

FAST AND NON-APPROXIMATIVE
LANGUAGE MODEL PREFIXQUERIES
FOR WORD PREDICTION USING
TOP-K JOINING TECHNIQUES

BACHELOR THESIS

Lukas Schmelzeisen

lukas@uni-koblenz.de

Institute for Web Science and Technologies,
University of Koblenz-Landau

August 6, 2015

OUTLINE

1. NEXT WORD PREDICTION

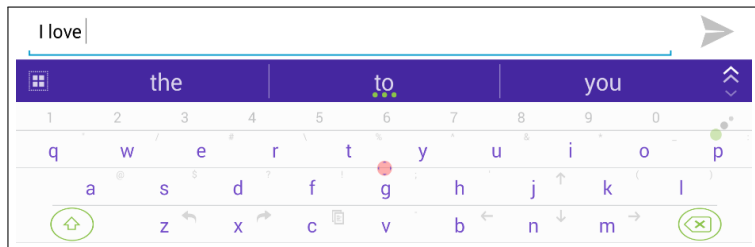
2. LANGUAGE MODELS

3. TOP-K JOINING TECHNIQUES

Section 1

NEXT WORD PREDICTION

WHAT IS NEXT WORD PREDICTION?



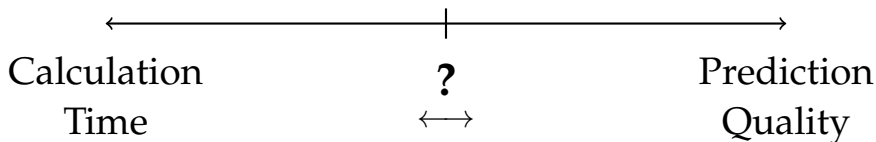
Guess the *next intended* word
from words already entered

BENEFITS OF GOOD WORD PREDICTION

- Faster typing
- Spelling / Grammar
- ...
- Metric for Language Models?

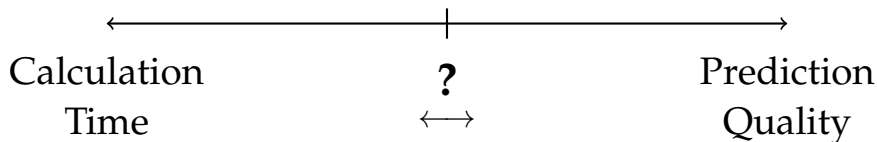
THESIS GOAL

Word prediction usually optimized for
speed **or** quality



THESIS GOAL

Word prediction usually optimized for
speed **or** quality



This thesis optimizes for speed while
maintaining best quality!

NEXT WORD PREDICTION

$$\text{NWP}(h) = \arg \max_{w \in \Sigma} P(w | h)$$

NEXT WORD PREDICTION

history of entered words



$$\text{NWP}(h) = \arg \max_{w \in \Sigma} P(w | h)$$

NEXT WORD PREDICTION

history of entered words

$$\text{NWP}(h) = \arg \max_{w \in \Sigma} P(w | h)$$

for each word in vocabulary

Section 2

LANGUAGE MODELS

HOW TO FIND PROBABILITIES?

$$P(w|h) ?$$

HOW TO FIND PROBABILITIES?

$$P(w|h) ?$$

Solution: Statistical Language Models

- analyze text corpora
- to estimate sequence probabilities

HOW TO FIND PROBABILITIES?

$$P(w | h) ?$$

Solution: Statistical Language Models

- analyze text corpora
- to estimate sequence probabilities

We will employ two state-of-the-art Language Models:

- Modified Kneser Ney Smoothing
(Chen and Goodman 1999)
- Generalized Language Model
(Pickhardt et al. 2014)

MODIFIED KNESER NEY SMOOTHING

$$P_{\text{MKN}}(w_n | w_1^{n-1}) = \frac{c^d(w_1^n) + \gamma(w_1^{n-1}) \hat{P}_{\text{MKN}}(w_n | w_2^{n-1})}{c(w_1^{n-1} \square)}$$

MODIFIED KNESER NEY SMOOTHING

$$P_{\text{MKN}}(w_n | w_1^{n-1}) = \frac{c^d(w_1^n) + \gamma(w_1^{n-1}) \hat{P}_{\text{MKN}}(w_n | w_2^{n-1})}{c(w_1^{n-1} \square)}$$

GENERALIZED LANGUAGE MODEL

$$P_{\text{GLM}}(w_n | w_1^{n-1}) = \frac{c^d(w_1^n) + \frac{\gamma(w_1^{n-1})}{\#_{\partial} w_1^{n-1}} \sum_{j=1}^{\#_{\partial} w_1^{n-1}} \hat{P}_{\text{GLM}}(w_n | \partial_j w_1^{n-1})}{c(w_1^{n-1} \square)}$$

