

Machine Translation Panel

Nadir Durrani, University of Edinburgh: Operation Sequence Model

Chris Dyer, CMU: Word Classes

Spence Green, Stanford: Sparse Feature Training

Kenneth Heafield, Stanford: Huge Language Models

Stephan Peitz, RWTH Aachen: Leave One Out Training

Philip Williams, University of Edinburgh: String-to-Tree Syntax

Nadir Durrani
University of Edinburgh

Operation Sequence Model



Introduction

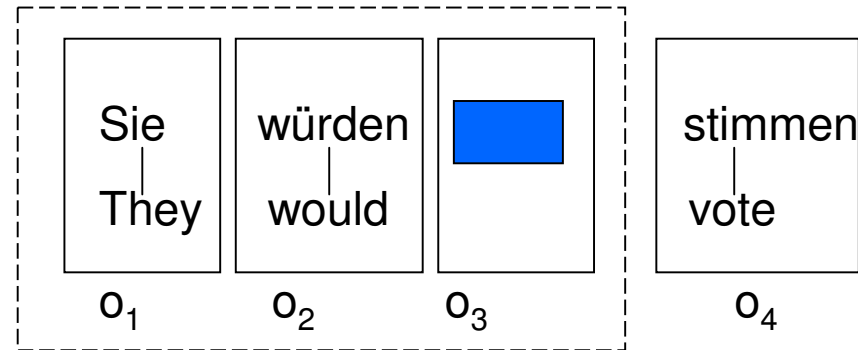
- A model that
 - combines benefits from Phrase-based and N-gram-based SMT
 - is based on minimal translation but memorizes like phrases
 - considers source and target contextual information across phrases
 - integrates translation and reordering into a single model
- Convert a bilingual sentence to a sequence of operations
 - Translate (Generate a minimal translation unit)
 - Reordering (Insert a gap or Jump)
- $P(e,f,a)$ = N-gram model over resulting operation sequences

Example

Sie würden gegen Sie stimmen
 They would vote against you

Operations

- o_1 Generate (Sie, They)
- o_2 Generate (würden, would)
- o_3 Insert Gap
- o_4 Generate (stimmen, vote)
- o_5 Jump Back (1)
- o_6 Generate (gegen, against)
- o_7 Generate (Sie, you)



Context Window

Model:

$$p_{\text{osm}}(F, E, A) = p(o_1, \dots, o_N) = \prod_i p(o_i | o_{i-n+1} \dots o_{i-1})$$



How does it improve Phrase-based SMT?

- Overcomes phrasal independence assumption
 - Considers source and target contextual information across phrases
- Better reordering model
 - Translation and reordering decisions influence each
 - Handles local and long distance reorderings in a unified manner
- No spurious phrasal segmentation problem
- Average gain of +0.40 on news-test2013 across 10 pairs

Thank You !!!

Chris Dyer
CMU

Word Classes

Using Word Clusters



1. Cluster monolingual data
2. 500-1000 clusters
3. Use for: LMs, features

percyliang / brown-cluster

C++ implementation of the Brown word clustering algorithm.

