

**CLAIRE Makes Machine Translation
BLEU No More**

by

Ali Mohammad

B.Sc. Mathematics and Physics, Kansas State University (2003)

B.Sc. Computer Engineering, Kansas State University (2003)

B.Sc. Computer Science, Kansas State University (2003)

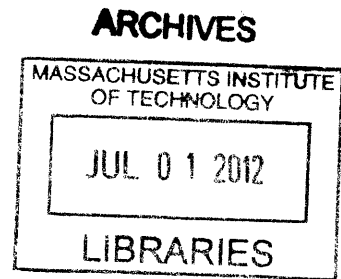
S.M., Massachusetts Institute of Technology (2006)

Submitted to the Department of
Electrical Engineering and Computer Science
in partial fulfillment of the requirements for the degree of
Doctorate of Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2012



© Massachusetts Institute of Technology 2012. All rights reserved.

Author
Department of Electrical Engineering and Computer Science
February 14, 2012

Certified by
Boris Katz
Principal Research Scientist in Computer Science
Thesis Supervisor

Accepted by
Leslie A. Kolodziejski
Chair of the Committee on Graduate Students
Department of Electrical Engineering and Computer Science

CLAIRE Makes Machine Translation

BLEU No More

by

Ali Mohammad

Submitted to the Department of Electrical Engineering and Computer Science
on February 14, 2012, in partial fulfillment of the
requirements for the degree of
Doctorate of Science

Abstract

We introduce CLAIRE, a mathematically principled model for inferring ranks and scores for arbitrary items based on forced-choice binary comparisons, and show how to apply this technique to statistical models to take advantage of problem-specific assistance from non-experts. We apply this technique to two language processing problems: parsing and machine translation. This leads to an analysis which casts doubts on modern evaluation methods for machine translation systems, and an application of CLAIRE as a new technique for evaluating machine translation systems which is inexpensive, has theoretical guarantees, and correlates strongly in practice with more expensive human judgments of system quality. Our analysis reverses several major tenants of the mainstream machine translation research agenda, suggesting in particular that the use of linguistic models should be reexamined.

Thesis Supervisor: Boris Katz

Title: Principal Research Scientist in Computer Science

Acknowledgments

MIT is an amazing place, and I probably would have chosen to never graduate if that were an option. So many people were involved in making it amazing that it would take another thesis to thank them all.

Boris Katz was a wonderful advisor. He gave consistently good advice, but allowed me to make the mistakes I was determined to make. I grew as a scientist under his tutelage, and I am incredibly grateful for the guidance he gave me.

I would like to thank my committee members Bob Berwick and Tommi Jaakkola, who bore with my meandering progress with patience and gave nuanced comments on several drafts of this thesis, and mentored me for the many years that I was at MIT.

Thanks to my master's thesis advisor Michael Collins and his research group for giving me some excellent early guidance and an introduction to language processing.

I owe a great deal to Leslie Pack Kaelbling, Tomas Lozano-Perez, Denny Freeman, Eric Grimson, and the rest of the 6.01 crew for providing a tranquil oasis right when I really needed it.

I am deeply indebted to Andrew Correa, Yuan Shen, and Gabriel Zaccak who provided feedback, encouragement, and helpful discussions, especially as I was wrapping up. This thanks for discussions extends to anyone that was ever my office-mate or sat anywhere near me.

I have to thank the CSAIL Olympics (which was the AI Olympics when I first came

to MIT!) for introducing me to so many other students and getting me into student activities right when I started and led me into the CSAIL Student Committee.

Thanks to the good folks at Serious Business and danger!awesome who have given me a much needed creative/artistic outlet, and to Ron Wiken and the folks at TIG for tolerating my presence at the machine shop and for preventing me from slowly demolishing the lab.

Thanks to my work-out buddies who sometimes got me out of the building and kept me from getting too frustrated at lab.

Thanks to all my new friends at Google that showed me that there are amazing places to be outside of MIT.

Most of all, I am grateful to my family which suffered burdens beyond description with fortitude and dignity. My parents and my sisters have been profoundly supportive and patient in spite of everything from start to finish, and my extended family has always been there to encourage me. You are my inspiration.