GENERALIZED LANGUAGE MODEL

$$P_{\text{GLM}}(w_n | w_1^{n-1}) = \frac{c^d(w_1^n) + \frac{\gamma(w_1^{n-1})}{\#_{\partial} w_1^{n-1}} \sum_{j=1}^{\#_{\partial} w_1^{n-1}} \hat{P}_{\text{GLM}}(w_n | \partial_j w_1^{n-1})}{c(w_1^{n-1} \square)}$$

EXAMPLE: PROBABILITY OF “YOU” AFTER “I LOVE”

$$P_{\text{MKN}}(\text{you} \mid \text{I love}) = \frac{c^d(\text{I love you}) + \gamma(\text{I love}) \hat{P}_{\text{MKN}}(\text{you} \mid \text{love})}{c(\text{I love } \square)}$$

$$\hat{P}_{\text{MKN}}(\text{you} \mid \text{love}) = \frac{N_{1+}(\bullet \text{ love you}) + \gamma(\text{love}) \hat{P}_{\text{MKN}}(\text{you})}{N_{1+}(\bullet \text{ love } \bullet)}$$

$$\hat{P}_{\text{MKN}}(\text{you}) = \frac{N_{1+}(\bullet \text{ you})}{N_{1+}(\bullet \bullet)}$$

EXAMPLE: PROBABILITY OF “YOU” AFTER “I LOVE”

$$P_{\text{MKN}}(\text{you} \mid \text{I love}) = \frac{c^d(\text{I love you}) + \gamma(\text{I love}) \hat{P}_{\text{MKN}}(\text{you} \mid \text{love})}{c(\text{I love } \square)}$$

$$\hat{P}_{\text{MKN}}(\text{you} \mid \text{love}) = \frac{N_{1+}^d(\bullet \text{ love you}) + \gamma(\text{love}) \hat{P}_{\text{MKN}}(\text{you})}{N_{1+}(\bullet \text{ love } \bullet)}$$


$$\hat{P}_{\text{MKN}}(\text{you}) = \frac{N_{1+}(\bullet \text{ you})}{N_{1+}(\bullet \bullet)}$$