18 commits 1 branch

branch: master brown-cluster / +

Common prefixes,
Common context

^1101100 (126689)	Merkel Obama Müller Schmidt Friedrich Steinbrück Fischer Koch Wulff Westerwelle Schäuble Schneider Löw Sarkozy Seehofer Schröder Vettel S. Putin Beck Berlusconi B. Wagner Gabriel Rösler Wolf Hoeneß Becker Weber Steinmeier Meyer Bush Clinton Bauer Schulz Schumacher Jung Romney Schäfer Klein M. Kaiser Gomez Roth K. W. Ali Hollande Gauck Heynckes
^1101101000 (2965)	Deutschland Madrid Leverkusen BSC Motors Stanley Brothers Woods bar Christus Ostdeutschland
^11011010010 (5147)	China Österreich Frankreich Russland Italien Spanien Israel Japan Großbritannien Polen Ägypten Indien Brasilien Schweden England Niedersachsen Brandenburg Portugal Irland Pakistan Australien Ungarn Kanada Belgien Nordkorea Südafrika Dänemark Thüringen Tschechien Rumänien Teheran Serbien Norwegen Südkorea Silber Kroatien Finnland Tunesien Argentinien Island Gladbach Bulgarien Bronze Osteuropa Hongkong Singapur Georgien Katar Thailand Holland
^11011010011 (4071)	Europa Griechenland Syrien Afghanistan Hessen NRW Amerika Libyen Afrika Asien Zypern Fukushima Rio Boston Kuba Gaza Mallorca Mali Wimbledon Bosnien Tibet Haiti Gorleben Sylt Palästina Hawaii Auschwitz Tschernobyl Sotschi Westeuropa Schach Lampedusa Jugoslawien Oberbayern Santiago Philadelphia Oberfranken Wembley Übersee Guantánamo Teilzeit Fort Guantanamo Ostwestfalen Tschetschenien Aufruhr Echtzeit Südeuropa Sierra Steueroasen
^1101101010 (6344)	München Frankfurt Stuttgart Köln Dortmund Hannover Bremen Düsseldorf Nürnberg Wolfsburg Augsburg Mainz Freiburg Bochum 04 Hoffenheim Mönchengladbach United Paderborn Cottbus Magdeburg Zürich Fürth Darmstadt Basel Braunschweig Lübeck Bamberg Ingolstadt Osnabrück Valencia Ried Friedrichshafen Herford Minden Koblenz Gütersloh Aalen Aue Florenz Mavericks Rottweil Schweinfurt Borken Donaueschingen Freudenstadt Calw Ravensburg Saarbrücken Reutlingen
^1101101011 (22412)	Berlin Hamburg Wien London Paris Washington Brüssel Moskau Sachsen Rom Peking Salzburg Athen Münster Leipzig Dresden Potsdam Kiel Bonn Linz Wiesbaden Bielefeld Kairo Mexiko Aachen Karlsruhe Duisburg Istanbul Mailand Tokio Damaskus Graz Luxemburg Regensburg Mannheim Kassel Rostock Flensburg Amsterdam Kalifornien Innsbruck Tripolis Würzburg Florida Jerusalem Offenbach Chicago Oberösterreich Erfurt Ulm

Using Word Clusters



Rule “shape” features (prefix length = 6)

$X \rightarrow (\textit{dass X angekommen ist, that X arrived})$

$C1001_X_C001101_C110100::C010111_X_C111111=1$

7-gram class-based LM

Absolute discounting ($d=0.5$)

Separate features for transitions and emissions

	BLEU	MET	TER
Baseline	25.3	30.4	52.6
+Rule shape	25.5	30.5	52.4
+7gm LM	26.4	31.0	51.9

Spence Green
Stanford

Sparse Feature Training

Large-scale Discriminative Tuning

#1: 2010s ML in MT tuning

Online convex optimization

Arbitrary, overlapping features

#2: Large tuning sets

Fast decoding and updating

Bitext tuning...

See our poster and talk for details

WMT14 Shared Task Results

Uncased BLEU results

	dense-dev	features-dev	2014 rank
Fr-En	19.6	20.0	1
En-De	32.0	32.5	1

Tune: 13.5k sentences (2008-2012)

Models have 200-300k features

Kenneth Heafield
Stanford

Huge Language Models

Impact of Big Language Models

Target	Base Rank		+LM Rank	ΔBLEU
Czech	5–6	→	1–3	+0.6
Hindi	4–5	→	3	+1.4
Russian	6–7	→	4–5	+1.2
German	8–10	→	3–6	+0.5

After the evaluation: Hindi–English +0.9 BLEU

Download multiple LMs and training data from
`statmt.org/ngrams`

English: 1.8 trillion tokens

n Unique *n*-grams

1 2,640,258,088

2 15,297,753,348

3 61,858,786,129

4 156,775,272,110

5 263,690,452,834

Current work: approximate LM storage.