وَقُلْ رَبِّ زِدْنِي عِلْمًا

And say, O My Lord! Increase me in Knowledge.

Contents

Introduction	12
I.1 Prologue	12
I.2 Background	19
I.2.1 Parsing	19
I.2.2 Machine Translation	20
I.2.3 Five IBM Models	21
I.2.4 Phrase-Based Models	26
I.3 Answering Kay	28
I.4 Outline	29
1 Human Assisted Parsing	30
1.1 Previous Work	33
1.2 Our Approach	33
1.2.1 Collins Parser	34
1.2.2 Algorithm	35
1.2.3 Tree Transduction Grammar	36
1.2.4 Sample Output	41
1.3 Further Augmenting the Tree	42

1.3.1	Method	44
1.3.2	Data Collection	44
1.3.3	The Model	46
1.4	Results	50
1.5	Future Work	51
2	Human Assisted Word-Sense Disambiguation for Translation	54
2.1	Previous Work	55
2.2	Methods	57
2.3	Mechanical Turk	61
2.4	Results	61
3	Human Assisted Translation Evaluation	64
3.1	Answering Kay	64
3.1.1	Rule-Based MT is Bad	64
3.1.2	Hybrid Systems Are OK	65
3.1.3	Minimum Error-Rate Training	65
3.1.4	Linguistics Doesn't Work	65
3.1.5	<i>n</i> -gram Language Models Are Best	66
3.1.6	More <i>n</i> -grams Help... Forever	67
3.1.7	Shrinking the Translation Model is OK	70
3.1.8	Word-Sense Disambiguation Does Not Help Machine Translation	70
3.1.9	Having Monolingual Consultants Revise Translations is Bad .	70
3.2	BLEU	71
3.2.1	BLEU Correlates Well With Human Judgments?	73

3.2.2	Re-evaluating BLEU	73
3.3	Methods	75
3.3.1	Active Ranking by Forced Comparisons	79
3.3.2	Previous Work	81
3.3.3	Consequences of Choosing a Bad Activation Function	81
3.4	Mechanical Turk	84
3.5	Results	84
3.6	Conclusions	86
4	Conclusions and Future Work	87
A	EM Algorithm Reference	90
A.1	Definitions	90
A.2	Lemmas	90
A.3	EM is Nondecreasing	91
A.4	EM on Multinomials	92
A.5	Parametrizing Multinomials	93
A.6	EM on IBM2+1dG	94
B	Information on the WSJ Corpus Tagset	98
B.1	Terminal Tags (Parts of Speech)	98
B.2	Members of Closed-Class Tags	100
B.3	Nonterminals	102
B.4	Extended Tags	103
B.5	Pseudo-Attachment/Null Element Markers	104

B.6	Examples of the Full Extended Tags	105
C	Modifying the Bikel Parser	119
D	A Massively Simpler Alignment Model	120
D.1	Motivation	120
D.1.1	Algorithm	124
D.1.2	Implementation and Evaluation	130
	Bibliography	133

List of Figures

I-1	IBM Models 1 and 2 pseudocode.	25
I-2	Typical states in the Köhn decoder. In the original state, depicted above, the decoder has hypothesised that the French phrase $f_3 f_4$ corresponds with the English phrase $e_1 e_2$ with score S . This state is updated by adding that the phrase $f_5 f_6$ corresponds to the English phrase $e_3 e_4$ with score terms corresponding to this assignment $(\text{Score}(e_3, e_4; f_5, f_6))$, log likelihood terms corresponding to the language model, and a score for the phrase reordering.	28
1-1	Types used in the Transducer	39
1-2	Transducer Pseudocode	40
2-1	A search graph produced by <code>moses</code> for the German sentence <i>Es gibt ein Haus.</i>	59
2-2	Algorithm for incorporating human assistance in phrase-based translation system.	62
2-3	Sample Mechanical Turk Questionnaire including a test question. . .	63

