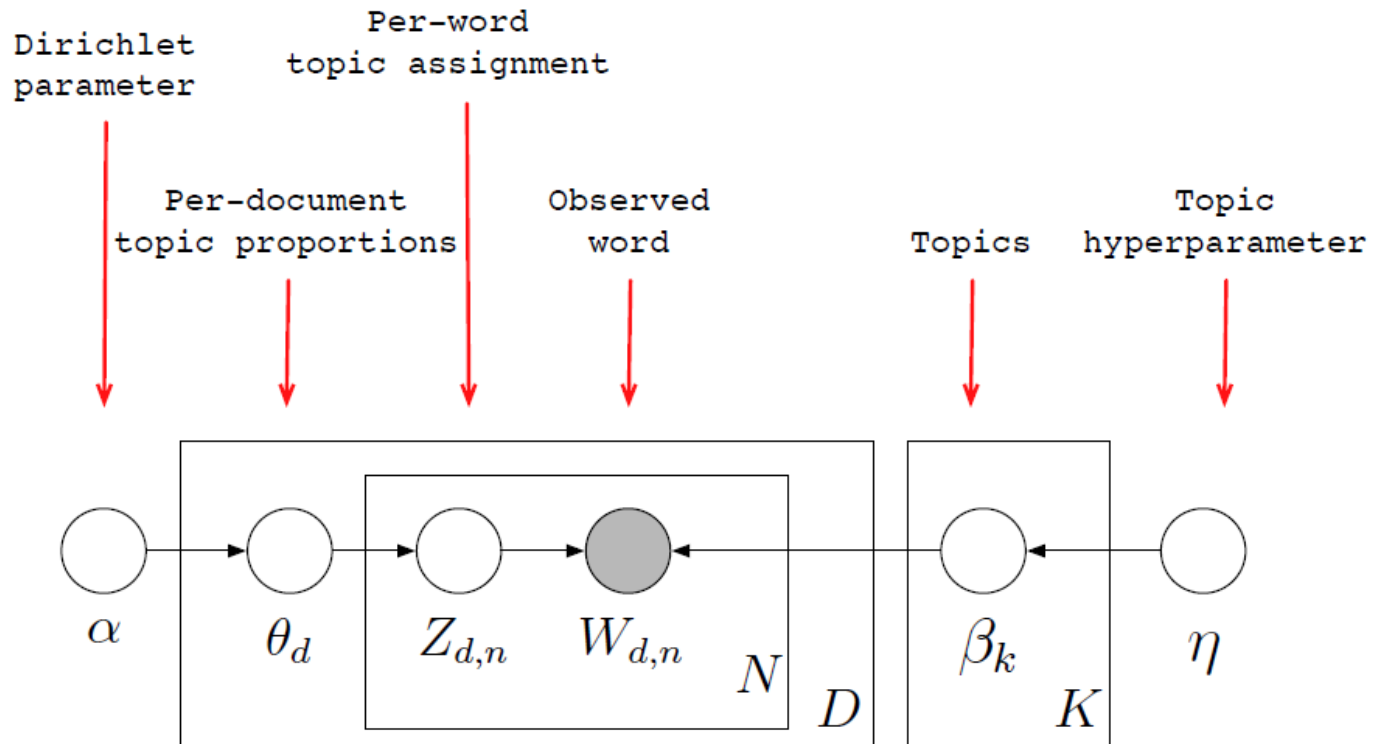


- Structure of the graph defines the pattern of conditional dependence between the ensemble of random variables
- E.g., this graph corresponds to

$$p(y, x_1, \dots, x_N) = p(y) \prod_{n=1}^N p(x_n | y)$$

(Blei 2008)

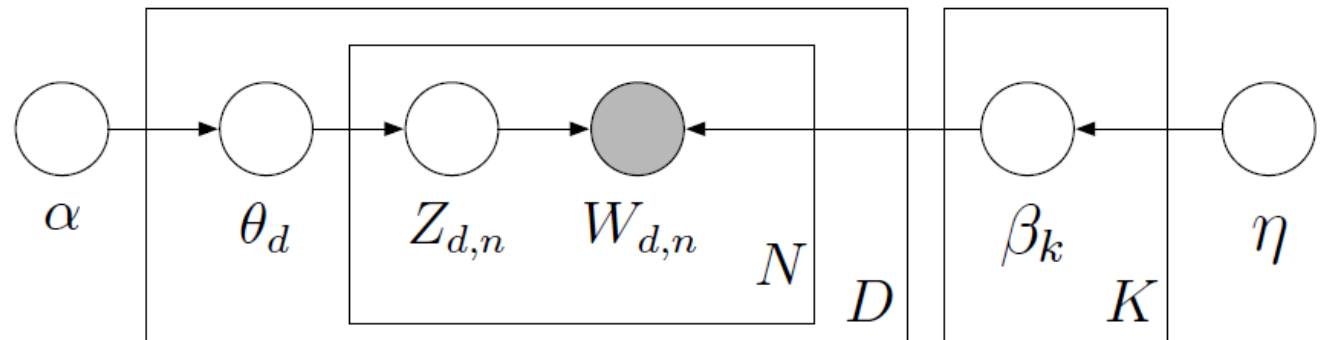
# LDA Graphical Model



(Blei 2008)

Each piece of the structure is a random variable.

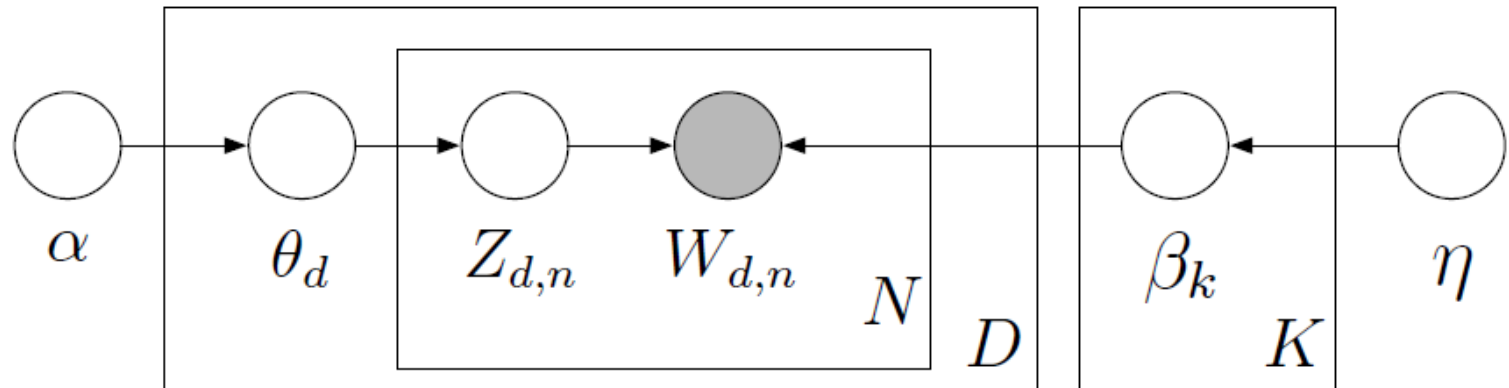
# LDA Generative Process



- 1 Draw each topic  $\beta_i \sim \text{Dir}(\eta)$ , for  $i \in \{1, \dots, K\}$ .
- 2 For each document:
  - 1 Draw topic proportions  $\theta_d \sim \text{Dir}(\alpha)$ .
  - 2 For each word:
    - 1 Draw  $Z_{d,n} \sim \text{Mult}(\theta_d)$ .
    - 2 Draw  $W_{d,n} \sim \text{Mult}(\beta_{Z_{d,n}})$ .

(Blei 2008)

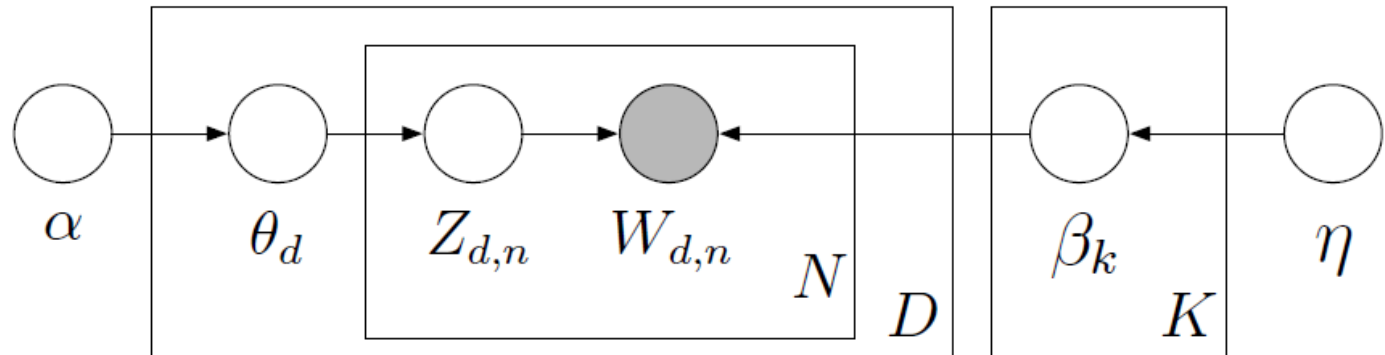
# LDA Estimation



- From a collection of documents, infer
  - Per-word topic assignment  $z_{d,n}$
  - Per-document topic proportions  $\theta_d$
  - Per-corpus topic distributions  $\beta_k$

(Blei 2008)

# LDA Approximation



- Computing the posterior is intractable:

$$\frac{p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta_{1:K})}{\int_{\theta} p(\theta | \alpha) \prod_{n=1}^N \sum_{z=1}^K p(z_n | \theta) p(w_n | z_n, \beta_{1:K})}$$

- Several approximation techniques have been developed.

(Blei 2008)



---

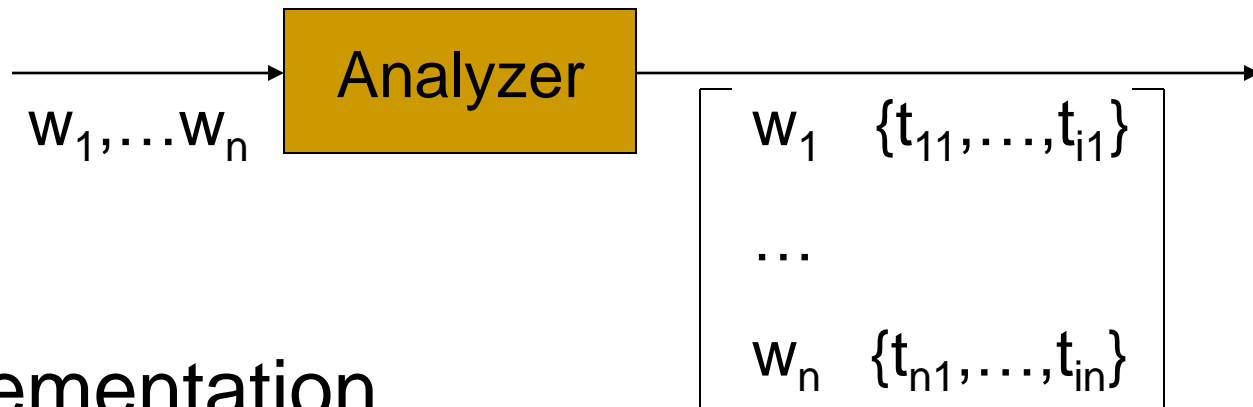
# Outline

- Objectives
- Topic Analysis with LDA
- Obtaining Precise Morphology in Hebrew
- Combining LDA and Morphological Analysis
- Using Topic Models for Search
- Evaluating Topic Models
- Next Steps

# Morphological Analysis

- בְּצִלָּם
  - בצלם proper-noun
- בְּצִלֵּם
  - בצלם verb, infinitive
- בְּצִלָּם
  - בצל-ם noun, singular, masculine
- בְּצִלָּם
  - בצל-ם noun, singular, masculine
- בְּצִלָּם בְּצִלָּם
  - בצל-ם noun, singular, masculine, absolute
  - בצל-ם noun, singular, masculine, construct
- בְּצִלָּם בְּצִלָּם
  - בצל-ם noun, definitive singular, masculine