IDEA: EXPRESS AS WEIGHTED SUMS

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

IDEA: EXPRESS AS WEIGHTED SUMS

depends on argument w


$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

IDEA: EXPRESS AS WEIGHTED SUMS

depends on argument w

independent of argument w

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

IDEA: EXPRESS AS WEIGHTED SUMS

depends on argument w

independent of argument w

$$P(w | h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

all words in vocabulary

IDEA: EXPRESS AS WEIGHTED SUMS

depends on argument w

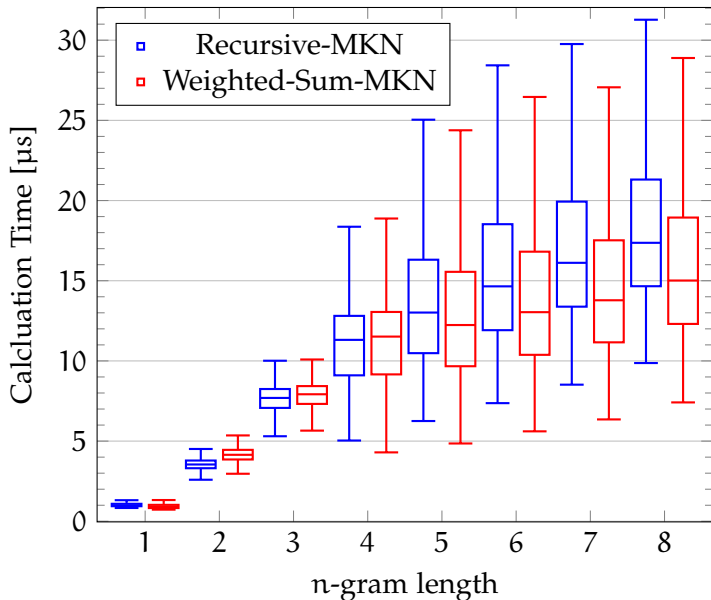
independent of argument w

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

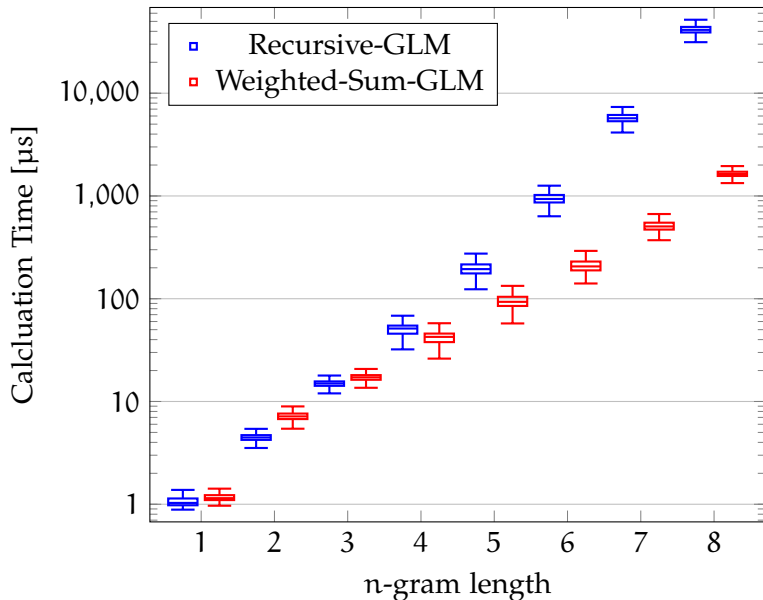
user given, thus fixed

all words in vocabulary

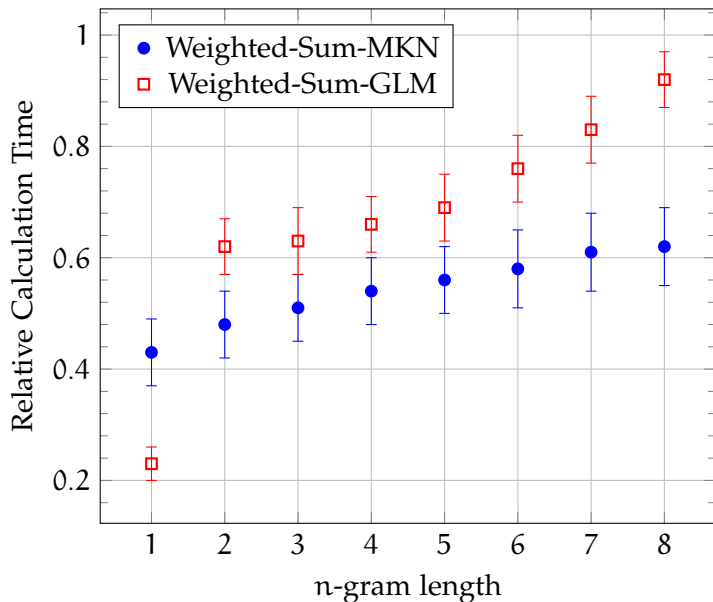
Calculation Time per Probability $P_{\text{MKN}}(w|h)$



Calculation Time per Probability $P_{\text{GLM}}(\mathbf{w} | \mathbf{h})$



Relative Weight λ_i^h Calculation Time per Probability

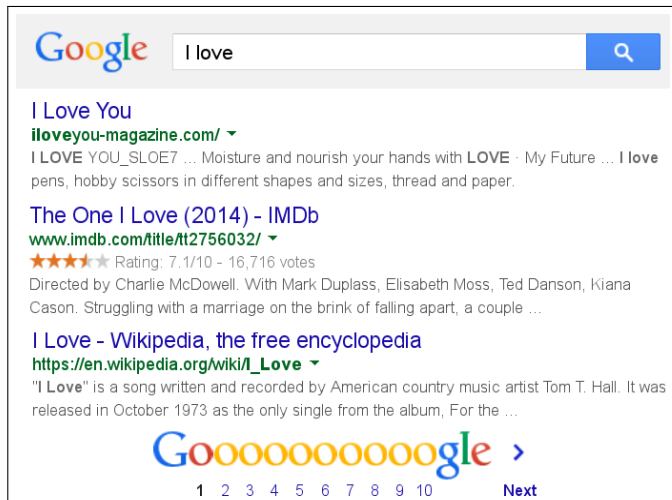


Section 3

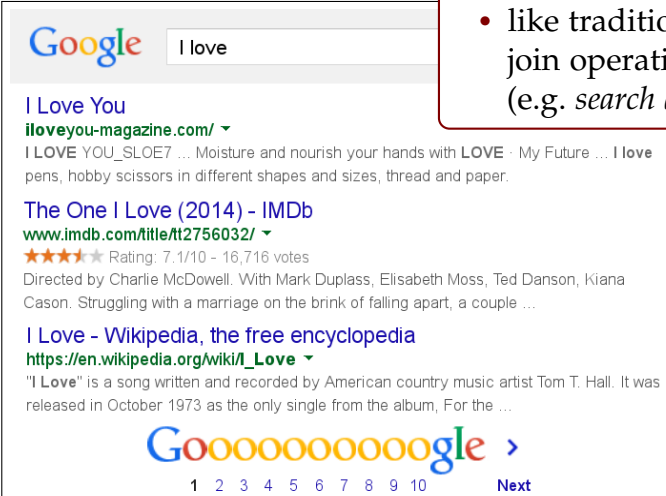
TOP-K JOINING TECHNIQUES

WHAT ARE TOP-K JOIN QUERIES?