3-1	BLEU scores for a translation system with varying amounts of data using Kneser-Ney (KN) and Stupid Backoff (SB). BP/ $\times 2$ indicates the increase in the BLEU score (in 0.01 BLEU points, or BP) for every doubling of the size of the language model (in tokens). Differences greater than 0.51 BP are significant at the 0.05 level. Adapted from [5].	69
3-2	Find the optimal scores based on counts; this is a nonlinear gradient descent algorithm. \mathcal{H} denotes the Hessian of \mathcal{L} .	82
3-3	Active ranker pseudocode to rank N items.	83
3-4	Sample Mechanical Turk Questionnaire including a test question.	85
D-1	The length of English sentences drawn from the English translation of the proceedings of the European Parliament appears to have a Gaussian distribution.	121
D-2	(a) The length of English sentences whose German translations are twenty-five words long appears to have a Gaussian distribution. (b) The length of sentences in English and German appears to have a bivariate Gaussian distribution. All things are indeed Gaussian.	122
D-3	Our motivation: The indices of aligned words can be approximated by a Gaussian.	123
D-4	The Vanilla Gaussian Algorithm.	127
D-5	The performance of the alignments induced by the one-dimensional gaussian model is similar to those induced by the IBM-2 model. This graph shows the BLEU metric of the two models when applied to the EUROPARL training data and standard test set for the German-English language pair.	131

List of Tables

1.1	Examples of voice-changing and clefting.	37
1.2	To distinguish between the different conjugated forms in Arabic, we append labels to the English verb in the gloss; M and F denote masculine and feminine, respectively, and S, D, and P denote singular, dual, and plural, respectively.	46
1.3	The agreement tags of conjugated verbs in the Arabic treebank. PV, IV, and CV are perfect, imperfect, and imperative (command) aspects, respectively.	47
1.4	Leaf feature values in Arabic and English. Tags not given here are evenly distributed between all possible labels. Tags placed in multiple categories are evenly distributed between those categories. Also, note that we have included verbal categories for Arabic; this is in anticipation of future work with statistically inducing these features on verb phrases.	48
1.5	The model optimization algorithm.	49
1.6	The active-to-passive transduction grammar	53
1.7	The cleft transduction grammar	53

3.1	Amount of information content per word in nits (according to a trigram model) needed to obtain a certain BLEU score on the EUROPARL corpus. A BLEU score of 0.28 is state-of-the-art on this corpus, and corresponds to the translation model providing 3.5 nits of information per word, compared to the 10.0 nits required to obtain a perfect BLEU score.	66
3.2	Examples from Systran and Google translate (at the time of this writing), elucidating some possible issues with BLEU.	68
3.3	Our proposed metric, CLAIRE, captures the results of the shared evaluation and correlates strongly across different types of translation systems.	85

Introduction

I.1 Prologue

This thesis reflects an interest in the pursuit of machine translation in particular, although the methods and findings have bearing on the natural language processing enterprise in general.

We might define the natural language processing enterprise as comprising of methods engaged in *engineering systems that produce or consume natural language text in some way that corresponds to aspects inherent to that modality*.

Why the interest in this field? Putting aside the scientific interest an anthropologist might have in the “language phenomenon”, there has been an absolute explosion in the amount of natural language text that is available.

Natural language is sought after as a desirable interface between man and machine more and more as we expect our devices to perform increasingly complex tasks.

Even before the information explosion, a globalization phenomenon has required us to communicate with others across a language barrier.

Even when communicating with others in the same language, we are each personally required to produce more prodigious quantities of text than ever before.

What are the greatest desires of the natural language practitioner? What is the most that can reasonably be expected? Let’s for a moment imagine a world in which natural language were solved—it becomes easy to convert back and forth

between natural language and various logical forms. Precise queries could be executed to extract specific information from large bodies of natural language text. Large documents could be summarized and brief summaries could be expanded into large documents! Documents written in other languages could be translated accurately and quickly, preserving meaning, tone, meter, etc.

The possibilities are staggering and indeed have bewitched researchers since the early days of computation. Warren Weaver[47], an early cryptanalyst and cryptographer renowned for his work decrypting the ENIGMA code in World War 2, famously compared translation to codebreaking in the 1940's. Early natural language practitioners achieved small victories in the 1950's, and with exciting advancements in linguistics (particularly by Noam Chomsky), they promised the dreams outlined above to funding agencies in the US and abroad. Anyone familiar with the modern fruits of machine translation, taking into consideration the great advancements in the theory of computation, learning theory, linguistics, as well as the massive improvements in supporting infrastructure (microprocessors and datasets), would hardly be surprised by the demoralizing failures that were to come.