# Morphological Analyzer



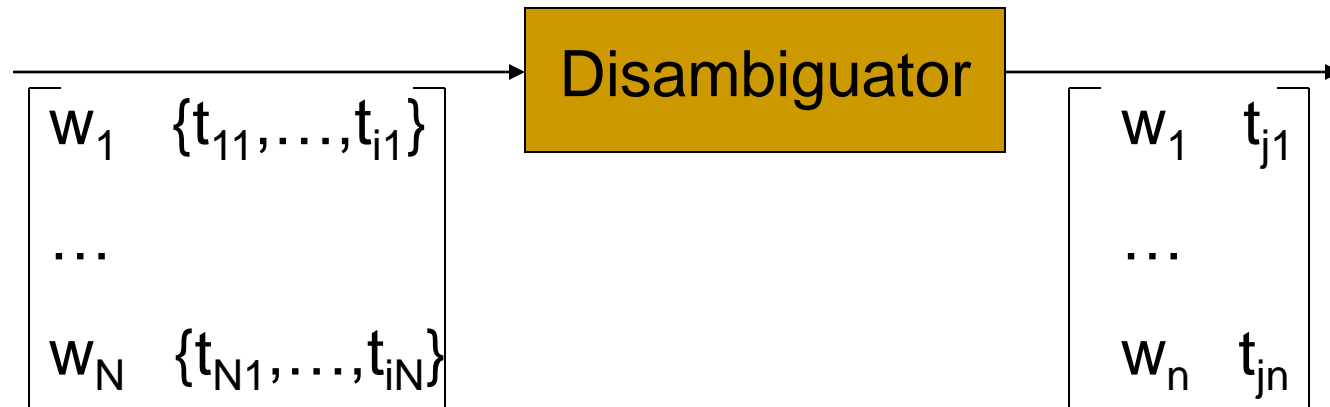
## ■ Implementation

- Corpus based
- Lexicon based
  - Analytic
  - Synthetic

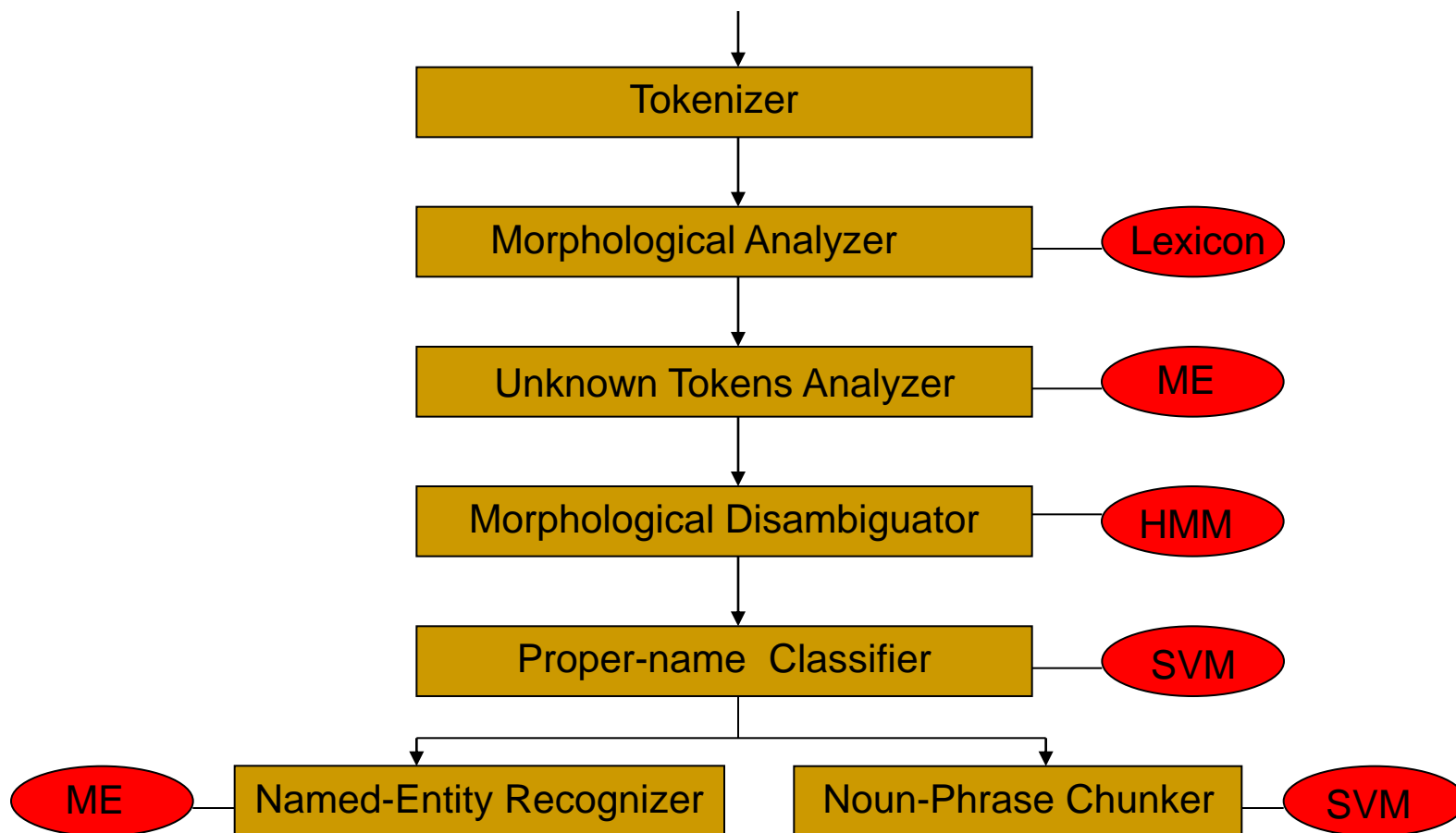
# Morphological Disambiguation

- ארגון בצלם הודיע
- בצלם את המשחק הבחנתי בעוזר המאמן
- בצלם נחטף בשווקים
- בצלם הנעים חסינו
- נתקלתי בצלם מקצועי
- פגשתי בצלם חתונות
- נתקלתי בצלם המקצועי

# Morphological Disambiguation



# Hebrew Text Analysis System



<http://www.cs.bgu.ac.il/~nlpproj/demo>

# Morphological Disambiguation - Methods

- Rule-based vs. Stochastic models
- Supervised vs. Unsupervised learning
- Exact vs. Approximate inference

# Hidden Markov Model

- $S$  – a set of states (= tags)
- $O$  – a set of output symbols (= tokens)
- $\mu$  – a probabilistic model
  - State transition probabilities  $A = \{a_{i,j}\}$
  - Symbol emission probabilities  $B = \{B_{i,k}\}$



# HMM- An Example

- $S = \{\text{start, noun, verb}\}$
- $O = \{\text{ירח, ילד}\}$
- $\mu = (A, B)$

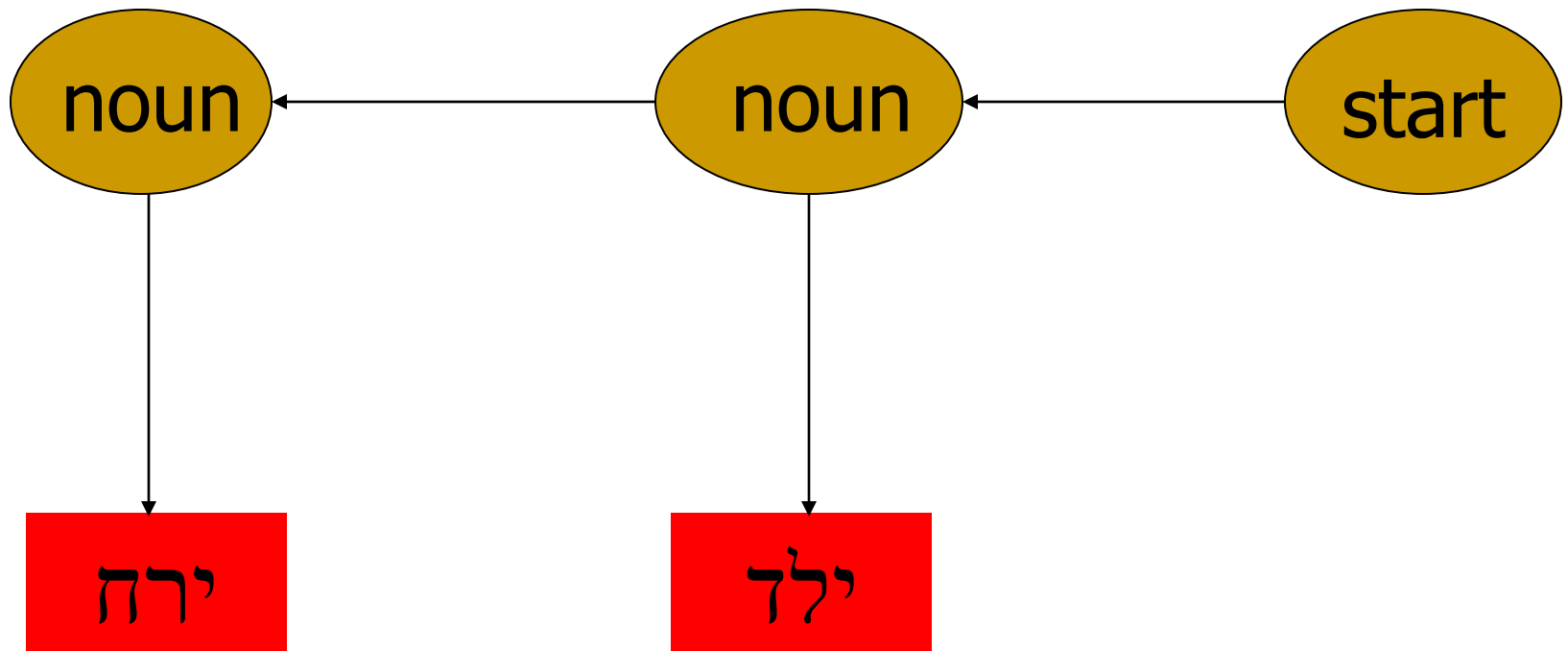
	noun	verb
start	0.8	0.2
noun	0.9	0.1
verb	0.9	0.1