WHAT ARE TOP-K JOIN QUERIES?



WHAT ARE TOP-K JOIN QUERIES?



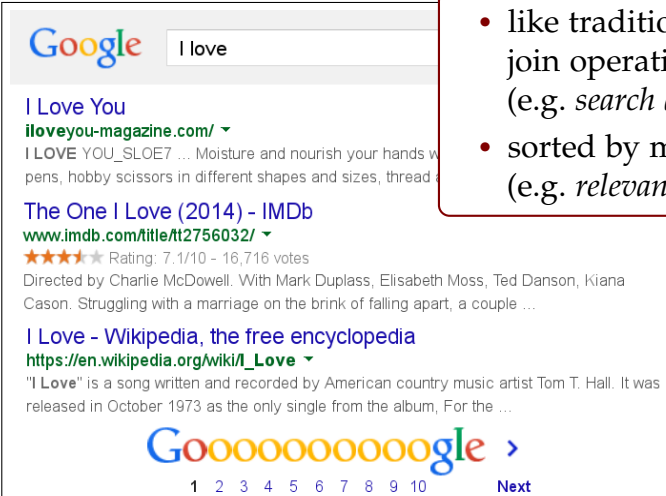
The screenshot shows a Google search interface with the query "I love" entered in the search bar. The results are as follows:

- I Love You**
iloveyou-magazine.com/ ▾
 I LOVE YOU_SLOE7 ... Moisture and nourish your hands with LOVE · My Future ... I love pens, hobby scissors in different shapes and sizes, thread and paper.
- The One I Love (2014) - IMDb**
www.imdb.com/title/tt2756032/ ▾
 ★★★★★ Rating: 7.1/10 - 16,716 votes
 Directed by Charlie McDowell. With Mark Duplass, Elisabeth Moss, Ted Danson, Kiana Cason. Struggling with a marriage on the brink of falling apart, a couple ...
- I Love - Wikipedia, the free encyclopedia**
https://en.wikipedia.org/wiki/I_Love ▾
 "I Love" is a song written and recorded by American country music artist Tom T. Hall. It was released in October 1973 as the only single from the album, For the ...

At the bottom of the results, the Google logo is displayed with the letters "Goooooooogle" where the "o"s are elongated, followed by a "Next" button.

- like traditional join operations (e.g. *search query*)

WHAT ARE TOP-K JOIN QUERIES?



A screenshot of a Google search interface. The search bar contains the text "I love". Below the search bar, several search results are displayed. The first result is "I Love You" from iloveyou-magazine.com/. The second result is "The One I Love (2014) - IMDb" with a rating of 7.1/10. The third result is "I Love - Wikipedia, the free encyclopedia" with a URL. At the bottom of the screenshot, the Google logo is shown with the letters "Goooooooogle" where the 'o's are elongated, and a "Next" button is visible.

Google

I love

I Love You
iloveyou-magazine.com/ ▼
 I LOVE YOU_SLOE7 ... Moisture and nourish your hands w
 pens, hobby scissors in different shapes and sizes, thread a

The One I Love (2014) - IMDb
www.imdb.com/title/tt2756032/ ▼
 ★★★★★ Rating: 7.1/10 - 16,716 votes
 Directed by Charlie McDowell. With Mark Duplass, Elisabeth Moss, Ted Danson, Kiana
 Cason. Struggling with a marriage on the brink of falling apart, a couple ...

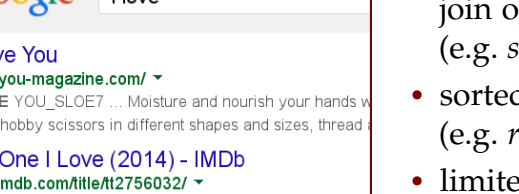
I Love - Wikipedia, the free encyclopedia
https://en.wikipedia.org/wiki/I_Love ▼
 "I Love" is a song written and recorded by American country music artist Tom T. Hall. It was
 released in October 1973 as the only single from the album, For the ...

Goooooooogle >

1 2 3 4 5 6 7 8 9 10 Next

- like traditional join operations (e.g. *search query*)
- sorted by metric (e.g. *relevance*)

WHAT ARE TOP-K JOIN QUERIES?