Stephan Peitz
RWTH Aachen

Leave One Out Training

Consistent Phrase Training

State of the art

- ▶ Heuristic extraction of phrases using word alignments
- ▶ Compute translation probabilities as relative frequencies

Issues of this heuristic

- ▶ Extract from likely alignment?
- ▶ Models used in decoding are not considered \Rightarrow **inconsistency**

Forced decoding

- ▶ Run decoder on training data
- ▶ Count used phrases, recompute probabilities
- ▶ Apply **leave-one-out** to counteract overfitting

Leave-One-Out

- ▶ Occurrences of a phrase in a sentence pair (f_n, e_n) are subtracted from the phrase counts obtained from the full training data

$$p_{l_{1o,n}}(\tilde{f}|\tilde{e}) = \frac{C(\tilde{f}, \tilde{e}) - C_n(\tilde{f}, \tilde{e})}{\sum_{\tilde{f}'} C(\tilde{f}', \tilde{e}) - C_n(\tilde{f}', \tilde{e})}$$

- ▶ Singleton phrases get a low probability

Consistent Phrase Training using Leave-One-Out

▶ Publications

- ▶ Phrase-based [Wuebker & Mauser⁺ 10, Wuebker & Ney 13]
- ▶ Hierarchical [Peitz & Mauser⁺ 12, Peitz & Vilar⁺ 14]

▶ Improvements: 0.5-1.5 BLEU

▶ Reducing phrase-table size to 5-20% of the original size

▶ Systems using phrase/rule training:

- ▶ WMT 2011 (RWTH, German→English, phrase-based)
- ▶ IWSLT 2011 (RWTH, German→English, phrase-based)
- ▶ IWSLT 2012 (RWTH, German→English, hierarchical)
- ▶ BOLT 2012 (RWTH, Chinese→English, hierarchical)
- ▶ OpenMT 2012 (NRC, Chinese→English, phrase-based)