A

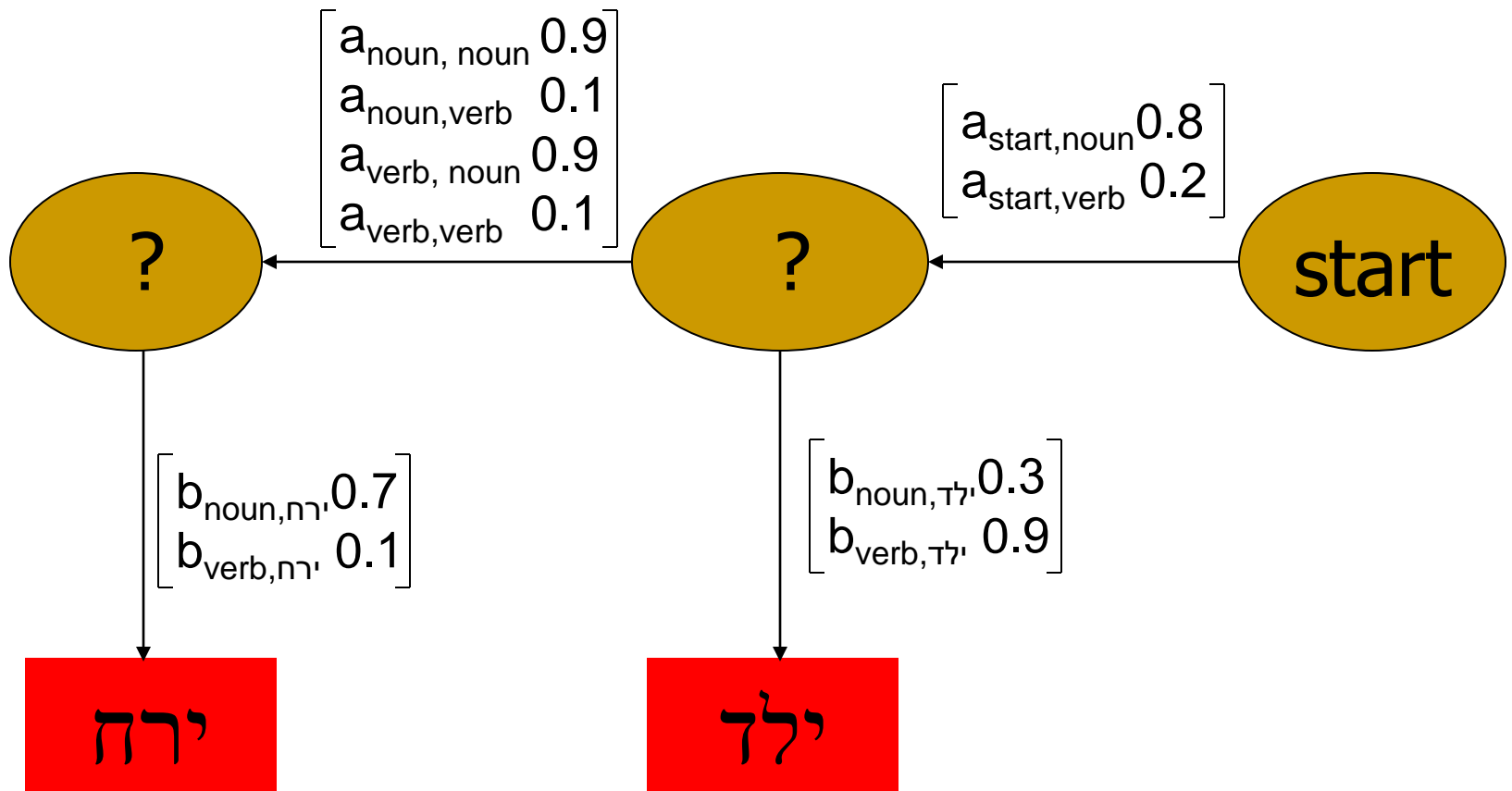
	ילד	ירח
noun	0.3	0.7
verb	0.9	0.1

B

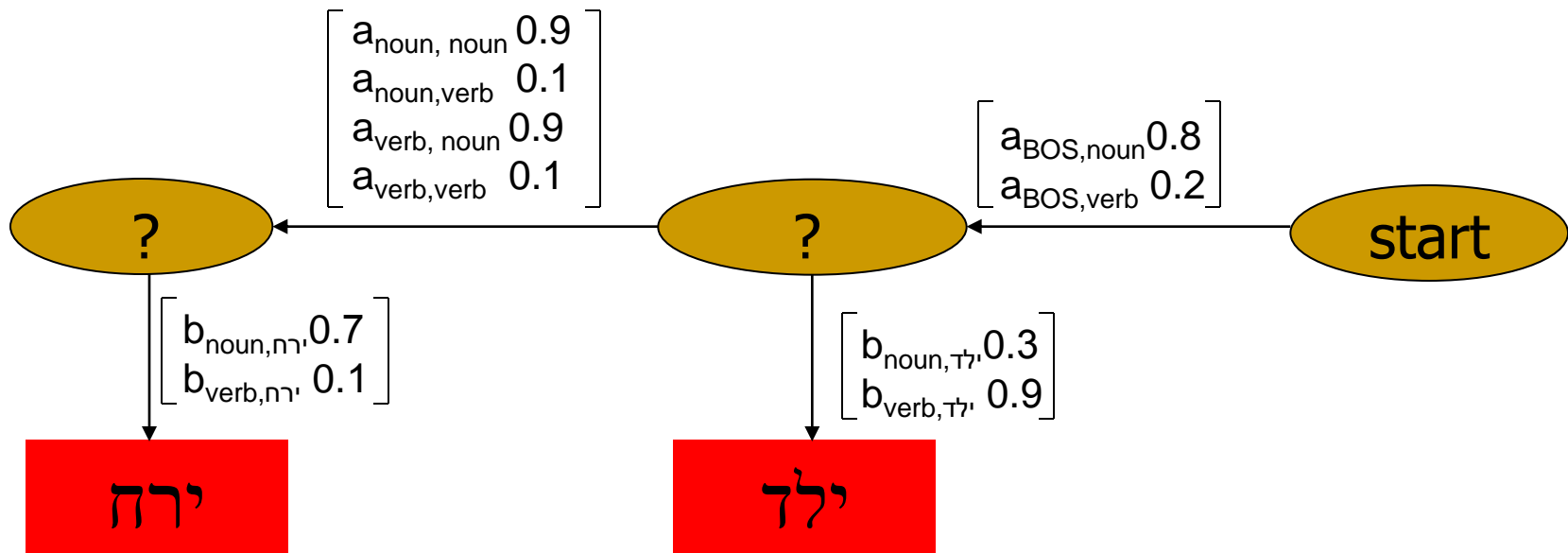
# Markov Process



# Decoding

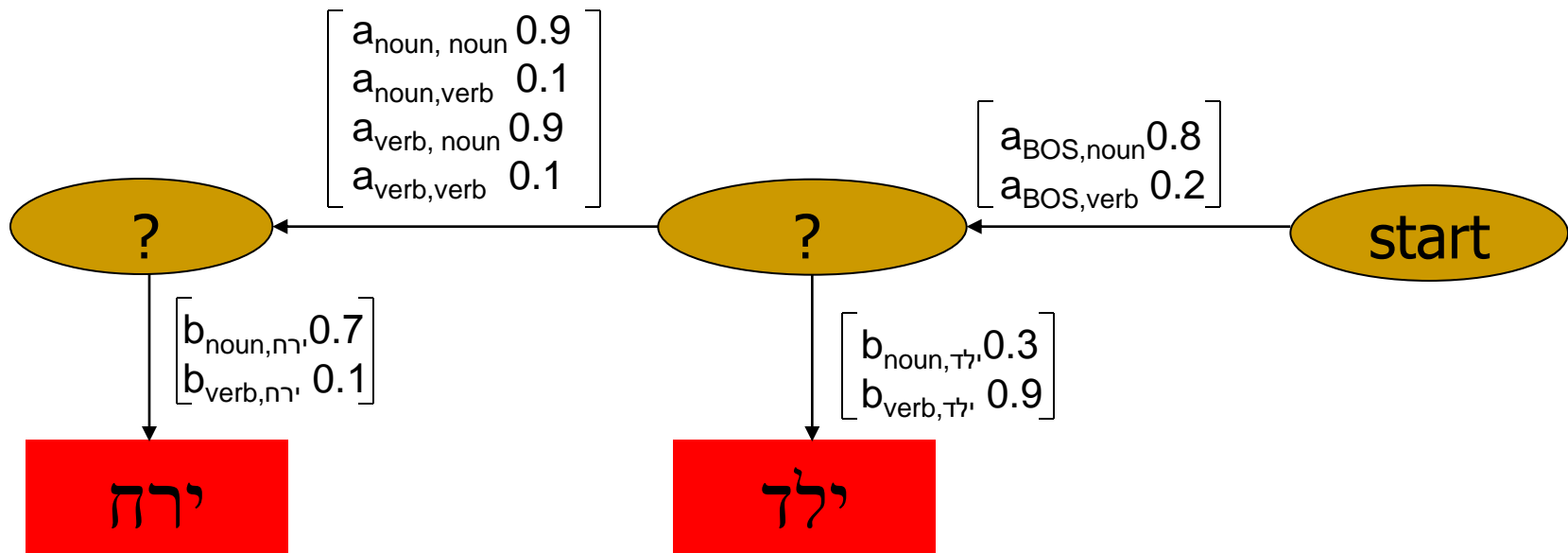


# Decoding



$$\begin{aligned} (\text{noun}, \text{noun}) &= a_{\text{start}, \text{noun}} b_{\text{noun}, \text{ילד}} a_{\text{noun}, \text{noun}} b_{\text{noun}, \text{ירה}} = 0.8 * 0.3 * 0.9 * 0.7 = 0.1512 \\ (\text{noun}, \text{verb}) &= a_{\text{start}, \text{noun}} b_{\text{noun}, \text{ילד}} a_{\text{noun}, \text{verb}} b_{\text{verb}, \text{ירה}} = 0.8 * 0.3 * 0.1 * 0.1 = 0.0024 \\ (\text{verb}, \text{noun}) &= a_{\text{start}, \text{verb}} b_{\text{verb}, \text{ילד}} a_{\text{verb}, \text{noun}} b_{\text{noun}, \text{ירה}} = 0.2 * 0.9 * 0.9 * 0.7 = 0.1134 \\ (\text{verb}, \text{verb}) &= a_{\text{start}, \text{verb}} b_{\text{verb}, \text{ילד}} a_{\text{verb}, \text{verb}} b_{\text{verb}, \text{ירה}} = 0.2 * 0.9 * 0.1 * 0.1 = 0.0018 \end{aligned}$$

# Decoding



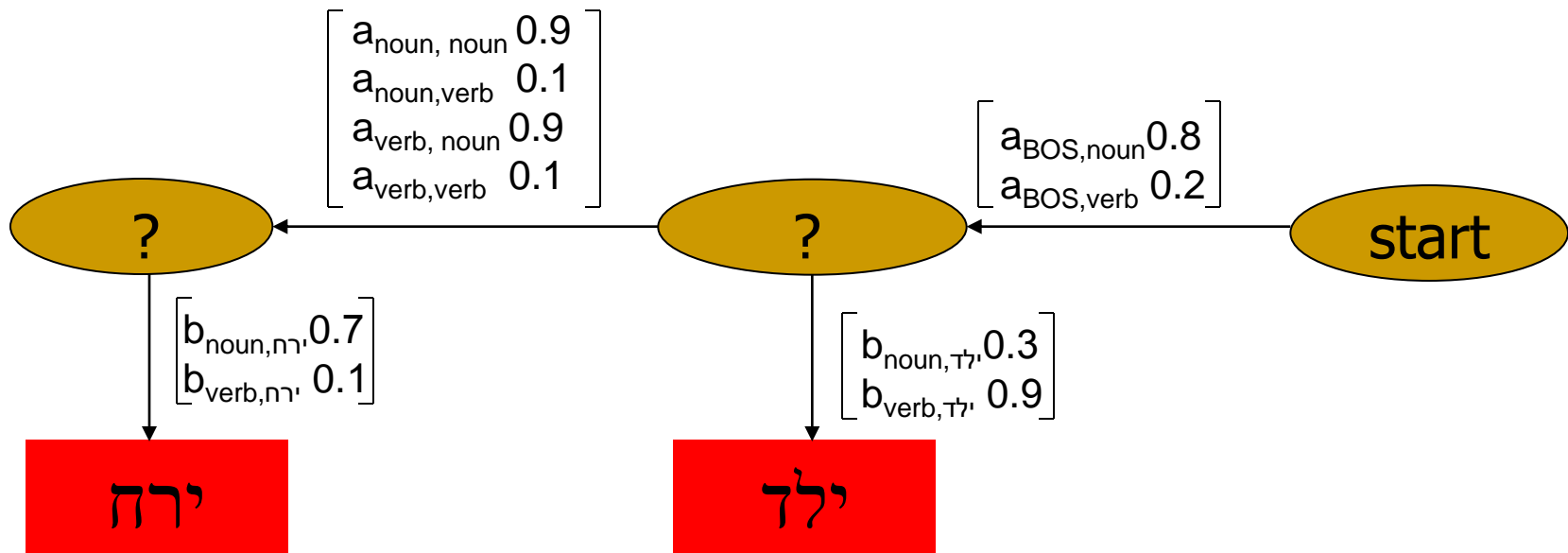
$$(noun, noun) = a_{\text{start}, \text{noun}} b_{\text{noun}, \text{ילד}} a_{\text{noun}, \text{noun}} b_{\text{noun}, \text{ירה}} = 0.8 * 0.3 * 0.9 * 0.7 = 0.1512$$

$$(noun, verb) = a_{\text{start}, \text{noun}} b_{\text{noun}, \text{ילד}} a_{\text{noun}, \text{verb}} b_{\text{verb}, \text{ירה}} = 0.8 * 0.3 * 0.1 * 0.1 = 0.0024$$

$$(verb, noun) = a_{\text{start}, \text{verb}} b_{\text{verb}, \text{ילד}} a_{\text{verb}, \text{noun}} b_{\text{noun}, \text{ירה}} = 0.2 * 0.9 * 0.9 * 0.7 = 0.1134$$

$$(verb, verb) = a_{\text{start}, \text{verb}} b_{\text{verb}, \text{ילד}} a_{\text{verb}, \text{verb}} b_{\text{verb}, \text{ירה}} = 0.2 * 0.9 * 0.1 * 0.1 = 0.0018$$

# Decoding



~~$$( \text{noun}, \text{noun} ) = a_{\text{start}, \text{noun}} b_{\text{noun}, \text{ילד}} a_{\text{noun}, \text{noun}} b_{\text{noun}, \text{ירח}} = 0.8 * 0.3 * 0.9 * 0.7 = 0.1512$$~~

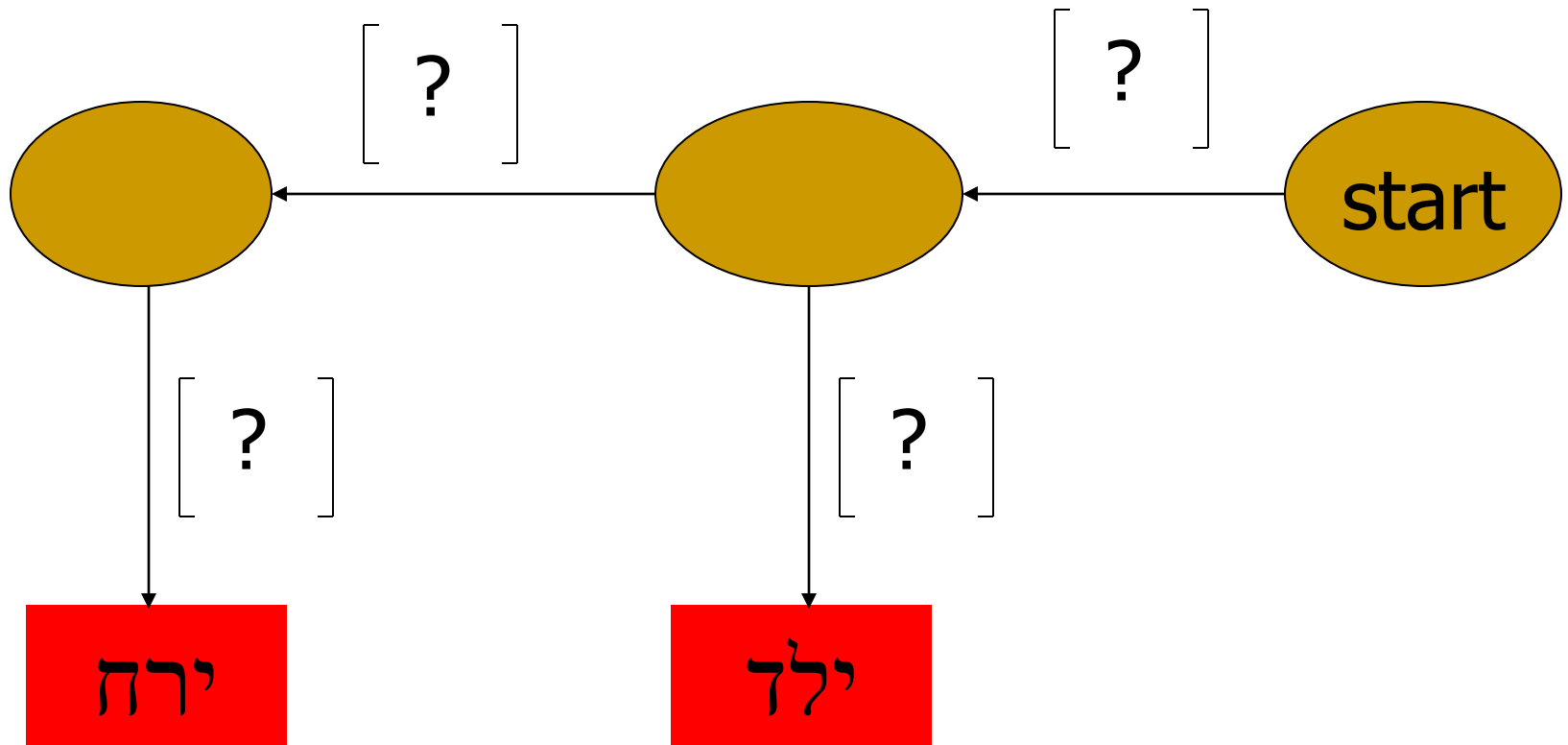
~~$$( \text{noun}, \text{verb} ) = a_{\text{start}, \text{noun}} b_{\text{noun}, \text{ילד}} a_{\text{noun}, \text{verb}} b_{\text{verb}, \text{ירח}} = 0.8 * 0.3 * 0.1 * 0.1 = 0.0024$$~~

~~$$( \text{verb}, \text{noun} ) = a_{\text{start}, \text{verb}} b_{\text{verb}, \text{ילד}} a_{\text{verb}, \text{noun}} b_{\text{noun}, \text{ירח}} = 0.2 * 0.9 * 0.9 * 0.7 = 0.1134$$~~

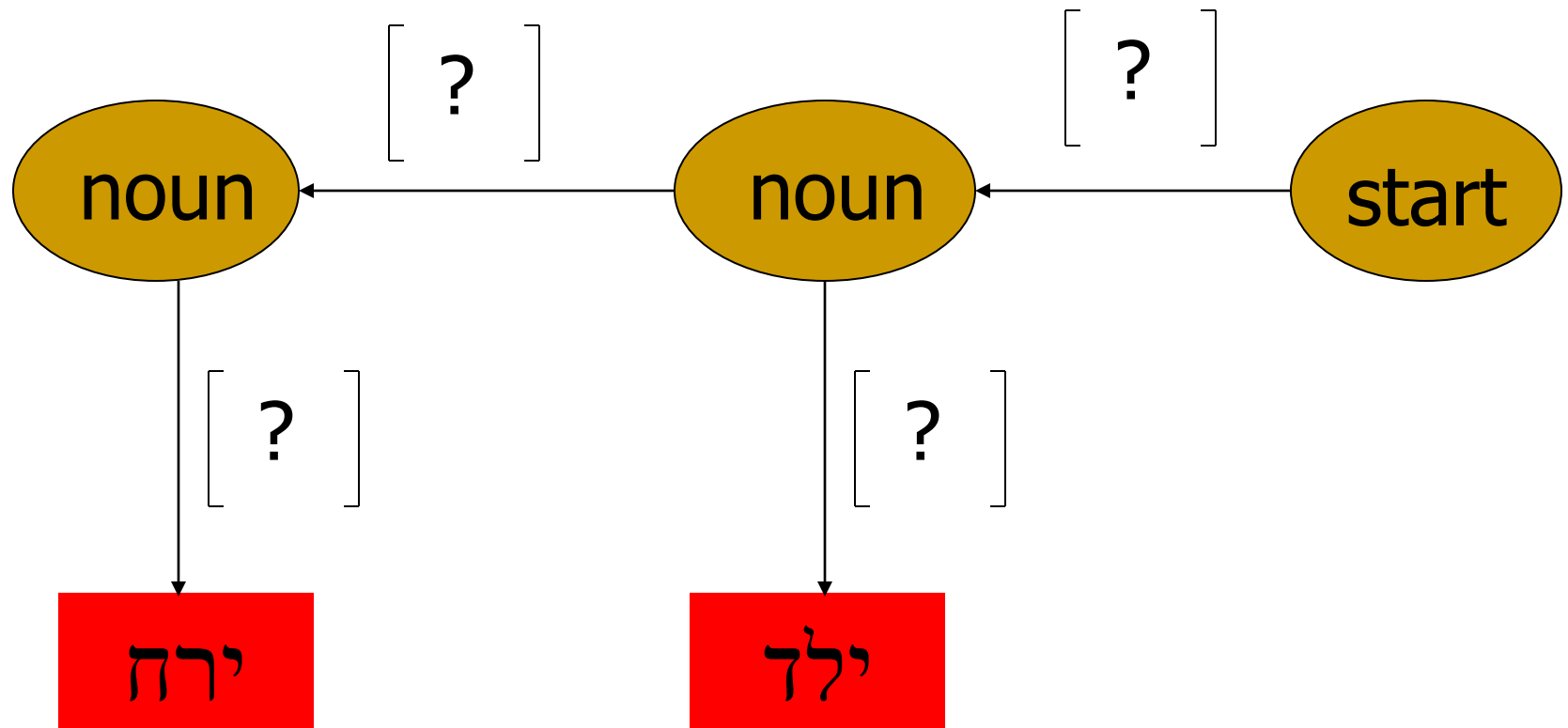
~~$$( \text{verb}, \text{verb} ) = a_{\text{start}, \text{verb}} b_{\text{verb}, \text{ילד}} a_{\text{verb}, \text{verb}} b_{\text{verb}, \text{ירח}} = 0.2 * 0.9 * 0.1 * 0.1 = 0.0018$$~~