The screenshot shows the Google search interface with the query "I love". The results are as follows:

- I Love You**
iloveyou-magazine.com/ ▾
 I LOVE YOU_SLOE7 ... Moisture and nourish your hands w
 pens, hobby scissors in different shapes and sizes, thread a
- The One I Love (2014) - IMDb**
www.imdb.com/title/tt2756032/ ▾
 ★★☆☆☆ Rating: 7.1/10 - 16,716 votes
 Directed by Charlie McDowell. With Mark Duplass, Elisabeth
 Cason. Struggling with a marriage on the brink of falling apart, a couple ...
- I Love - Wikipedia, the free encyclopedia**
https://en.wikipedia.org/wiki/I_Love ▾
 "I Love" is a song written and recorded by American country music artist Tom T. Hall. It was
 released in October 1973 as the only single from the album, For the ...

At the bottom, the Google logo is displayed with the text "1 2 3 4 5 6 7 8 9 10" and a "Next" button.

- like traditional join operations (e.g. *search query*)
- sorted by metric (e.g. *relevance*)
- limited to $k = 3$ best results

CALCULATING PROBABILITIES FOR ALL WORDS IS SLOW!

CALCULATING PROBABILITIES FOR ALL WORDS IS SLOW!

Observation: Words with high probabilities
have high occurrence counts
(probabilities are monotone)

CALCULATING PROBABILITIES FOR ALL WORDS IS SLOW!

Observation: Words with high probabilities
have high occurrence counts
(probabilities are monotone)

Idea: Keep sorted lists of occurrence counts

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

CALCULATING PROBABILITIES FOR ALL WORDS IS SLOW!

Observation: Words with high probabilities
have high occurrence counts
(probabilities are monotone)

Idea: Keep sorted lists of occurrence counts

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

We apply two Top-k Joining Techniques to NWP:

- Threshold Algorithm (Fagin et al. 2001)
- No Random Access Algorithm (Fagin et al. 2001)

COMPARISON

Threshold Algorithm

- Requires Sorted and Random Access
- + Faster
- + Less Memory

No Random Access Algorithm

- + Only Sorted Access necessary

DATA STRUCTURE

DATA STRUCTURE

- Sorted lists are unique for each history h

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

DATA STRUCTURE

- Sorted lists are unique for each history h

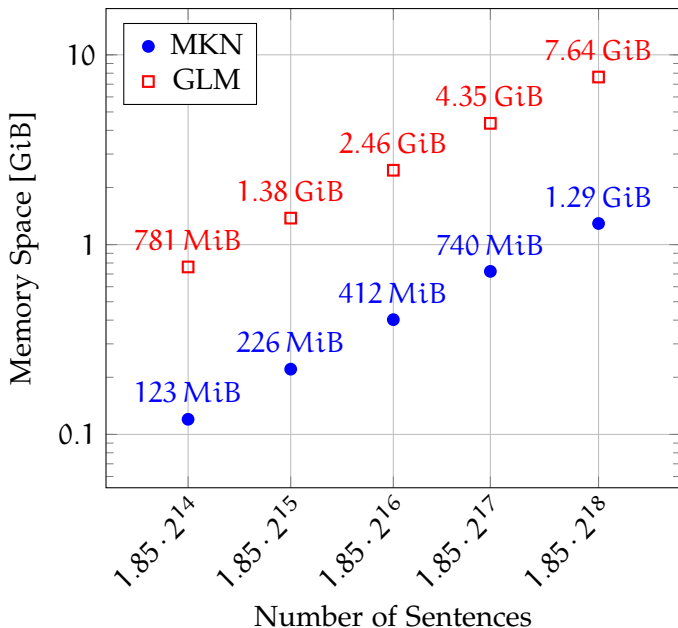
ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

- Optimized data structure: *Completion Trie* (Hsu and Ottaviano 2013)
 - Optimized for prefix-retrieval
 - High data compression

