

Bayesian Language Modelling of German Compounds

Jan Botha* **Chris Dyer[†]** **Phil Blunsom***

*Department of Computer Science
University of Oxford

[†]School of Computer Science
Carnegie Mellon University

COLING 2012

Why Statistical Language Modelling?

- Central to tasks where sentences are hypothesised
 - Machine Translation
 - Speech Recognition
 - Optical Character Recognition

Token-based n-gram models

Task: $P(\text{Wir fahren mit der Bahn}) = ?$

We're going by train

Token-based n-gram models

Task: $P(\text{Wir fahren mit der Bahn}) = ?$

We're going by train

Markov assumption

Makes parameter estimation feasible

e.g. $P(\text{Bahn} \mid \text{mit der}) = 0.3$

Token-based n-gram models

Task: $P(\text{Wir fahren mit der Bahn}) = ?$ *We're going by train*

Markov assumption

Makes parameter estimation feasible

$$\text{e.g. } P(\text{Bahn} \mid \text{mit der}) = 0.3$$

Vocabulary Assumptions

Closed vocabulary, independent words

$$P(\begin{array}{c} \text{Bahn} \\ \text{Familie} \end{array} \mid \text{mit der}) = ?$$

Token-based n-gram models

Task: $P(\text{Wir fahren mit der Bahn}) = ?$ *We're going by train*

Markov assumption

Makes parameter estimation feasible

$$\text{e.g. } P(\text{Bahn} \mid \text{mit der}) = 0.3$$

Vocabulary Assumptions

Closed vocabulary, independent words

$$P(\begin{array}{c} \text{Bahn} \\ \text{Familie} \\ \text{Mauer} \end{array} \mid \text{mit der}) = ?$$

Token-based n-gram models

Task: $P(\text{Wir fahren mit der Bahn}) = ?$ *We're going by train*

Markov assumption

Makes parameter estimation feasible

$$\text{e.g. } P(\text{Bahn} \mid \text{mit der}) = 0.3$$

Vocabulary Assumptions

Closed vocabulary, independent words

$$P(\begin{array}{l} \text{Bahn} \\ \text{Familie} \\ \text{Mauer} \\ \text{Straßenbahn} \\ \text{U-Bahn} \\ \text{Bobbahn} \end{array} \mid \text{mit der}) = ?$$

Compounding

Example

Wir fahren mit der	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
	U-bahn .	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Compounding

Example

Wir fahren mit der	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
	U-bahn	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Highly Productive Process

Regal

shelf

Compounding

Example

	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
Wir fahren mit der	U-bahn	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Highly Productive Process

Regal	<i>shelf</i>
Buchregal	<i>bookshelf</i>

Compounding

Example

	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
Wir fahren mit der	U-bahn	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Highly Productive Process

Regal	<i>shelf</i>
Buchregal	<i>bookshelf</i>
Buchregalhersteller	<i>bookshelf manufacturer</i>

Compounding

Example

	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
Wir fahren mit der	U-bahn	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Highly Productive Process

Regal	<i>shelf</i>
Buchregal	<i>bookshelf</i>
Buchregalhersteller	<i>bookshelf manufacturer</i>
Buchregalherstellername	<i>bookshelf manufacturer's name</i>
...	

Compounding

Example

Wir fahren mit der	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
	U-bahn	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Regularity

Compounding

Example

	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
Wir fahren mit der	U-bahn .	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Regularity

- compound = modifiers & head

Compounding

Example

	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
Wir fahren mit der	U-bahn .	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Regularity

- compound = modifiers & head
- head determines syntactic properties

Compounding

Example

	Bahn	<i>train</i>
	Straßenbahn	<i>tram</i>
Wir fahren mit der	U-bahn .	<i>metro</i>
	Bobbahn	<i>bobsled</i>
	Achterbahn	<i>rollercoaster</i>

Regularity

- compound = modifiers & head
- head determines syntactic properties
- generalising over modifiers would reduce sparsity