Viterbi Algorithm (dynamic programming)

# Parameter Estimation



# Supervised Parameter Estimation



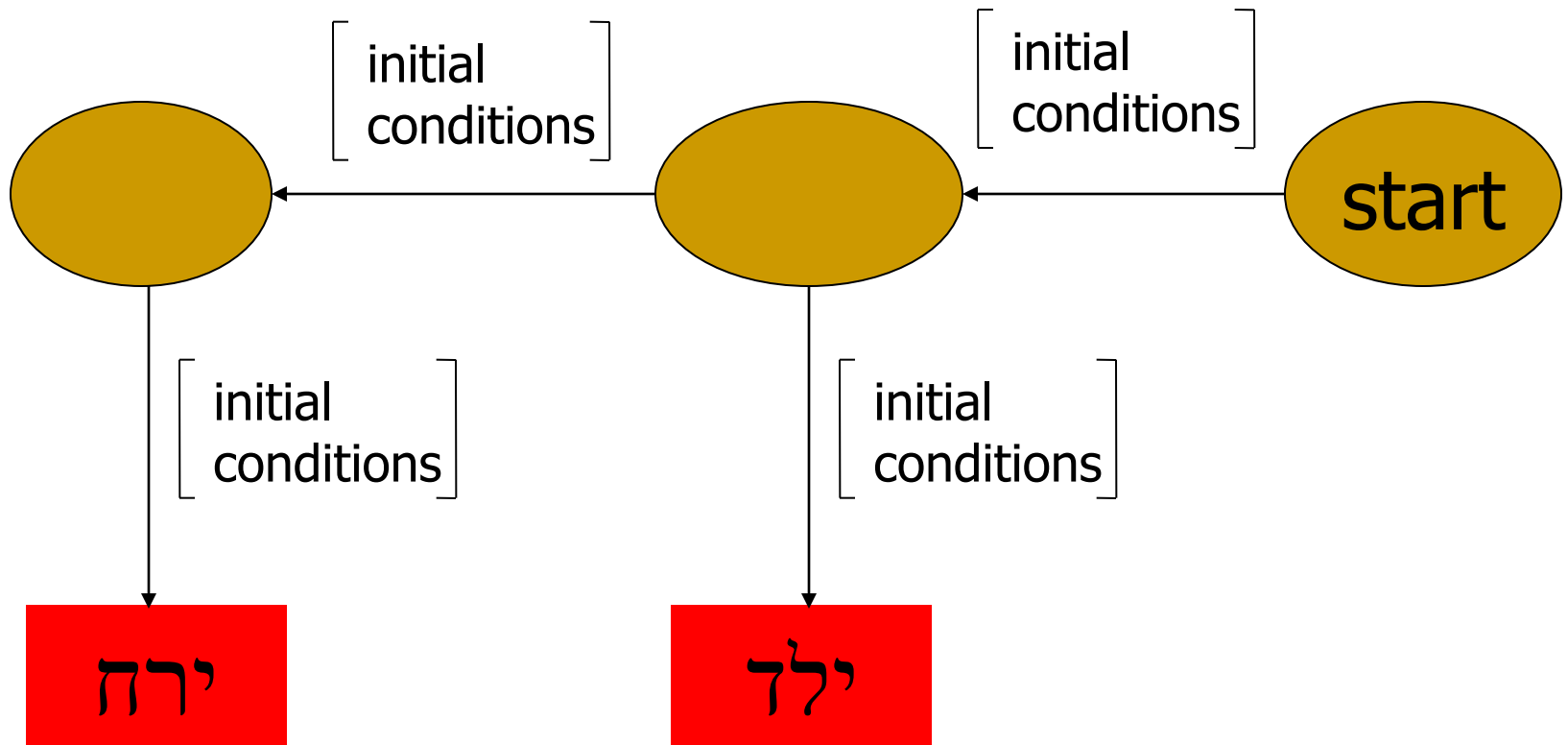


# Supervised Parameter Estimation

$$a_{i,j} = \frac{\text{number of transitions from state } i \text{ to state } j}{\text{number of transitions from state } i}$$

$$b_{i,k} = \frac{\text{number of lexical transitions from state } i \text{ to symbol } k}{\text{number of transitions from state } i}$$

# Unsupervised Parameter Estimation



# Unsupervised Parameter Estimation

$$a_{i,j} = \frac{\text{expected number of transitions from state } i \text{ to state } j}{\text{expected number of transitions from state } i}$$

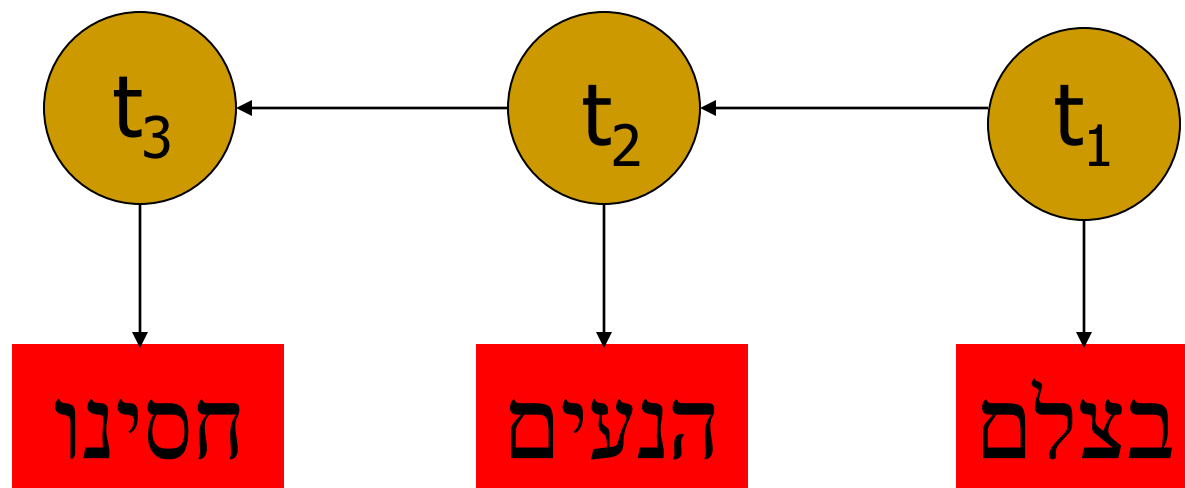
$$b_{i,k} = \frac{\text{expected number of lexical transitions from state } i \text{ to symbol } k}{\text{expected number of transitions from state } i}$$

---

# Parameter Estimation

- Baum-Welch algorithm
  - Start with a model with **initial conditions**
  - do
    - Expectation: Calculate the expected number of transitions according to the corpus and the current model.
    - Maximization: Maximize the parameters of the model according to the expected transition number.

# Token-based first order HMM

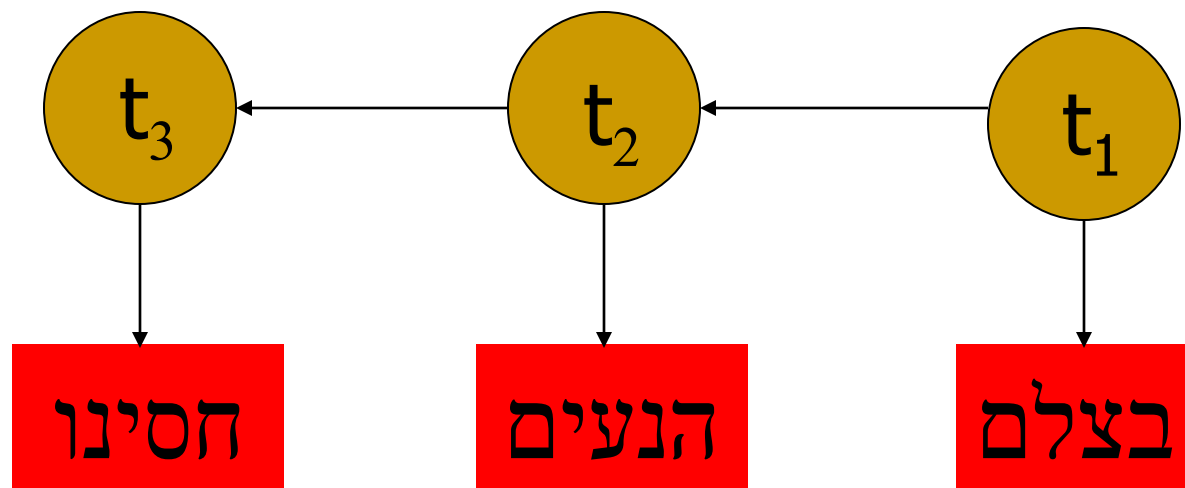


$t_1$  prep + noun.sing.masc + possessive

$t_2$  def + adj.sing.masc

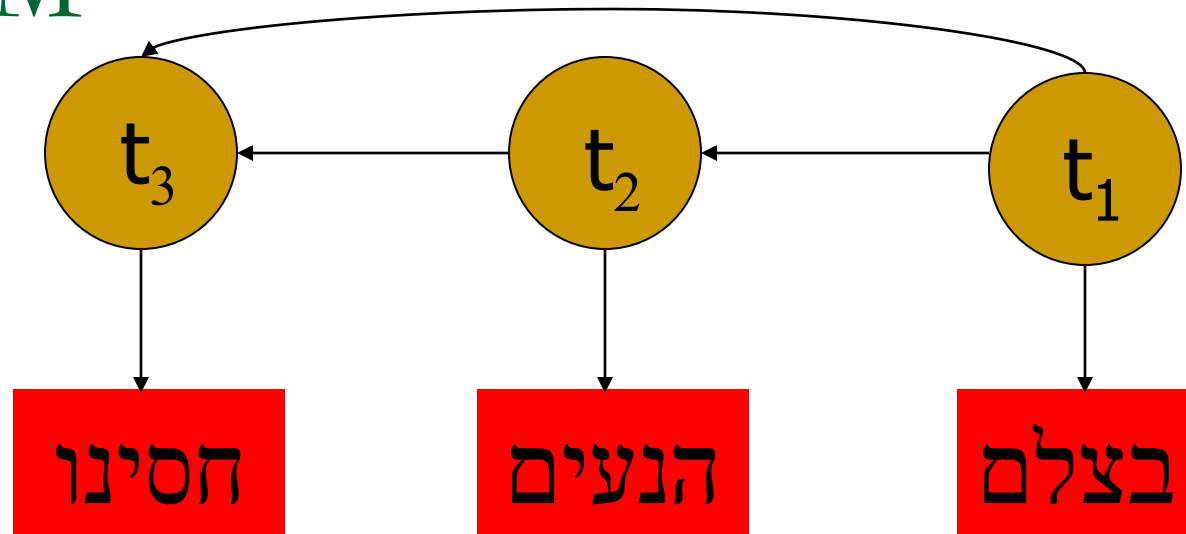
$t_3$  verb.plural.1.past

# Token-based first order HMM



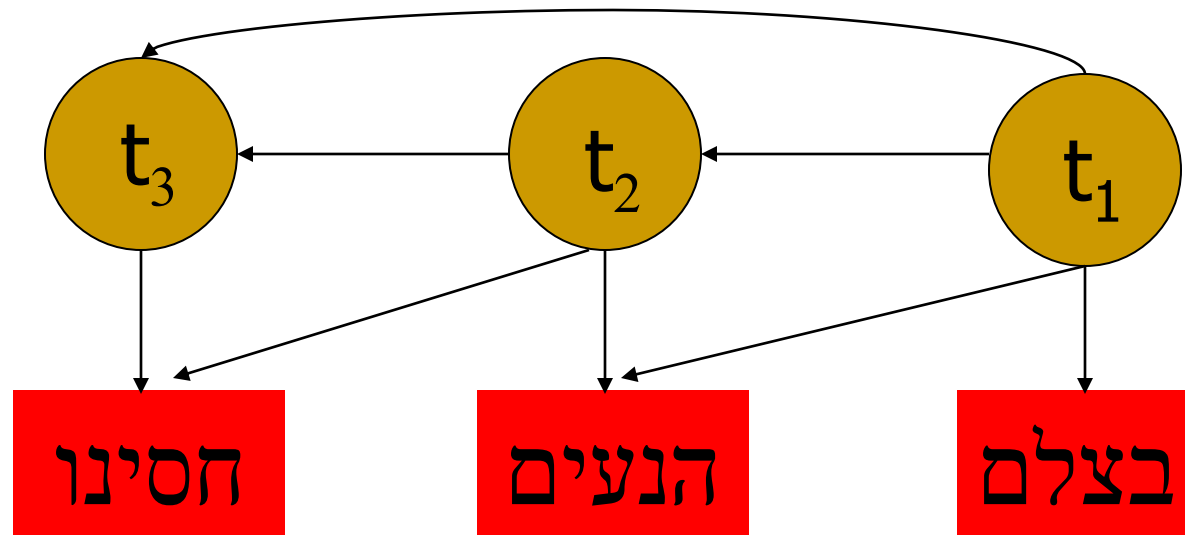
- Tags: English 48, Hebrew 3561
- State Transitions: English 1.8K, Hebrew 855K
- Lexical Transitions: English 57K, Hebrew 3.2M