After pouring funding into one promising project after another for decades, the NSF, DOD, and CIA commissioned a report from the National Academy of Sciences to advise them on sensible next steps in the field. The National Academy of Sciences formed an Automatic Language Processing Advisory Committee (ALPAC) in 1964 to answer this need. ALPAC shared their findings in 1965 [42]. As a result of the report, support of machine translation projects was suspended in the United States for twenty years.

This report and its findings have a negative reputation in the statistical language community as an unjustified attack on a nascent field and which subsequently held the artificial intelligence community in general, and the natural language processing community in particular, back for decades.

ALPAC was chaired by John R. Pierce, a seasoned electrical engineer and physicist from Bell Labs who supervised the team that invented the transistor. ALPAC's

other members were the American psychologist John B. Carroll, known for his contributions to educational linguistics and psychometrics, then at Harvard; the American linguist Eric P. Hamp, who remains a professor of linguistics at the University of Chicago; David G. Hays of RAND, an important pioneer of computational and corpus linguistics (and the inventor of the term); American linguist Charles F. Hockett; American mathematician Anthony G. Oettinger, who is currently the chair of the Center for Information Policy Research at Harvard; and American computer scientist Alan Perlis, a pioneer in programming languages and the recipient of the first Turing Award. In short, the committee consisted of academic heavyweights of the highest caliber and included members who were interested in the continued support of fields relating to computation and linguistics.

The report is supported by extensive research, not only of the progress of the natural language processing enterprise itself, but also to establish the needs that the projects could reasonably be expected to satisfy.

The committee concluded that:

- Machine-aided translation tools (such as user-friendly technical dictionaries) improved both the speed of a translator's work and the quality of his output. This is an area of research worth supporting.
- The quality of automatic machine translation was too poor to be used without extensive post-editing by a translator, and this was more costly than simply having the translator translate from the source text directly. This result rested on extensive comprehension tests on scientific documents translated by then state-of-the-art systems.
- Furthermore, there was no prospect of improved translation systems under current research agendas.
- The only real, immediate need that could not be met by the available human translators is for rapid, accurate translation for a small circulation (too few to justify the cost of expedited human translation), and the research directions

that were being pursued were too superficial to achieve it. They went so far as to suggest that document typesetting was a bottleneck that could be more readily alleviated by development to improve the throughput of the conventional translation teams. Also, recurring needs of individual researchers for translations from a particular language could be alleviated with basic language training for that researcher in the language in question.

- Research had resulted in a number of useful results in related fields, including stimulating energetic research in corpus linguistics. This research should be supported with funding as a science with expectations of long-term improvement to automatic systems.
- There was an immediate need for improved evaluation methodologies.

The funding agencies responded by following the recommendations to eliminate funding to fully automatic systems, but did not follow recommendations to provide funding to supporting fields.

Perhaps the funding agencies deserve some derision for not taking the long view on translation. At the same time, it is clear that the popular view at the time the funding proposals were written was that translation would be solved with minimal government expenditure and the benefits were manifest; this was certainly reflected in the report. Series of failed projects could not continue to receive funding; if administrators in the funding agencies felt inclined to be sympathetic, they would certainly have been replaced by more skeptical colleagues at this point.

Nearly twenty years after the ALPAC report, the British computational linguist and computer scientist Martin Kay wrote a brief statement [27] expressing his view on the future of automatic machine translation. The statement is brief enough to include here in its entirety:

Large expenditures on fundamental scientific research are usually limited to the hard sciences. It is therefore entirely reasonable to suppose that, if

large sums of money are spent on machine translations, it will be with the clear expectation that what is being purchased is principally development and engineering, and that the result will contribute substantially to the solution of some pressing problem.

Anyone who accepts large (or small) sums on this understanding is either technically naive or dangerously cynical. It may certainly be that

1. machine translation could provide a valuable framework for fundamental research;
2. texts in highly restricted subsets of natural language could be devised for particular purposes and texts in [*sic*] translated automatically;
3. computers have an important role to fill in making translations;
4. translations of extremely low quality may be acceptable [*sic*] on occasions.