Related Work: Factored Language Model

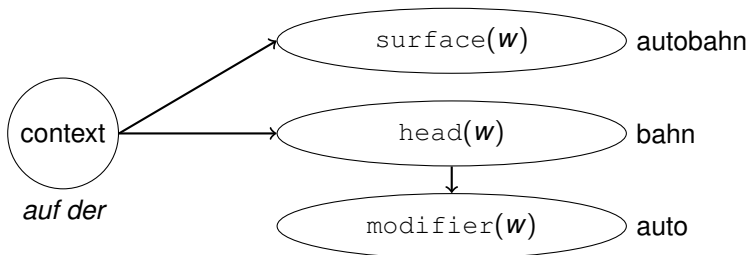
Each token w consists of k features, fixed dependencies

(Bilmes & Kirchhoff, 2003)

Related Work: Factored Language Model

Each token w consists of k features, fixed dependencies

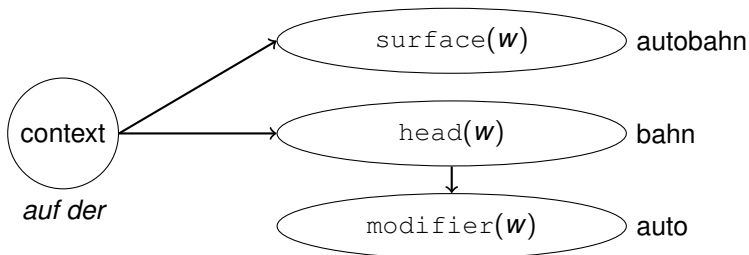
(Bilmes & Kirchhoff, 2003)



Related Work: Factored Language Model

Each token w consists of k features, fixed dependencies

(Bilmes & Kirchhoff, 2003)



Shortcoming: Compound lengths vary freely

Auto	+	bahn	+	geschwindigkeit	+	s	+	grenze
<i>motor</i>		<i>way</i>		<i>speed</i>				<i>limit</i>

Related Work: Splitting and Merging

Preprocess data: split compounds

(e.g. Koehn & Knight, 2003; Stymne 2008; Macherey et al. 2011)

Related Work: Splitting and Merging

Preprocess data: split compounds

(e.g. Koehn & Knight, 2003; Stymne 2008; Macherey et al. 2011)

Example

		Context	Predicted Token
<i>Standard:</i>	Wir fahren	auf der	Autobahn
<i>Split:</i>	Wir fahren auf	der Auto	bahn

Related Work: Splitting and Merging

Preprocess data: split compounds

(e.g. Koehn & Knight, 2003; Stymne 2008; Macherey et al. 2011)

Example

		Context	Predicted Token
<i>Standard:</i>	Wir fahren	auf der	Autobahn
<i>Split:</i>	Wir fahren auf	der Auto	bahn

Shortcoming: Important sentential context lost

Model Aims

- Condition head on previous words
- Condition compound-internal structure on head

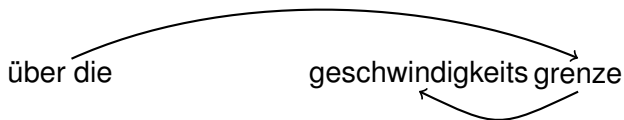
Model Aims

- Condition head on previous words
- Condition compound-internal structure on head



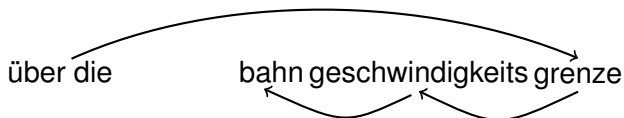
Model Aims

- Condition head on previous words
- Condition compound-internal structure on head



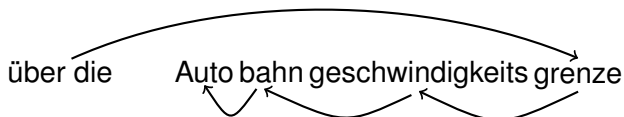
Model Aims

- Condition head on previous words
- Condition compound-internal structure on head



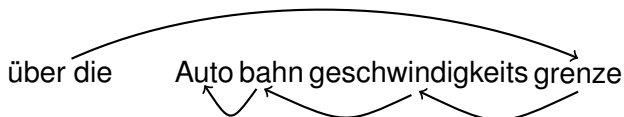
Model Aims

- Condition head on previous words
- Condition compound-internal structure on head



Model Aims

- Condition head on previous words
- Condition compound-internal structure on head



- *but* keep full surface form in model

Hierarchical Pitman-Yor Language Model (HPYLM)