# Token-based partial second order HMM



- Tags: English 48, Hebrew 3561
- State Transitions: English 38K, Hebrew 41M
- Lexical Transitions: English 57K, Hebrew 3.2M

# Token-based second order HMM



- Tags: English 48, Hebrew 3561
- State Transitions: English 38K, Hebrew 41M
- Lexical Transitions: English 300K, Hebrew 40M



---

# Word-based Model

- Computational Considerations
  - Sparse data
  - Complexity (Number of parameters)
- Linguistic Motivation
  - Adequate representation
  - Dynamic nature of the language

# Hebrew Word Definition

- Preposition prefix ב כ ל מ
  - בית, בבית
- Conjunctions ו ש כש לכש בכש
  - בית, ובית שבית
- Definite article ה
  - בית, הבית
  - הלא כל כך נחמדים
- Pronoun suffix
  - בית, ביתו

---

# Hebrew Word Definition

## ■ Prepositions

- לפני, על ידי
- בעקבות, מטעם, בגדר, במסגרת

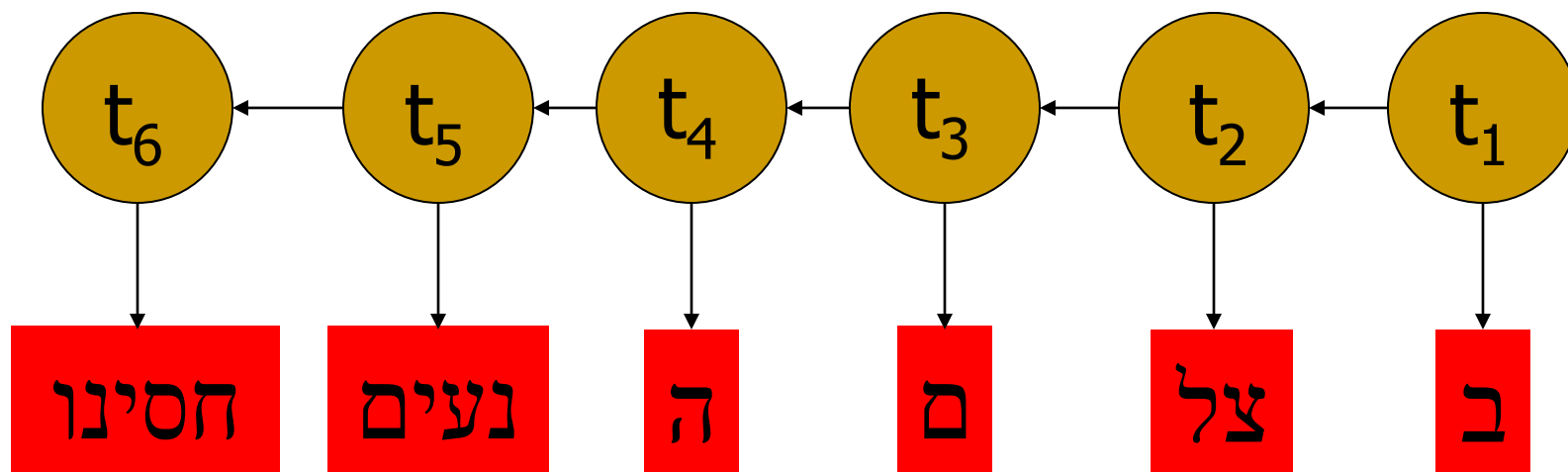
## ■ Adverbs

- במהירות, בחזרה, במשותף

## ■ Inter token words

- פה אחד, עורך דין
- אף על פי, על ידי, לבד מ, פרט ל, בנוסף ל

# Word-based first order HMM



$t_1$  prep

$t_2$  noun.sing.masc

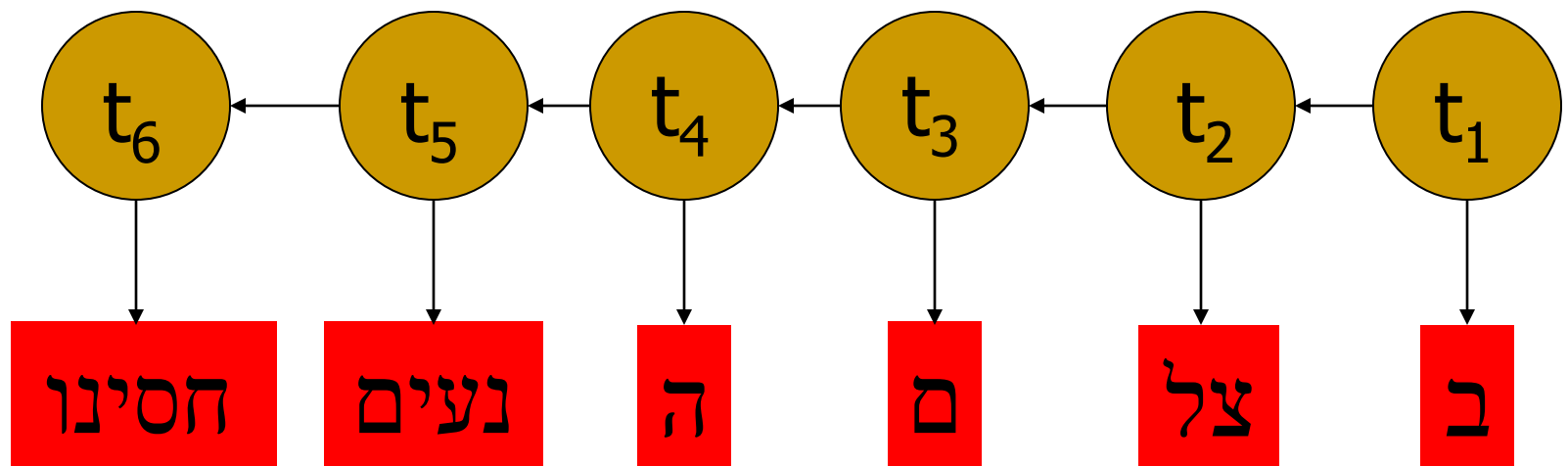
$t_3$  possessive

$t_4$  def

$t_5$  adj.sing.masc

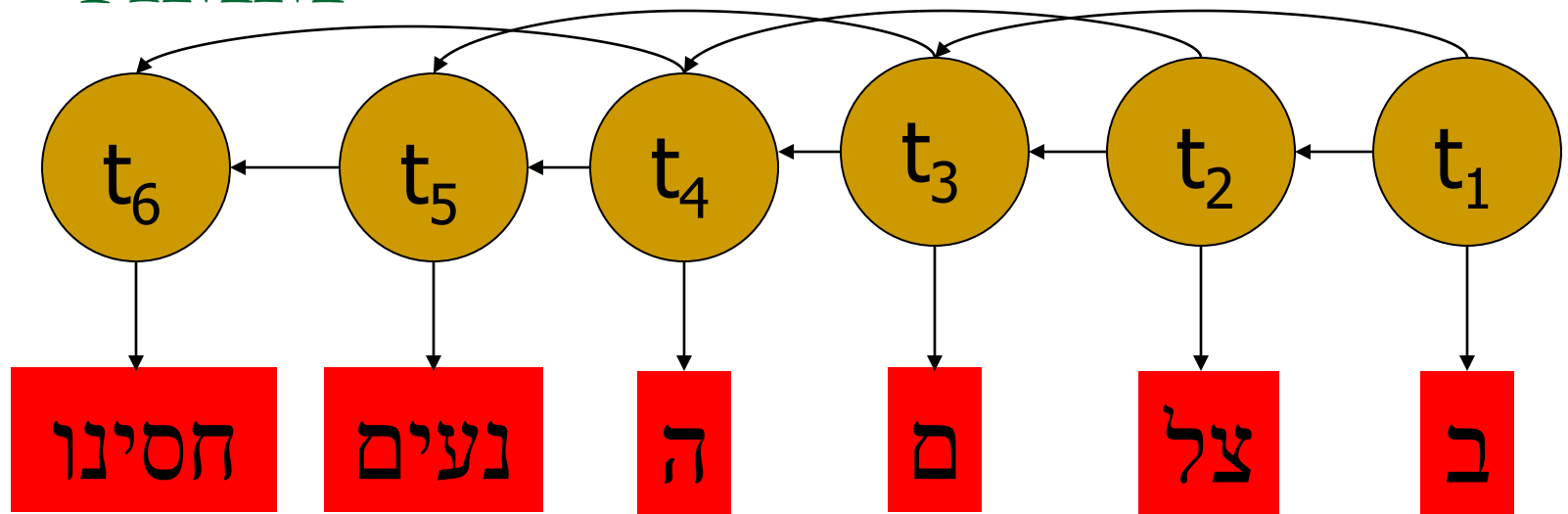
$t_6$  verb.plural.1.past

# Word-based first order HMM



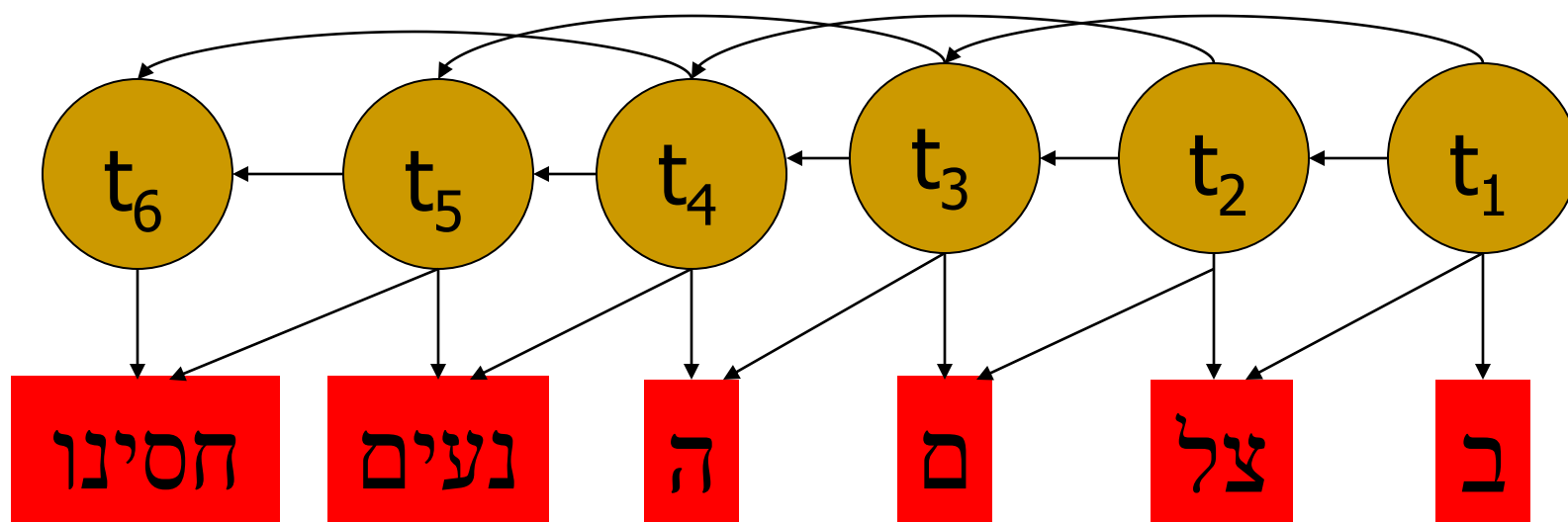
- Tags: English 48, Hebrew 362 (3561)
- State Transitions: English 1.8K, Hebrew 54K (855K)
- Lexical Transitions: English 57K, Hebrew 2.3M (3.2M)

# Word-based partial second order HMM



- Tags: English 48, Hebrew 362 (3561)
- State Transitions: English 38K, Hebrew 2.5M (41M)
- Lexical Transitions: English 57K, Hebrew 2.3M (3.2M)

# Word-based second order HMM

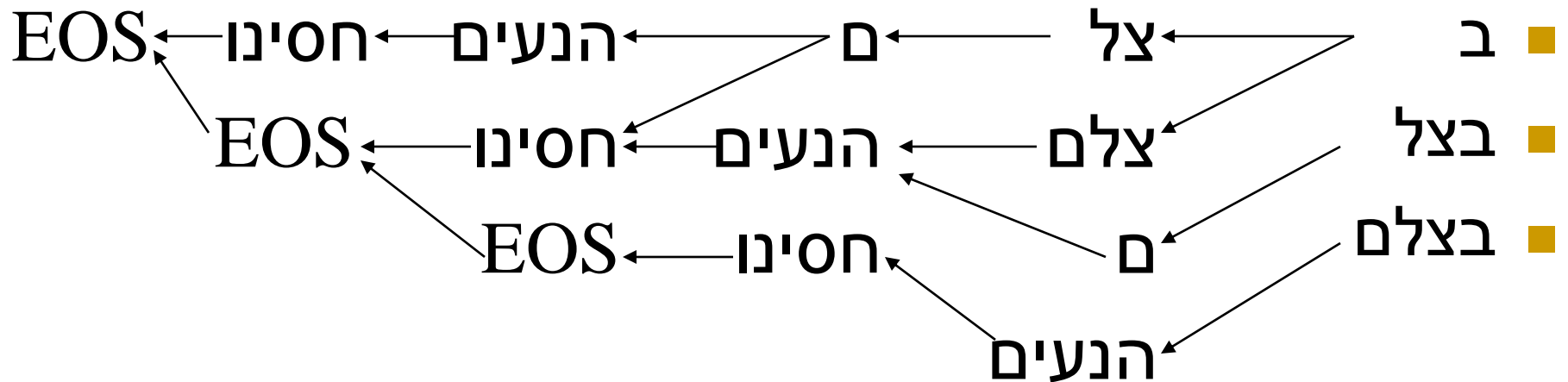


- Tags: English 48, Hebrew 362 (3561)
- State Transitions: English 38K, Hebrew 2.5M (41M)
- Lexical Transitions: English 300K, Hebrew 16M (40M)