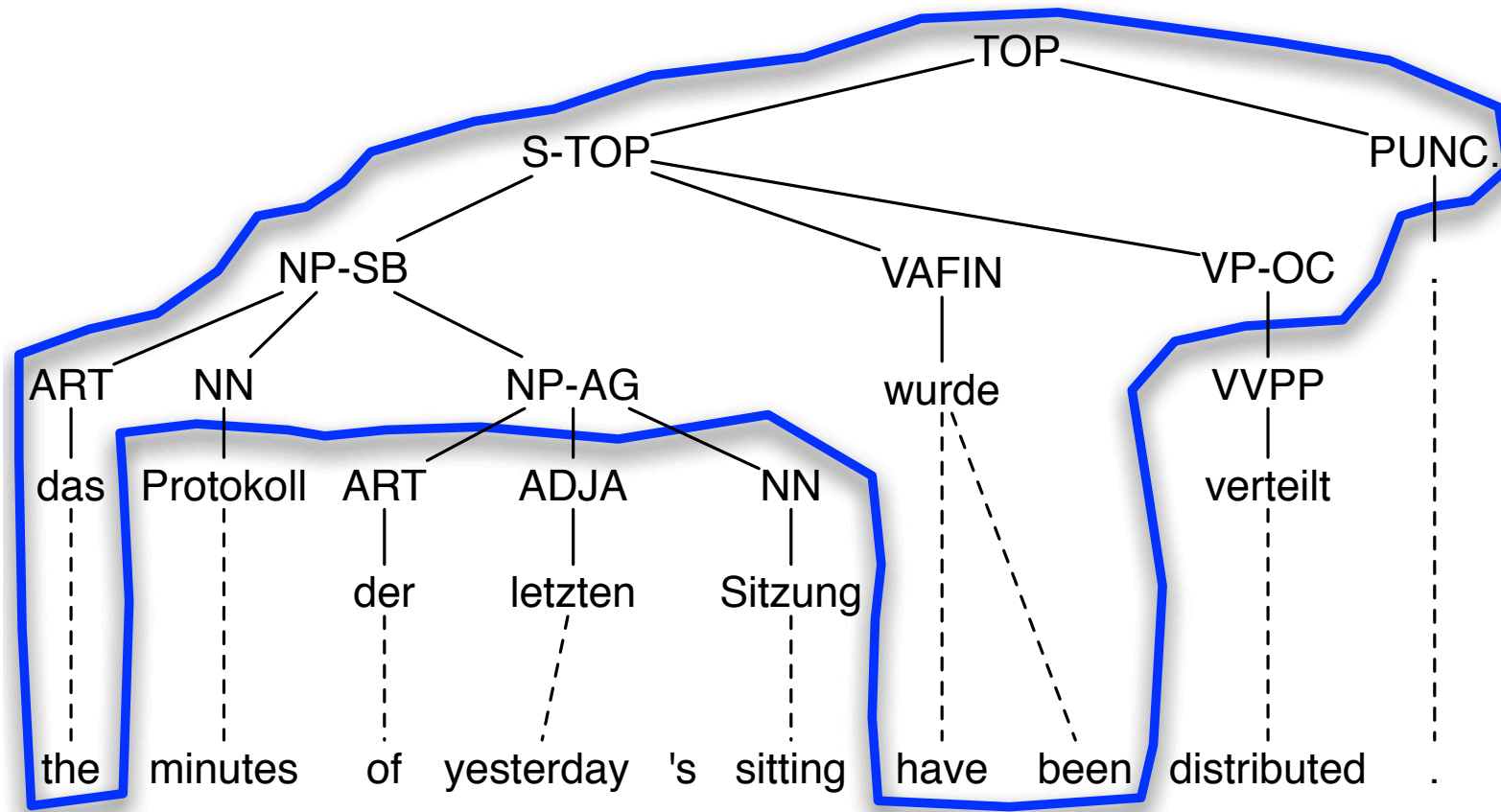
▶ Implemented in RWTH's translation toolkit **Jane**

<http://www.hltpr.rwth-aachen.de/jane>

Philip Williams
University of Edinburgh

String-to-Tree Syntax

uedin-syntax: string-to-tree



TOP → the X_1 X_2 have been X_3 X_4 | das NN_1 $NP-AG_2$ wurde $VP-OC_3$ $PUNC.4$



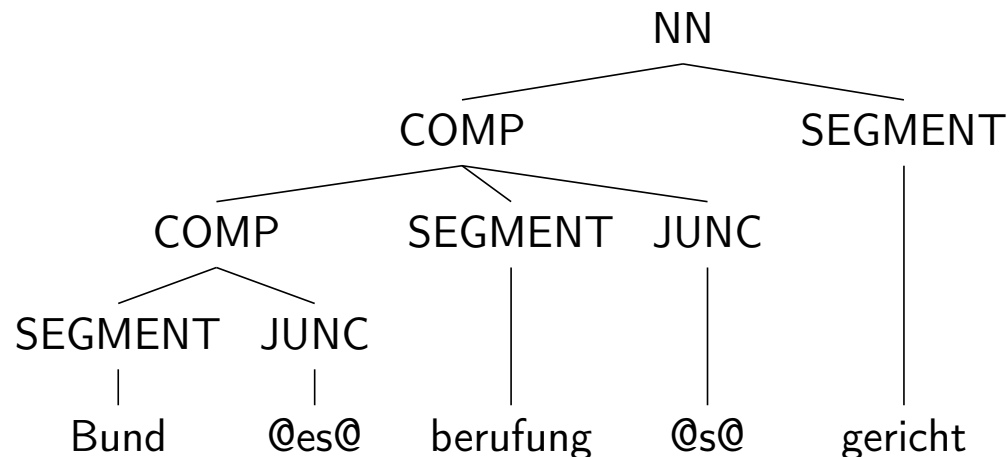
uedin-syntax: string-to-tree extensions

Use syntactic structure to help model other aspects of target-side grammar.