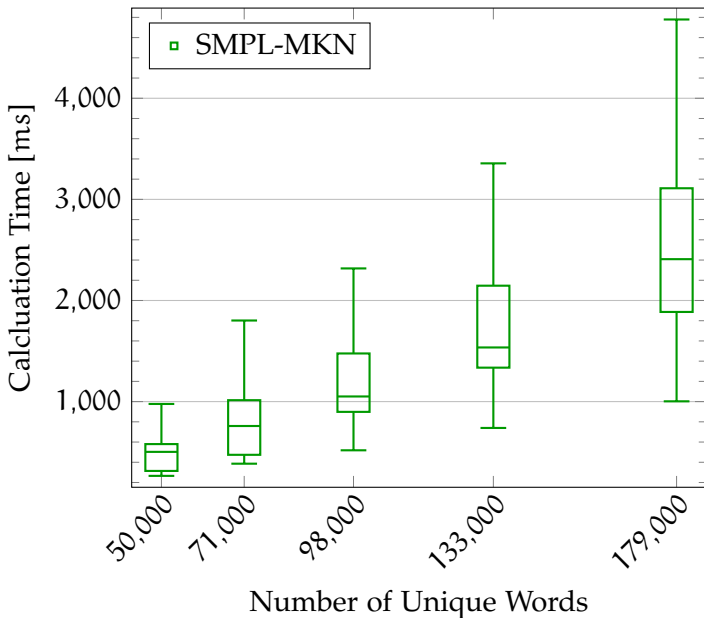
PREDICTION QUALITY

n-gram length	Keystroke Savings	
	MKN	GLM
2	0.44	0.44
3	0.50	0.50
4	0.51	0.51
5	0.51	0.52

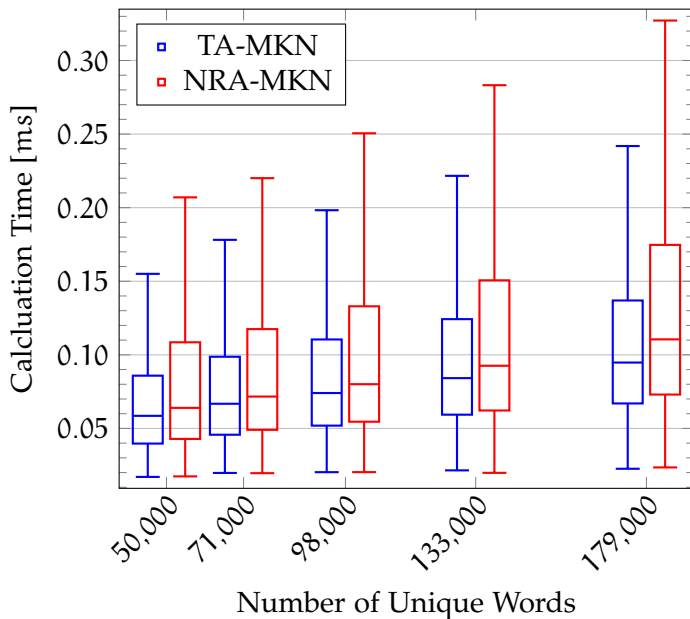
Completion Trie Size for 5-gram Prediction



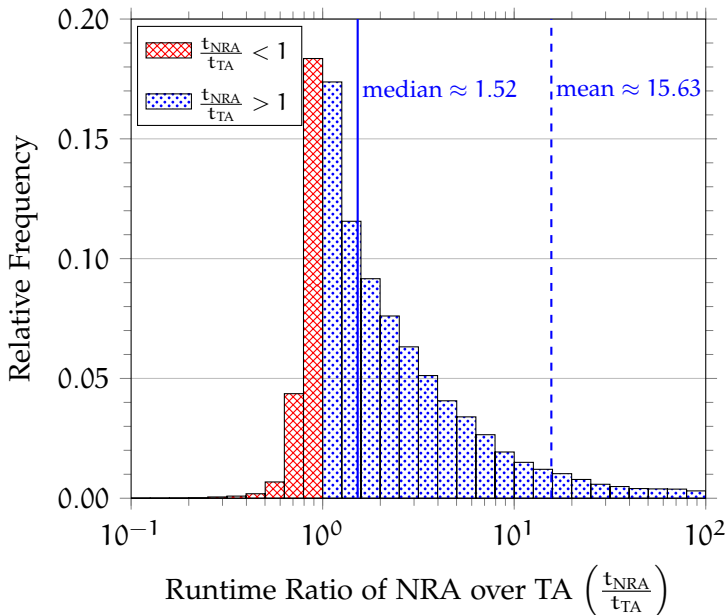
Calculation Time per top-1 Prediction using 5-grams



Calculation Time per top-1 Prediction using 5-grams



Comparison of NWP with TA and NRA over varying n-grams lengths, top-k predictions and corpus sizes



Section 4

CONCLUSION

CONCLUSION

- Weighted Sum representation for probabilities
 - Calculation equally fast for MKN, double as fast for GLM
 - Speeds up word prediction by 50% (MKN) to 70% (GLM)
- Top-k Joining
 - Speeds up word prediction by multiple orders of magnitude
 - Threshold Algorithm preferable to No Random Access Algorithm

THANKS FOR YOUR ATTENTION!