# Text Representation

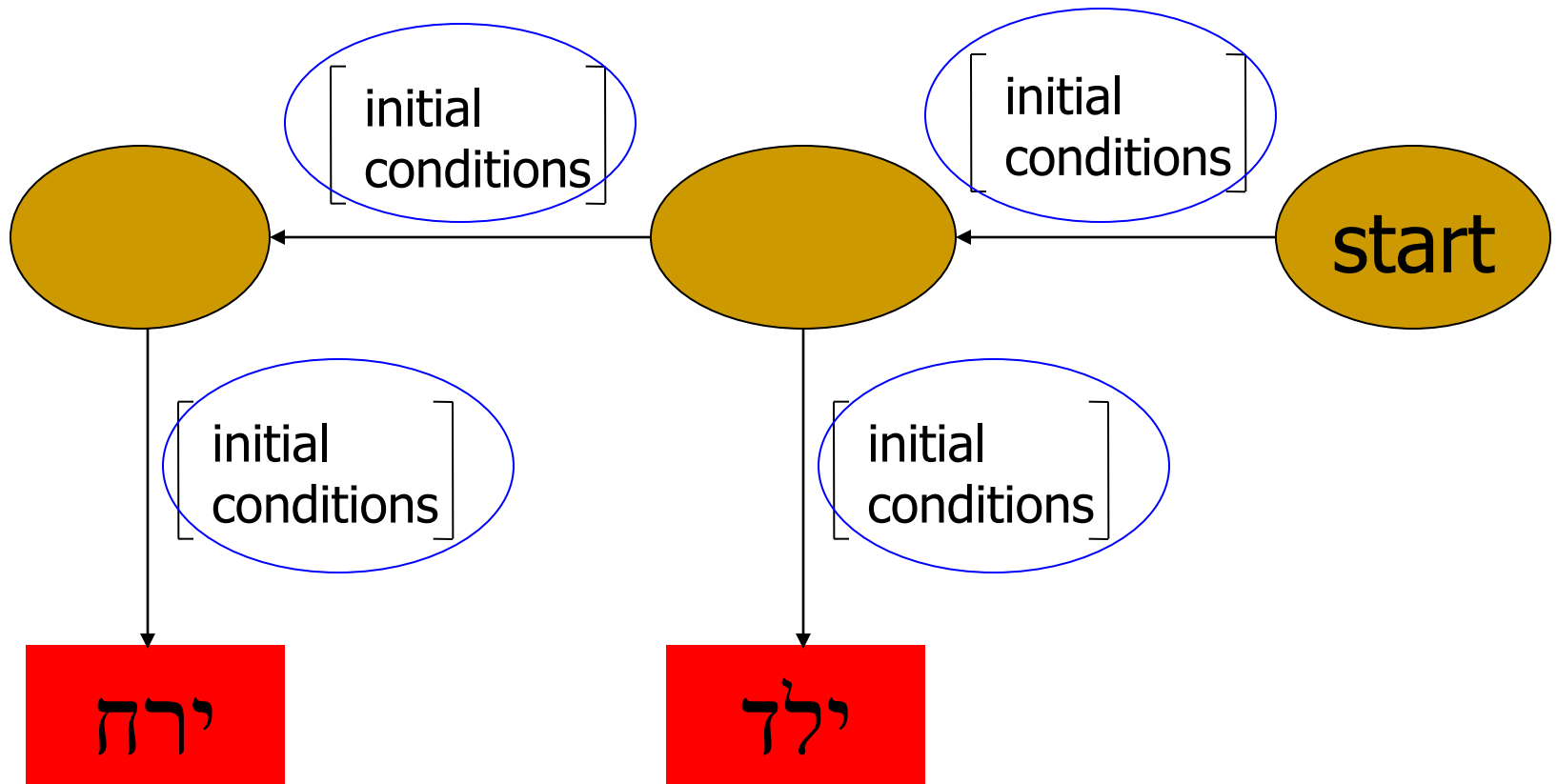
■ ב צל ם הנעים חסינו

■ בצלם הנעים חסינו





# Initial Conditions



---

# Initial Condition Types

- Morpho-lexical

- $p(t|w)$

- את: noun, preposition, pronoun

- Syntagmatic

- $p(t_i|t_{i-2}t_{i-1})$

- בצלם הנעים חסונו: probability of three consecutive verbs

---

# Morpho-lexical approximations

- Morphology-based [Levinger et al. 95]

את יכולה להעביר לי את את החפירה?

- pronoun
- preposition
- noun

# Morphology-based Approximations

- Similar word sets

- noun: האת, אתי, אתים: noun
- pronoun: אתה, אתם, אתן: pronoun
- Preposition: /

- The approximation of  $p(t|w)$  is based on the frequencies of the similar words of  $t$  in the corpus

# Linear context-based approximations

## ■ Motivation

### □ שלך

- preposition: שלך שלכם שלכן
- noun: שלך שלכם שלכן

## ■ Method

### □ לכובע שלך שלוש פינות

- $p(\text{preposition} \mid \text{שלוש}, \text{לכובע})$
- $p(\text{noun} \mid \text{שלוש}, \text{לכובע})$

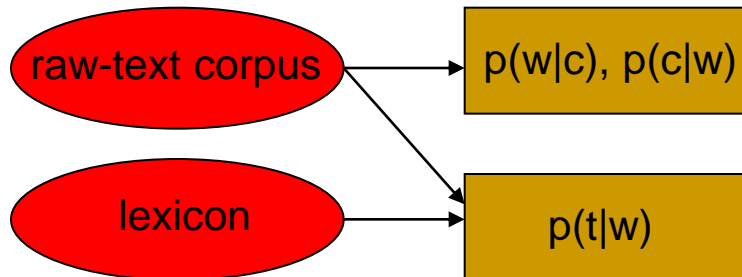
### □ Observed Data

- $p(w|c), P(c|w)$  שלך שלוש
- $p(t|w)$  - similar words algorithm

### □ Expectation/Maximization over $p(t|w)$ and $p(t|c)$

# Linear-context Model

- Notation:  $w$  – word,  $c$  – context of a word,  $t$  - tag
- Initial Conditions



- Expectation

$$p(t|c) = \sum_w p(t|w)p(w|c)$$

- Maximization

$$p(t|w) = \sum_c p(t|c)p(c|w)$$

# Syntagmatic Approximations

## ■ Syntagmatic Constraints

- a construct state form cannot be followed by verb, preposition, punctuation, existential, modal, or copula  
שרות התקדמה, פרס את נעמי בע"מ
- a verb cannot be followed by prep של  
ילד של חלום
- copula and existential cannot be followed by a verb  
יש ילד בגן, דני הוא ילד חמוד
- a verb cannot be followed by another verb (with some exceptions)  
ילד חלם חלום

# Syntagmatic Approximations

## ■ Initial Transitions

- A small seed of randomly selected sentences (10K annotated tokens)
- Tag trigram and bigram counts (ignoring the tag-word annotations) are used for initialization of  $p(t|t_{-2}, t_{-1})$  distribution



Syntagmatic	Morpho-lexical	Full	POS
Unif	Unif	87.1	91.9
	Morph	88.0	92.1
	Linear	87.5	92.4
	Morph+Linear	88.0	92.8
Pair Const	Unif	87.8	92.0
	Morph	88.1	91.8
	Linear	87.5	92.4
	Morph+Linear	88.0	92.8
Init Trans	Unif	89.4	92.6
	Morph	90.0	93.0
	Linear	89.7	93.0
	Morph+Linear	<b>90.0</b>	<b>93.1</b>

Baseline:  
Token-based,  
EM learning

# Hebrew Tagging - Analysis

- EM learning
  - Unsupervised HMM learning on word model
  - EM is very effective for Hebrew: error reduction of 65% over uniform initial conditions
- Morphology-based Initial Conditions:
  - Error reduction of 7.7% upon uniform distribution
- Syntagmatic Initial Conditions:
  - Pair constraints only have minor impact.
  - Initial transition frequencies: error reduction of 16.5% for full analysis; 12.5% for POS tagging

---

# Outline

- Objectives
- Topic Analysis with LDA
- Obtaining Precise Morphology in Hebrew
- Combining LDA and Morphological Analysis
- Using Topic Models for Search
- Evaluating Topic Models
- Next Steps

---

# Combining LDA and Morphology

- LDA picks up patterns of word co-occurrence in documents.
- Heavy variations in Hebrew could mean we “miss” co-occurrence if we do not first analyze morphology.

→ What is the best method to combine LDA and Morphological analysis?

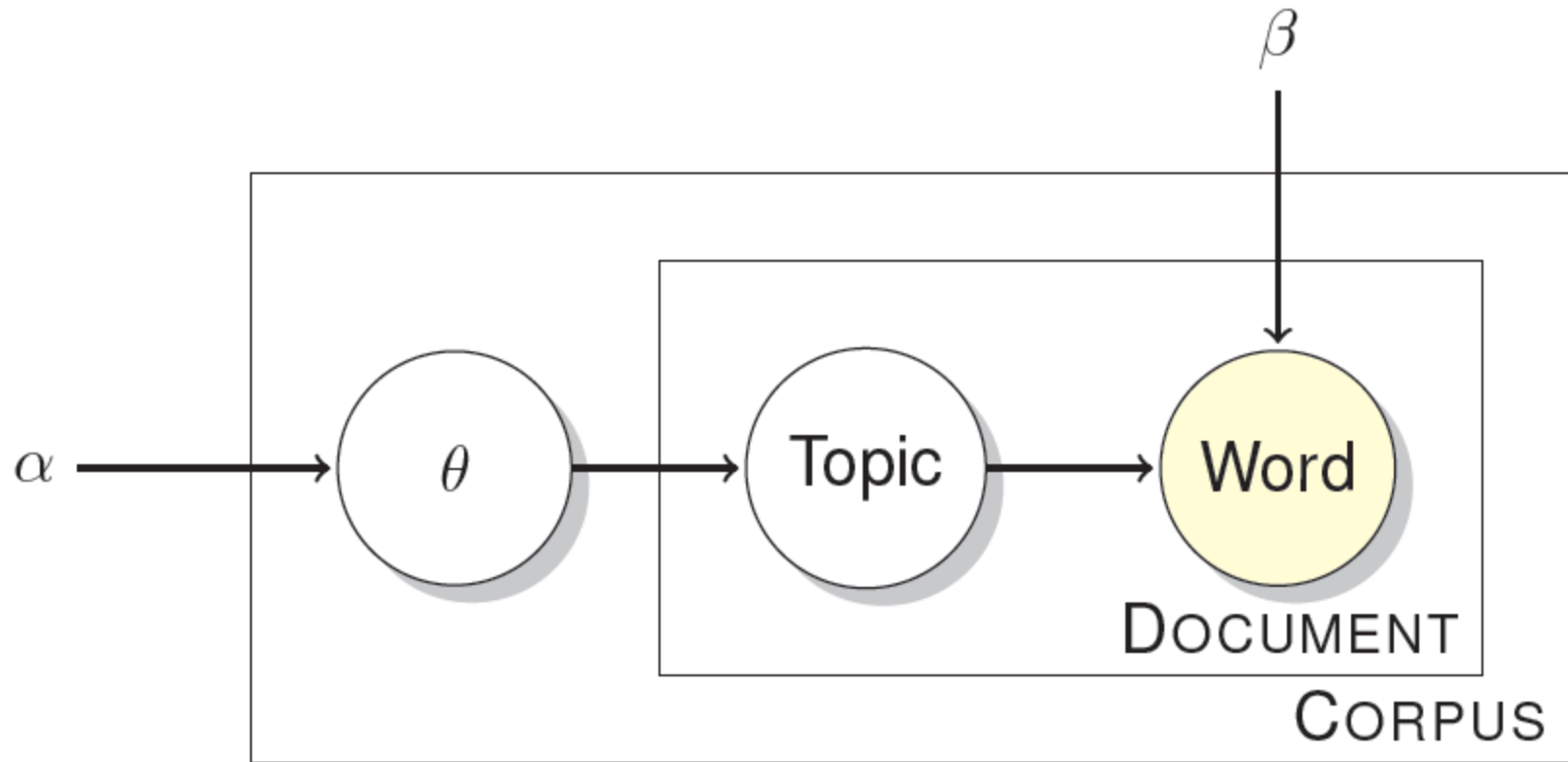
---

# Combining LDA and Morphology

3 options:

- **Ignore morphology** – token-based LDA
- **Pipeline** – resolve morphological ambiguities, then learn LDA.
- **Joint** – learn LDA on distributions of possible morphological analyses

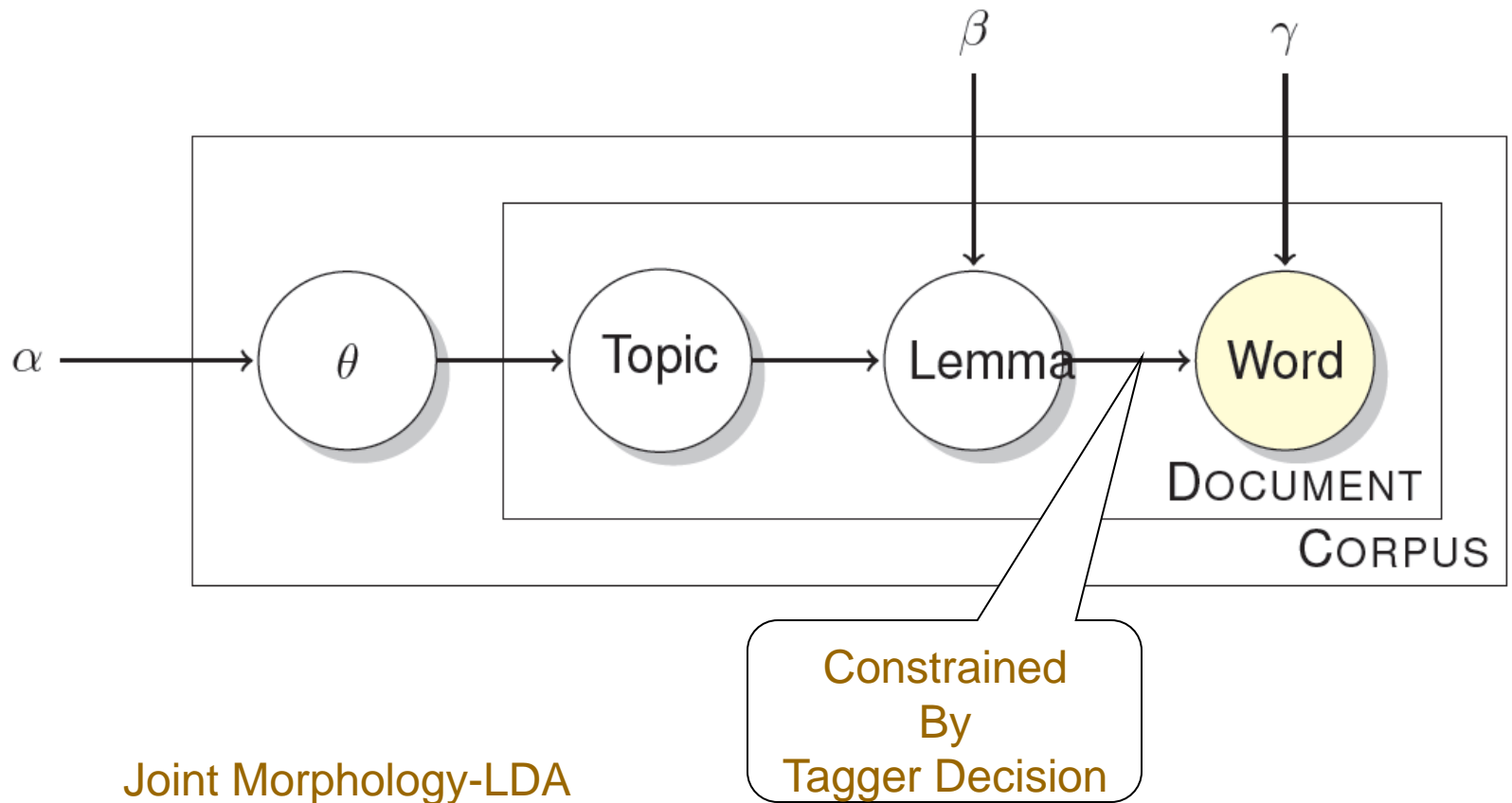
# Joint LDA-Morphology Learning



Standard token-based LDA

Combining LDA and Morphology

# Joint LDA-Morphology Learning



---

# Outline

- Objectives
- Topic Analysis with LDA
- Obtaining Precise Morphology in Hebrew
- Combining LDA and Morphological Analysis
- Using Topic Models for Search
- Evaluating Topic Models
- Next Steps