Example 1. Agreement

TOP \rightarrow the X_1 X_2 have been X_3 X_4 | das NN_1 NP-AG $_2$ wurde VP-OC $_3$ PUNC. $_4$
 $\langle NN_1$ AGR $\rangle = \langle$ wurde AGR \rangle

Example 2. Compound Splitting



Best constrained system for:
English-German
German-English
Hindi-English (tied with CMU)



References

- Durrani, N., Schmid, H., and Fraser, A. (2011). A Joint Sequence Translation Model with Integrated Reordering. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1045–1054, Portland, Oregon, USA.
- Durrani, N., Fraser, A., and Schmid, H. (2013a). Model With Minimal Translation Units, But Decode With Phrases. In *The 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Durrani, N., Fraser, A., Schmid, H., Hoang, H., and Koehn, P. (2013b). Can Markov Models Over Minimal Translation Units Help Phrase-Based SMT? In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, Sofia, Bulgaria. Association for Computational Linguistics.
- Durrani, N., Koehn, P., Schmid, H., and Fraser, A. (2014a). Investigating the Usefulness of Generalized Word Representations in SMT. In *Proceedings of the 25th Annual Conference on Computational Linguistics (COLING)*, Dublin, Ireland.

- Buck, C., Heafield, K., and van Ooyen, B. (2014). N-gram Counts and Language Models from the Common Crawl. In *Proceedings of the Language Resources and Evaluation Conference*, Reykjavik, Iceland.
- S. Peitz, A. Mauser, J. Wuebker, H. Ney: Forced Derivations for Hierarchical Machine Translation. In *International Conference on Computational Linguistics*, pp. 933–942, Mumbai, India, Dec. 2012.
- S. Peitz, D. Vilar, H. Ney: Simple and Effective Approach for Consistent Training of Hierarchical Phrase-based Translation Models. In *Conference of the European Chapter of the Association for Computational Linguistics*, Gothenburg, Sweden, April 2014.
- J. Wuebker, A. Mauser, H. Ney: Training Phrase Translation Models with Leaving-One-Out. In *Annual Meeting of the Assoc. for Computational Linguistics*, pp. 475–484, Uppsala, Sweden, July 2010.
- J. Wuebker, H. Ney: Length-incremental Phrase Training for SMT. In *ACL 2013 Eighth Workshop on Statistical Machine Translation*, pp. 309–319, Sofia, Bulgaria, Aug. 2013.
- Victor Chahuneau, Eva Schlinger, Noah A. Smith, and Chris Dyer. 2013. Translating into Morphologically Rich Languages with Synthetic Phrases. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*.

Heafield, K., Pouzyrevsky, I., Clark, J. H., and Koehn, P. (2013). Scalable Modified Kneser-Ney Language Model Estimation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, Sofia, Bulgaria.

Spence Green, Sida Wang, Daniel Cer, and Christopher D. Manning. 2013. Fast and adaptive online training of feature-rich translation models. In *ACL*.

Spence Green, Daniel Cer, Kevin Reschke, Rob Voigt, John Bauer, Sida Wang, et al. 2013. Feature-rich phrase-based translation: Stanford University's submission to the WMT 2013 translation task. In *WMT*.

Philip Williams and Philipp Koehn. 2012. GHKM Rule Extraction and Scope-3 Parsing in Moses. In *Proceedings of the Seventh Workshop on Statistical Machine Translation*, pages 388–394, Montréal, Canada, June.

Philip Williams and Philipp Koehn. 2011. Agreement Constraints for Statistical Machine Translation into German. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 217–226, Edinburgh, Scotland, July.