<http://141.26.208.6/>

Section 5

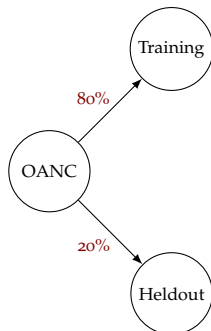
EXPERIMENTAL SETUP

EVALUATION CORPUS

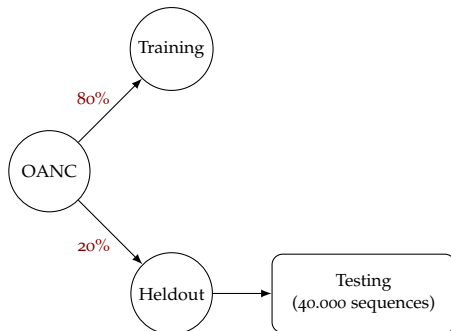
Open American National Corpus (Ide and Suderman 2007)

- open collection of American English
- historical and contemporary
- written text of all genres
- around 600.000 sentences / 14.000.000 words

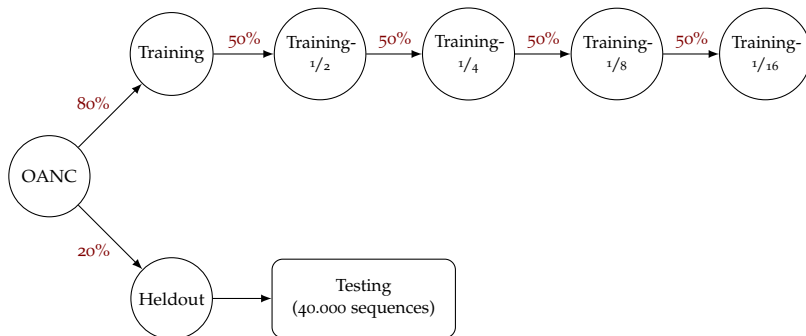
EXPERIMENTAL SETUP



EXPERIMENTAL SETUP



EXPERIMENTAL SETUP



Section 6

TOP-K JOIN EXAMPLES

TOP-K JOINS COMBINED WITH WEIGHTED SUMS

TOP-K JOINS COMBINED WITH WEIGHTED SUMS

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

TOP-K JOINS COMBINED WITH WEIGHTED SUMS

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

- Weights λ_i^h can be precomputed

TOP-K JOINS COMBINED WITH WEIGHTED SUMS

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

- Weights λ_i^h can be precomputed
- Sorted lists can store terms $\alpha_i^h(w)$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

TOP-K JOINS COMBINED WITH WEIGHTED SUMS

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

- Weights λ_i^h can be precomputed
- Sorted lists can store terms $\alpha_i^h(w)$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

- In examples we will assume $\lambda_1^h, \dots, \lambda_N^h = 1$
(not possible in practice!)

TOP-K JOINS COMBINED WITH WEIGHTED SUMS

$$P(w|h) = \sum_{i=1}^N \lambda_i^h \alpha_i^h(w)$$

- Weights λ_i^h can be precomputed
- Sorted lists can store terms $\alpha_i^h(w)$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

- In examples we will assume $\lambda_1^h, \dots, \lambda_N^h = 1$
(not possible in practice!)
- $P(\text{you} | \text{I love}) = 35 + 30 + 25 = 95$

THRESHOLD ALGORITHM

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w \mid \text{I love})$$

THRESHOLD ALGORITHM

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w \mid \text{I love})$$

ε	love	I love

THRESHOLD ALGORITHM

1. *Sorted Access* to lists
in any order
(e.g. round robin)

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w \mid \text{I love})$$

ε	love	I love

THRESHOLD ALGORITHM

1. *Sorted Access* to lists
in any order
(e.g. round robin)
2. For new words,
Random Access to
all other lists

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w \mid \text{I love})$$

ϵ	love	I love

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w \mid \text{I love})$$

ϵ	love	I love

- FAST WORD PREDICTION USING TOP-K JOINS

ϵ	love	I love

FAST WORD PREDICTION USING TOP-K JOINS

- LUKAS SCHMELZEISEN

$$P_{\text{threshold}} = \infty + \infty + \infty = \infty$$

ϵ	love	I love

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50		

$$P_{\text{threshold}} = \infty + \infty + \infty = \infty$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50		
	the 20	
		the 3

$$P_{\text{threshold}} = \infty + \infty + \infty = \infty$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50		
	the 20	
		the 3

$$P_{\text{threshold}} = \infty + \infty + \infty = \infty$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. *Compute threshold of highest possible unseen probability*
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50	the 20	the 3

$$P_{\text{threshold}} = 50 + \infty + \infty = \infty$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50	you 30	
	the 20	
		the 3

$$P_{\text{threshold}} = 50 + \infty + \infty = \infty$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50	you 30	you 25
	the 20	
you 35		the 3

$$P_{\text{threshold}} = 50 + \infty + \infty = \infty$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50	you 30	you 25
	the 20	
you 35		the 3