# Searching with Topics

Combine search and browse:

- Word  $\rightarrow$  Topics **Search**
- Topic  $\rightarrow$  Documents **Browse**

Word  $\rightarrow$  Topics **Disambiguate**

Cluster unseen documents based on topics

# Mishne Torah Topics

- 100 topics (K parameter)
- Each word covers many variations

- 0 עשור קדה סלע מעה כסף שני ירושלים פרי שוה ח' ל אָמַר מע"ש ח'טש נָצָא הוֹסִיף נָתַן פֶּרַט דִּינָר זָקֵב חֲלָל
- 1 כּוֹכַב עֲבוּדָה הִנָּאָה עֲבָד עוֹבֵד אָסֵר צוּרָה עֲשֶׂה אֲבָן לֶה זֶה דָרָךְ בִּעַל יִשְׂרָאֵל עוֹלָם נִסְקַל בֵּה נֶאֱסָר לָקֵחַ בְּנֵה
- 2 כִּתְבָה בִּעַל לֶה אִישׁ נָכַס נָתַן הוֹצִיא מִזֶּוֹן דִּין אוֹתָהּ הִכָּה מֵת ל' א נָטַל יֵשׁ בֵּית עֶקֶר בֵּין אַחֵר מִכָּר
- 3 מַעֲלָה מִאָּה הִנֵּה חֶלֶק שֵׁשׁ הוּא עֲשָׂרִים שְׁלֹשׁ לֹשׁ יִהְיֶה עֶשֶׂר מִנֵּה יָדַע אַחַת הוֹסִיף אַרְבָּעִים אֲמָצַע כ' ל שְׁנֵי חֲמֹשׁ
- 4 שָׁלַח הַבִּיא קִרְבָּן חֲטָאָה שְׁלֹם דָּם הִפְרִישׁ נָפֵל אָשָׁם מֵעַה עוֹף מֵת זֶה נִדְבָה קִדָּשׁ אַחֵר מֵעַה שְׁתֵּים אָמַר קָרֵב
- 5 עַד אָמַר הַעִיד נֶאֱמַן פְּלוֹנִי דִין בָּא דָבָר פְּנִים יָדַע הִכָּה מֵת פֶּה זֶה אַחֵר עִדוֹת שָׁנָה הִנֵּה אֶחָד כֵּךְ
- 6 בִּעַל הַחֲזִיר גִּזְלוּ לוֹ נָטַל שֶׁבַח גִּזְלוֹן מַעְצָא יָד לָקַח חֵיב סִימָן כְּלֵי גִזְלוֹהַ נָתַן מִלֶּךְ מִמֶּנּוּ חֲזָקָה מִן הַשׁוֹת
- 7 הַצָּר בֶּן מִבּוּי עָרֵב בֵּית הַשׁוֹת אֶחָד פִּתַח חֲלוֹן אֲנָה הִכָּה חֲבִירוֹ יָכֹל כ' תֵּל עָלָיו עָרֹב עִם צָרִיךְ אִישׁ בִּעַל
- 8 טֵמֵא מֵת טַמְאָה טְהוֹר נָטַע הוּא אָדָם כְּלֵי טֵמֵא מִגֵּעַ א' הֵל אֲנָה זָב מַשָּׂא בְּגָד עֲצָם בֵּין מִשְׁכַּב כ' ל בּוֹ
- 9 כֵּד הוֹצִיא שִׁיעוֹר חֵיב תוֹלְכָה שׁוֹם פֶּה כַּגְרוּגֶרֶת חֵי נָתַן קָטָן גַּב חֲבַל בְּהֵמָה עֵץ בִּשְׁלַל עוֹר חֵץ בִּישָׁה שְׁתֵּים
- 10 תְּשׁוּבָה עֲשֶׂה עֲבָרָה כִּשְׁעַת חֲבִירוֹ חִיטָא יָד סוֹר קָבַל יִשְׂרָאֵל מִמּוֹן בִּעַל בֵּין סִמְךָ לוֹ עוֹן חוֹטֵא חֲשׂוֹד מִימָה גָּמוֹר
- 11 לֶחֶם הַקִּרְיָב קִרְבָּן כִּבְשׁ שָׁלַח הוּא בָּא כ' ל מִן הוֹדָה עֶשֶׂר שְׁלֹם מִנְחָה קוֹרִים עִם נֶאֱכַל חֲטָאָה עֲבָר מוֹסֵף הַסִּכִּים
- 12 שְׁעָר חוֹב לוֹ כִּתְבָ גֵבָה אָמַר הִלְוֵה פְּרַע נָתַן הַמְלוּהַ יָד מִן בִּעַל לָוֶה עָרֵב דִּין יֵשׁ נָכַס מִנֵּה עֶרֶף
- 13 טֵים נָפֵל כְּצוֹן תוֹךְ סֵאה אֵת אֵינֹה הַכְּשֶׁר יָכֵד אָבָה בּוֹ הַמְקוּהַ מְקוּהַ מִשְׁק תְּלוּשׁ הוֹגֵה הִנֵּה עֲלִיהֵן מֵהֵן גָּשָׁם
- 14 יוֹם דָּם הִכָּה נָד טֵמֵא זָב טְהוֹר אִישׁ זֶה אַחֵר עֵת בְּדָק לֶה תִּלְהַ כְּתָם זִיבָה מַעְצָא שְׁמֵשׁ נִמְצָא וְסֵת
- 15 טַמְאָה הִנֵּה טֵמֵא מִפֶּה בֵּית טְהוֹר הוּא תוֹךְ גַּב בֵּין כְּלֵי פִתַח א' הֵל כ' ל תַּחַת אֵת נֶגֶד הַר פְּנִים אֶרֶץ
- 16 עֶשֶׂר פֶּאֶרֶק הִלְכָה יָצָא אַרְבָּעָה קֶצָא חֲמֹשׁ מִנֵּה נִשְׁאַר עֲשִׂירֵי תְּשִׁיעָה דִּיר שֵׁשׁ סִפֵּר פִּתַח שֵׁם מִגְלָה סִלִּיק עֲשָׂרִים תוֹךְ
- 17 שְׁבִיעֵי פְרִי אֶרֶץ שֶׁכַח אֵילָן ע' מָר נָטַע שָׁנָה עָרַל שָׁנָה אֵילָן עֲשֶׂה מוֹצָא רִיבַע הַתֵּר עֲבוּדָה זִית קֶצֶר ר"ה עֶקֶר

# Mishne Torah Topics

- Each word covers many variations

0 עֶשֶׂר פְּדָה סְלַע מַעֵה כֶּסֶף שְׁנֵי יְרוּשָׁלַיִם פְּרִי שׁוֹה ח' ל' אָמַר מֵע"ש ח' מִשׁ יֵצֵא הוֹסִיף נָתַן פֶּרֶט דִּינָר זָקֵב חֲלָל

1 כּוֹכַב עֲבוּדָה הַנֶּאֱמָר  יֵצֵא יוֹצֵא שׁוֹצֵא וְיֵצֵא תֵצֵא לְצֵאת יוֹצֵאִין יֵצְאוּ הַיּוֹצֵא וְיֵצֵא יוֹצֵאָה שׁוֹצְאוּ וְיֵצְאוּ צֵא הַיּוֹצֵאִין וְתֵצֵא וְיֵצְאוּ יוֹצֵאִין יֵצֵאת שְׁתֵצֵא שׁוֹצֵאת וְיֵצֵאת יוֹצֵאת וְהַיּוֹצֵא וְיֵצְאָה יוֹצֵאִים יוֹצֵאת כְּשׁוֹצֵא וְיֵצְאָה שׁוֹצְאִין שׁוֹצֵאה מְשׁוֹצֵא וְכְשׁוֹצֵא יֵצְאָה וְיֵצְאָה הַיּוֹצֵאִים

עוֹלָם נִסְקַל בֵּה נֶאֱסָר לָקַח בְּנָה

2 כְּתִיבָה בַּעַל לֵה אִישׁ נָכַס נָתַן הוֹצִיא קִזּוֹן דִּין אוֹתָהּ קָצָה מֵת ל' א' נָטַל יֵשׁ בֵּית עֶקֶר בֵּין אַחֵר מְבַר

# Query Disambiguation

## ■ Various topics for 1 word

שור

- [65](#) שלם הזיק בעל נזק שור בהטת השות תב נזוק חץ פטר בור מועד נפל הנה אדם את לו חבירו דין 0.0217290799815
- [83](#) הביא מנחה מן כלי אָמַר עֲלֶיהָ נָתַן שָׁמֶן לֹאג עֲשָׂרוֹן שֶׁר לְבוֹנָה מִיִּן קִטְץ וְצָא שֵׁם חֵץ נָדָר מִזְבַּח גֵּב 0.00679501698754
- [72](#) בהטת המור בה עֲלָה פָּרִיךְ אָנָּה אָדָם מִחֲשֵׁבָה עֲגָלָה כִּשְׁר חֲשָׁב חֵי מִשְׁךְ הַכְּנִים שׁוֹר מְלֹאכָה אֵין כְּלָב בַּעַל 0.00636205899364
- [73](#) כָּאָה מִזֹּל זָרַח עֲשָׂר יֵהִי־הָאֵינָה מַעֲלָה שְׁנֵי מִן ר' אִשׁ צָפוֹן יָדַע רֵאשׁוֹן קָשָׁה גָּרַע דְּרוֹם לֵיל עֹלָם חֶלֶק מִמֶּנּוּ 0.00604686318972
- [54](#) בְּשֶׁר מִן חֵי חֶלֶב בְּהֵטָה עֹרֵף דָּם אָבֵר וַיִּת דָּג אָכַל טֵמָא טְהוֹר בִּינְיָה בַּה אָסֵר מִיִּן עֲצָם תּוֹכָה שְׁחַט 0.002
- [35](#) מוֹם בְּכוֹר בְּהֵטָה בַּעַל תְּמוּכָה זֹו אָמַר כִּי הֵן קִדְשׁ בַּה בּו אָנָּה תְּמִים עֲשֶׂה נֹולֵד קְבוּעַ בֵּית הַקִּדְשׁ מִזְבַּח ח' ל 0.00176574455562
- [92](#) גַּר אִשׁ הַדְּלִיק אֹור גֵּב תְּמֹור עֵץ כְּלֵי קִדְרָה בַּשָּׁל תְּנַכָּה- גִּחְלַת תְּבַשִּׁיל הַנִּיחַ הוֹסִיף תוֹךְ עֲלָה צָרִיךְ חֲשָׁבָה חֵם 0.00129449838188
- [66](#) קְתוּחַ סְתוּם כִּי אֵל דְּבַר מִשָּׁה אָמַר ' שְׁתִּים בֵּן יִי אִישׁ אֵת שְׁלֹשׁ עֲשֶׂה וַיְהִי כֹולן יוֹם הוּא עִישׁ 0.00125680770842
- [53](#) שְׁחַט פֶּסַח שְׁחִיעָה אָכַל שְׁנֵי אֵת סִכִּין שֵׁם עֲשֶׂה רֵאשׁוֹן שֶׁר אַחַר חֲבוּכָה עֲזָכָה דָּחָה דָּם יִשְׁחוּט עֲלֵיו ח' ל אֲרֻבָּעָה 0.00125588697017
- [82](#) נָפַל כִּי ל מָאָה ח' ל סָאָה עֵעַם אָסֵר תְּרוּמָה תוֹךְ הַתְּעַרְבַּב יִין מִיִּן אָסֹור כֶּכֶם אִי נָתַן כְּלָא עֲלָה דְּבַר עִישׁ 0.00098167539267
- [81](#) קִדְשׁ מַעַל הַקִּדְשׁ נִהְנָה דָּם הַקִּדְשׁ בְּדֵק בֵּית פְּנֵה פֶרֶט נָתַן מִזְבַּח וְצָא בְּהֵטָה ח' מִשׁ יוֹבֵל דְּבַר אוֹתָה נָפַל הוֹסִיף 0.000669792364367
- [96](#) שְׁלָם גֵּב תְּשִׁלוּם קָרָן אֲרֻבָּעָה חֲמִשׁ דָּם דִּין גֵּב כֶּפֶל נָתַן בַּעַל תְּבַב אֵת ח' מִשׁ פֶּרֶר מִן מְכַר קִנְס מִמּוֹן 0.000576701268743
- [49](#) מִים יָד הַגֵּל חֲצִץ גּוֹף ר' אִשׁ עֵין אָדָם נָטַל יָדִים עֲלָה : נָתַן עַל בְּשֶׁר עֵבֶל צָרִיךְ עֲבִילָה גֵּב שַׁעַר 0.000497636227917
- [84](#) ( ) אֵי כֹל חוּץ אָכַל : נִכְנָס דּוּ ב ג א אֵל יָדֵי נָתַן ה ז ח עֲנִי 0.000467508181393
- [64](#) קִדְשׁ אִישׁ אָמַר לִי לָה זֹו קִדְשׁ קְדוּשִׁין אֵת מְקִדְשׁ נָתַן לוֹ דְּבַר אָנָּה פֶרֶט סִפְק צָרִיךְ בַּת נִמְצָא מָנָה 0.000322684737012