Standard back-off smoothing

Hierarchical Pitman-Yor Language Model (HPYLM)

Standard back-off smoothing

mit der Autobahn



$P(\text{Autobahn} \mid \text{mit der})$

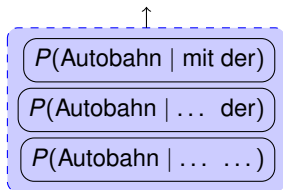
$P(\text{Autobahn} \mid \dots \text{der})$

$P(\text{Autobahn} \mid \dots \dots)$

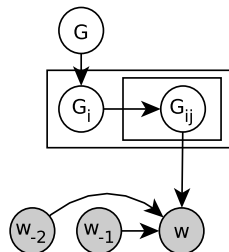
Hierarchical Pitman-Yor Language Model (HPYLM)

Standard back-off smoothing

mit der Autobahn

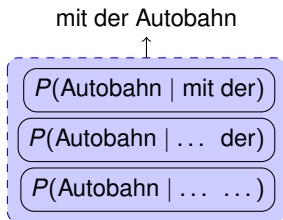


Hierarchical Prior

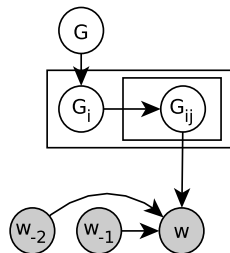


Hierarchical Pitman-Yor Language Model (HPYLM)

Standard back-off smoothing



Hierarchical Prior



Conditional distributions drawn from Pitman-Yor process:

$$G \sim PY(d, \theta, G_0)$$

discount

strength

~~base distribution~~

Our Model: HPYLM with Compounds

Proposed back-off

mit der Autobahn

$P(\text{Autobahn} \mid \text{mit der})$

$P(\text{bahn} \mid \text{mit der})$

$P(\text{bahn} \mid \dots \text{der})$

$P(\text{bahn} \mid \dots \dots)$

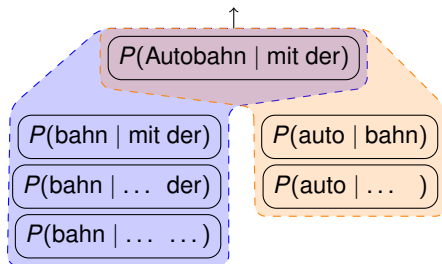
$P(\text{auto} \mid \text{bahn})$

$P(\text{auto} \mid \dots)$

Our Model: HPYLM with Compounds

Proposed back-off

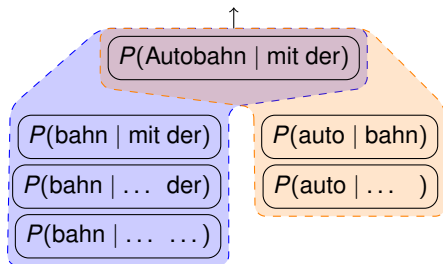
mit der Autobahn



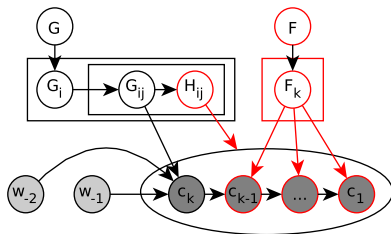
Our Model: HPYLM with Compounds

Proposed back-off

mit der Autobahn



Extended Hierarchical Prior



Overview of Setup

- Data from WMT '11 shared task
- Language Models
 - 4-gram, unless otherwise stated
 - trained on 59m tokens (Europarl, news & commentary)
- English→German Translation System
 - trained on 1.7m parallel sentences (Europarl)

Overview of Setup

- Data from WMT '11 shared task
- Language Models
 - 4-gram, unless otherwise stated
 - trained on 59m tokens (Europarl, news & commentary)
- English→German Translation System
 - trained on 1.7m parallel sentences (Europarl)
- Vocabulary sparsity – ratio of English/German:
 - baseline: 3.13
 - decompounded: 1.36

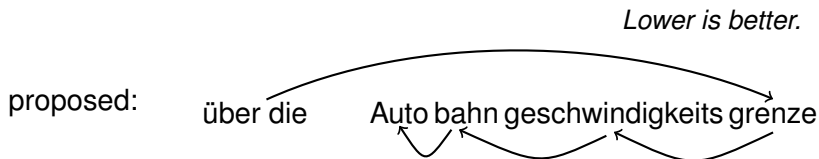
Comparison of Methods for Predicting Unseen Text

	Perplexity
Modified Kneser-Ney	299.9
HPYLM	294.1

Lower is better.