$$P_{\text{threshold}} = 50 + \infty + \infty = \infty$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. *Compute threshold of highest possible unseen probability*
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50	you 30	you 25
	the 20	
you 35		the 3

$$P_{\text{threshold}} = 50 + 30 + \infty = \infty$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50	you 30	you 25
	the 20	
you 35		the 3

$$P_{\text{threshold}} = 50 + 30 + \infty = \infty$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. *Compute threshold of highest possible unseen probability*
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ	love	I love
the 50	you 30	you 25
	the 20	
you 35		the 3

$$P_{\text{threshold}} = 50 + 30 + 25 = 105$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ		love	I love
the	50	you 30	you 25
a	40	the 20	
you	35		the 3

$$P_{\text{threshold}} = 50 + 30 + 25 = 105$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

$$P_{\text{threshold}} = 50 + 30 + 25 = 105$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

$$P_{\text{threshold}} = 50 + 30 + 25 = 105$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

$$P(\text{a} | \text{I love}) = 40 + 10 + 5 = 55$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. *Compute threshold of highest possible unseen probability*
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

$$P_{\text{threshold}} = 40 + 30 + 25 = 95$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

$$P(a | \text{I love}) = 40 + 10 + 5 = 55$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

$$P_{\text{threshold}} = 40 + 30 + 25 = 95$$

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

$$P(a | \text{I love}) = 40 + 10 + 5 = 55$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. *Compute threshold of highest possible unseen probability*
4. Terminate when k probabilities greater threshold have been found

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P_{\text{threshold}} = 40 + 20 + 25 = 85$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

$$P(a | \text{I love}) = 40 + 10 + 5 = 55$$

THRESHOLD ALGORITHM

1. *Sorted Access* to lists in any order (e.g. round robin)
2. For new words, *Random Access* to all other lists
3. Compute *threshold* of highest possible unseen probability
4. **Terminate when k probabilities greater threshold have been found**

$$\arg \max_{w \in \Sigma} P_{\text{MKN}}(w | \text{I love})$$

ϵ		love		I love	
the	50	you	30	you	25
a	40	the	20	a	5
you	35	a	10	the	3

$$P(\text{you} | \text{I love}) = 35 + 30 + 25 = 90$$

$$P_{\text{threshold}} = 40 + 20 + 25 = 85$$

$$P(\text{the} | \text{I love}) = 50 + 20 + 3 = 73$$

$$P(a | \text{I love}) = 40 + 10 + 5 = 55$$

NO RANDOM ACCESS ALGORITHM

NO RANDOM ACCESS ALGORITHM

1. *Sorted Access* to lists in any order

NO RANDOM ACCESS ALGORITHM

1. *Sorted Access* to lists in any order
2. Keep track of all seen counts

NO RANDOM ACCESS ALGORITHM

1. *Sorted Access* to lists in any order
2. Keep track of all seen counts
3. Compute *upper* and *lower bounds* for each probability

NO RANDOM ACCESS ALGORITHM

1. *Sorted Access* to lists in any order
2. Keep track of all seen counts
3. Compute *upper* and *lower bounds* for each probability
4. Terminate when k lower bounds greater than all other upper bounds have been found

Section 7

REFERENCES

REFERENCES I



Chen, Stanley F. and Goodman, Joshua T. (1999). “An empirical study of smoothing techniques for language modeling.” In: *Computer Speech & Language* 13.4, pp. 359–393 (cit. on pp. 12 sqq.).



Fagin, Ronald, Lotem, Amnon, and Naor, Moni (2001). “Optimal Aggregation Algorithms for Middleware.” In: *Proceedings of the Twentieth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*. PODS '01. Santa Barbara, California, USA: ACM, pp. 102–113 (cit. on pp. 35–38).

REFERENCES II



Hsu, Bo-June (Paul) and Ottaviano, Giuseppe (2013).

“Space-efficient Data Structures for Top-k Completion.” In:
Proceedings of the 22Nd International Conference on World Wide Web. WWW '13. Rio de Janeiro, Brazil: International World Wide Web Conferences Steering Committee, pp. 583–594
(cit. on pp. 40 sqq.).



Ide, Nancy and Suderman, Keith (2007). *The Open American National Corpus* (OANC). URL:

<http://www.anc.org/data/oanc/> (visited on July 5, 2015)
(cit. on p. 52).

REFERENCES III



Pickhardt, René, Gottron, Thomas, Körner, Martin, Wagner, Paul Georg, Speicher, Till, and Staab, Steffen (2014). “A Generalized Language Model as the Combination of Skipped n-grams and Modified Kneser Ney Smoothing.” In: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vol. abs/1404.3377. Baltimore, Maryland: Association for Computational Linguistics, pp. 1145–1154 (cit. on pp. 12 sqq.).