---

# Outline

- Objectives
- Topic Analysis with LDA
- Obtaining Precise Morphology in Hebrew
- Combining LDA and Morphological Analysis
- Using Topic Models for Search
- Evaluating Topic Models
- Next Steps

---

# How Good are Discovered Topics?

- Difficult to evaluate LDA topics
  - Many parameters (many words, many topics)
  - Each run gives slightly different results
  - How to compare topic models?
- Task-based evaluation
  - Use topics for Summarization
- Ontology alignment evaluation
  - Compare topics with existing ontology
- Data-oriented evaluation

---

# Ontology Alignment

- Mishne Torah has existing structure:
  - Hierarchy of Book/Section/Chapter
  - Excellent indexes exist
  - Compare indexes with existing ontology.
- We find excellent alignment topic/Book-Section
  - Some topics are “cross-concerns” (witnesses – topic 5)

# Topic → Documents

- Fits the Rambam's classification

65 שלם הדיק בעל נזק שור בהמה קשות תנב נזוק חץ פער בור מועד נפל קנה אדם את לו חבירו דין □

Submit

p(doc|topic):

[\('all', \('NZIKIN', 'nzki-mmown', '000012td00017000793000000prk\\_ib.txt'\)\),0.0822931114193](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000007td00017000788000000prk\\_z.txt'\)\),0.0721220527046](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000009td00017000790000000prk\\_t.txt'\)\),0.0686546463245](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000003td00017000784000000prk\\_g.txt'\)\),0.0674988441979](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000002td00017000783000000prk\\_b.txt'\)\),0.0665742024965](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000001td00017000782000000prk\\_a.txt'\)\),0.0605640314378](#)  
[\('all', \('NZIKIN', 'chowbl-owmzik', '000007td00017000832000000prk\\_z.txt'\)\),0.0485436893204](#)  
[\('all', \('NZIKIN', 'chowbl-owmzik', '000006td00017000831000000prk\\_ow.txt'\)\),0.0483125288951](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000010td00017000791000000prk\\_i.txt'\)\),0.043458159963](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000013td00017000794000000prk\\_ig.txt'\)\),0.0429958391123](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000008td00017000789000000prk\\_ch.txt'\)\),0.042071197411](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000004td00017000785000000prk\\_d.txt'\)\),0.0390661118816](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000014td00017000795000000prk\\_id.txt'\)\),0.0314378178456](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000006td00017000787000000prk\\_ow.txt'\)\),0.0298196948682](#)  
[\('all', \('NZIKIN', 'nzki-mmown', '000011td00017000792000000prk\\_ia.txt'\)\),0.0235783633842](#)  
[\('all', \('SHOFTIM', 'snhdrrin', '000005td00017001013000000prk\\_h.txt'\)\),0.0228848821082](#)



# Other Domain: More Noise

- Applied LDA to InfoMed corpus ([www.infomed.co.il](http://www.infomed.co.il)) corpus of “popular medicine” articles.

- Need a different evaluation method

- 0 מ'ח' ע'צב של ר' אש' התקף נז'רולוג פגיעה או הפקעה מיג'רנה מערכת פרכוס ארוץ' ח' סר' עצבי תחושה מוחי אפילפסיה מספר אבוד
- 1 יום שבוש' אחר כני התחיל היה יומי ש'לושה משך ע'דן אחר תחלה ? הראשון כן כמה רק הופיע שוב נמשך
- 2 מעי הוא בע'ן צו'עה שלשול של או מערכת עיכול יצי'עה ל' א' קבל גם עצירות ג'רם ע'שה קולונוסקופיה ג'ו יש עמיל
- 3 בדיקה בצ'ע מקי'ן הוצ'עה דם ב'דק אולטרסאונד בצו'ע ע'שה אבחון ברור נב'דק משו'בה זו אבחון הר'עה נע'שה בה מקי'ף אבחנה
- 4 הלך ג'כ' ל' שלך המל'יץ רופא על קנה התאים הוא רמואי לי את קנה שאת קבל אבחנה כד טפל היסטוריה המך
- 5 ###NUMBER### - % : עד ל' ב' כ' מ' ק'על ק'ה ו' בכ' להל'ן אחו'ו ב'א נתון נחשב א' פ'ן טוב
- 6 שימוש על ל' א' השקעה השתמש מחקר אין ז'מן המל'יץ ג'ב סכון נוק' כי ע'שון כן י'דע ומר' אך אם כ'טום
- 7 הסק'עה מ'עב או של נפשי חיים עם פסיכיאטר הכו' אדם טפול יש ח' סר' ק' שי' עת תקוקה תפקוד הוא ד'כאו'ן תרוקתי
- 8 ###NUMBER### גיל הוא בת ש'נה ג'דילה ג' בה עד של עוד התפתחות נ'טוך שלי מקדם מ'ע'ע כן מ'עבר התבגרות ע'סיר ק'על
- 9 תיה ג'כ' ל' אם ל' א' ? זה א'מור המשיך ע'דן בהחלט י'תכן בגלל אבל אחר ק'עה לי צ'ון גם ע'ריך אחר
- 10 כדור לקח את הוא אם נ'על נטי'לה התחיל ? הפסיק יום ל' א' הפסקה המל'יץ א'חד כן ל'קושה משך ח'ז המל'עה
- 11 רופא קנה משקפה ע'שה ב'דק אל ברור כ'ד אותי כ'דאי המך בדיקה עליך לי ע'לה שלך א'צל מטטה מה שלח
- 12 א'כל ק'זון ש'תנה אכילה ש'תה א'לכוהול חלב כמות נ'מנע מי'ם מוצ'ר המל'יץ הכיל א'רוקה פ'רי מ'ן מאכל מ'זונה י'קקות ל' א'
- 13 כ'בד של ע'לה תפקוד ג'בות אנויים ד'ם ש'היר ( כ'טה - ) נוסף בדיקה גם יש CPK פגיעה הסק'עה על
- 14 ק'עה הוא של או קבל ו'ש'ט ג'רם את טפול כיב ת'נדק אינו בע'ן ע'הבת פילורי כן פ'ה הליקובקטר עיכול סימפטום
- 15 גוף חומר הוא ש'נה סוג של את חלק תהליך אל כ' ל' ש'ם נ'מ'עא הכיל על בהם א'ש'ר ישנם מ'רכיב כ'ך

---

# Data-oriented Evaluation

- Method derived from our previous work on ontology evaluation
- How well does an automatically constructed partition of entities into classes represent reality?
- How can a classification be improved via merge or split operations?
- **Example:** movies genres, taken from IMDB: *drama, comedy, war, sport, action...*

# Text Classifier to Evaluate Classes

- **Idea:** use a set of texts from a given domain as a proxy of the “reality” represented by the ontology.
- **Hypothesis:**  
*If* the ontology indicates that some movies are "clustered" according to one of the dimensions  
*Then* documents associated to these movies should also be found to be associated by a text-classification engine that has been trained on the classification induced by the ontology
- **Procedure:**  
Train a classifier for each class we want to evaluate  
If a classifier can decide with good accuracy whether a given text belongs to a class – the class is well-defined.  
In the movies domain, reviews are used as representative



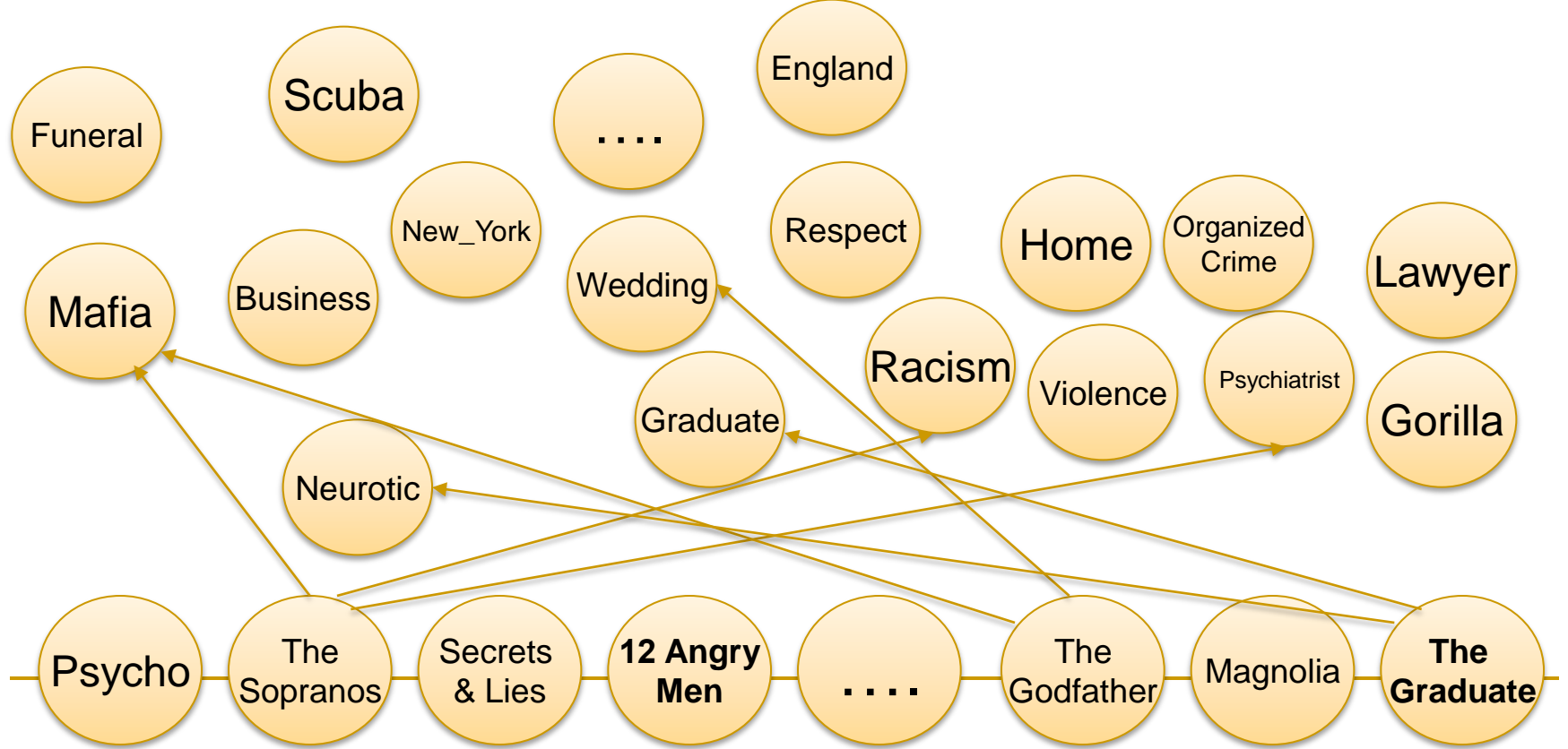
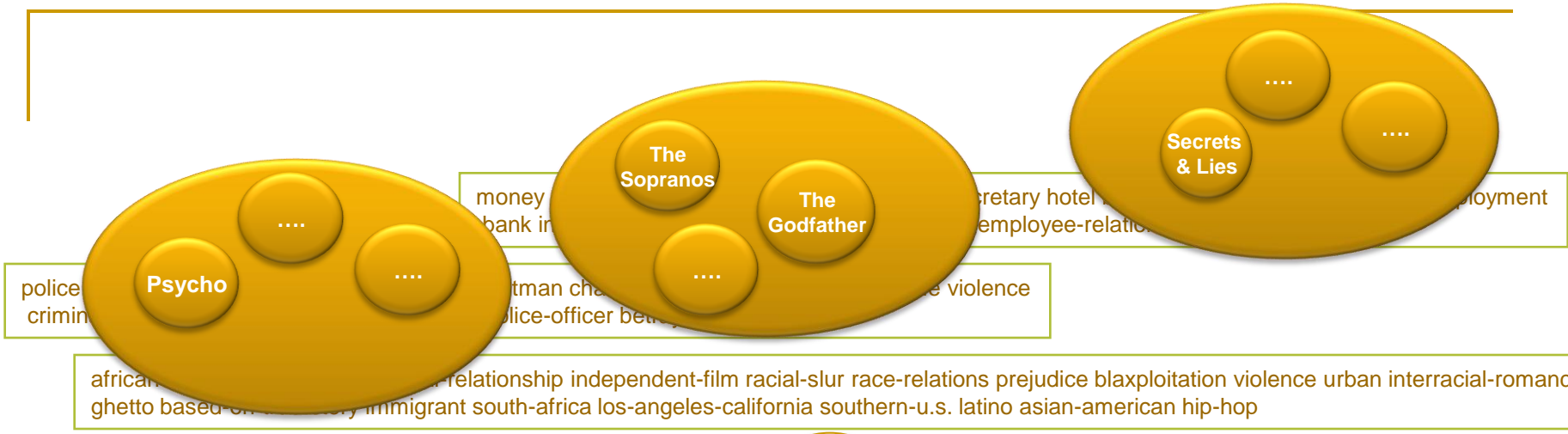
# Noisy Classes

- What can be done if the classes are very noisy?
- Example: Keywords in IMDB
  - Keyword lists are open to additions/deletions by users (yet moderated).
  - Examples of Keywords: *Mafia Business New\_York Wedding Respect Home Organized\_Crime Lawyer Violence ... (and the film?)*
  - Too few movies are associated with any single keyword
  - Movies associated with a given keyword are not necessarily related
- ***Could be the same for LDA topics in noisy domains***

---

# Cluster Keywords using LDA

- Apply LDA on the keywords
  - Divide movies into classes according to LDA distributions of the movies reviews
  - Construct classifiers to evaluate the quality of the LDA classes
  - Example of the top keywords of a “good” class: *england inheritance london-england based-on-novel mansion london maid class-differences servant 19th-century period-piece butler orphan estate aunt uncle heir victorian-era love marriage*
-



# Conclusions

- Morphological analysis is a critical pre-processing step for most text mining applications in Hebrew
- We obtain accuracy in the 90-95% range on full segmentation/POS/morphological analysis – robust on unknown words.
- Enables effective LDA Topic Analysis in Hebrew
- More work needed on Topic Analysis Evaluation