Comparison of Methods for Predicting Unseen Text

	Perplexity
Modified Kneser-Ney	299.9
HPYLM	294.1
Our model	293.6

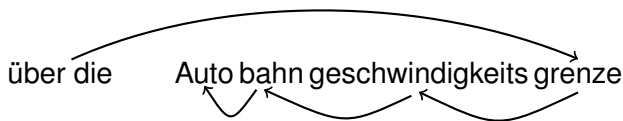


Comparison of Methods for Predicting Unseen Text

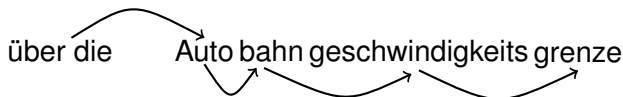
	Perplexity
Modified Kneser-Ney	299.9
HPYLM	294.1
Our model	293.6
Our model (inverted)	305.5

Lower is better.

proposed:

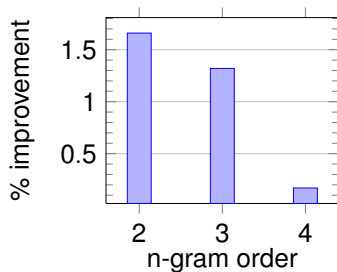


inverted:



Effect of Scaling Context Size

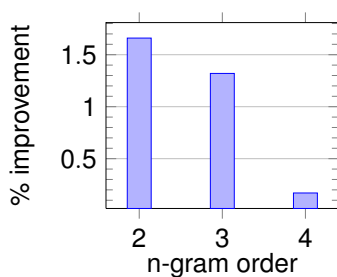
% Relative improvement over HPYLM baseline



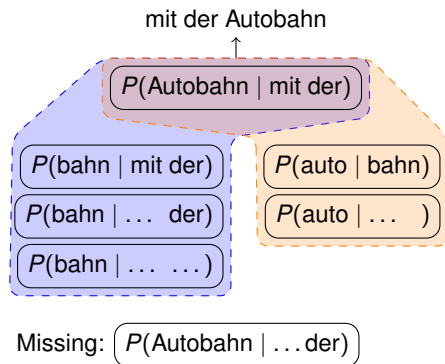
Higher is better.

Effect of Scaling Context Size

% Relative improvement over HPYLM baseline



Higher is better.



Machine Translation Results

English→German

	BLEU
Modified Kneser-Ney	13.9
HPYLM	13.9
Our model	13.9
Our model (inverted)	13.7

Higher BLEU is better.

Machine Translation Results

English→German

only 3.7% of reference tokens are compounds

	BLEU
Modified Kneser-Ney	13.9
HPYLM	13.9
Our model	13.9
Our model (inverted)	13.7

Higher BLEU is better.

Compound Translation Accuracy

Compare compounds in output against reference sentences

	Precision	Recall	F-score
Modified Kneser-Ney	25.4	17.1	20.5
HPYLM	24.3	17.5	20.4
Our model	27.5	17.3	21.3
Our model (inverted)	23.7	17.2	19.9

Higher is better.

Summary

- Productive compounding is an important source of sparsity.
- Proposed n-gram language model that accounts for productive compounding.
- Improved accuracy of translated compounds, while matching baseline BLEU score

Summary

- Productive compounding is an important source of sparsity.
- Proposed n-gram language model that accounts for productive compounding.
- Improved accuracy of translated compounds, while matching baseline BLEU score
- Future Work
 - Evaluate with an MT system that outputs *novel* compounds
 - Model compounds that occur within context

Thank you.

	BLEU
Modified Kneser-Ney	13.9
HPYLM	13.9
Our model	13.9

	LM-training	BLEU
Modified Kneser-Ney	< 1 hour	13.9
HPYLM	3 days	13.9
Our model	6 days	13.9

	LM-training	BLEU
Modified Kneser-Ney	< 1 hour	13.9
HPYLM	3 days	13.9
Our model	6 days	13.9
Our model (fastapprox)	4 hours	13.6