However,

1. the fundamental research is so far from applicability,
2. the language subsets are so restricted,
3. the useful computer technologies are so different from machine translation,
4. the quality of the translations that can be produced of natural texts by automatic means is so low, and
5. the occasions on which those translations could be useful are so rare,

that the use of the term in these cases can only result in confusion if not deception.

A determined attempt was made to bring machine translation to the point of usability in the sixties. It has become fashionable to deride these as “first generation” systems and to refer to what is being done now as belonging to the second or third generation. It should surely be possible for those who think that the newer systems can succeed where the earlier ones failed, to point to problems that have been solved since the sixties that are so crucial as substantially to change our assessment of what can be achieved. We know a good deal more about programming techniques and have larger machines to work with; we have more elegant theories of syntax and what modern linguists are pleased to call semantics; and there has been some exploratory work on anaphora. But, we still have little idea how to translate into a closely related language like French or German, English sentences containing such words as “he”, “she”, “it”, “not”, “and”, and “of”. Furthermore, such work as has been done on these problems has been studiously ignored by all those currently involved in developing systems.

Unfortunately, the sums that are being spent on MT in Europe and Japan are large enough to make virtually inevitable the production of a second ALPAC report sometime in the next few years. This will inevitably have a devastating effect on the whole field of computational linguistics, everywhere in the world. The report will be the more devastating for the fact that much of the money has in fact been spent frivolously, and much of the work has been incompetent, even by today’s limited standards.

Fortunately for us, Kay’s predictions on funding have not come true. The needs that ALPAC reported have grown by leaps and bounds, and machine translation systems have found use in niche applications. For example, the Canadian weather service has famously used machine translation (with a human post-editing phase) to translate weather bulletins from English to French for 30 years (including using the same system, METEO, for two decades 1981–2001) [45]. In 2009, President

Obama released a “Strategy for American Innovation” which named “automatic, highly accurate and real-time translation between the major languages of the world” an ambitious goal that will “improve our quality of life and establish the foundation for the industries and jobs of the future.” [36]

How can we distinguish our work from the pseudo-science that ALPAC and Kay described? It is difficult to declare natural language processing to be a science since our goal is not to learn about an existing system, but is instead to build useful systems of our own. I would argue that it is still possible to do science in this arena, but that it requires care to understand the limitations of our results. It’s difficult to make broad statements about the value of a particular approach when so much is still unknown about language in general and considering how far state-of-the-art systems fall short of the dream. It is not at all inconceivable that the best research systems extant would bear little resemblance to a future system that fulfills the promises of our predecessors.

In spite of this uncertainty, there is a definitive mainstream thrust of research in the natural language community: **Quantity Leads to Quality**. State-of-the-art performance is more surely and readily achieved by an appeal to a massive dataset than an appeal to linguistic theory. An extension of this is that improvements in quality are achieved by more advanced statistical models that are capable of modelling more exotic relationships between input and output, with careful regularization to account for the sparsity of data.

Indeed, researchers entering natural language processing are often warned of the consequences of expecting and promising too much, and veteran researchers are careful to establish metrics that can illustrate sustained, gradual improvement and to point out the immediate uses of the state-of-the-art translation systems, notwithstanding their manifestly poor quality. Church and Hovy go as far as stating that it may be more important to improve performance by seeking an appropriate application than by improving the system in question [15].

I.2 Background

I.2.1 Parsing

Our baseline parsing model is Collins Model 2 parser. We describe it here briefly; please see [16] for more details.

Just as a statistical speech recognition system aims to maximize $\Pr(\mathbf{e} | \mathbf{a})$, given a sentence S , the statistical parser aims to find the parse tree T that maximizes $\Pr(T | S)$. The Collins parsers are generative models based on the probabilistic context free grammar formalism (PCFG), except that they are lexicalized. An unlexicalized PCFG would be broken down as follows:

$$\begin{aligned} \arg \max_T \Pr(T | S) &= \arg \max_T \Pr(T, S) \\ &= \arg \max_T \prod_i \Pr(RHS_i | LHS_i), \end{aligned}$$

where RHS_i and LHS_i denote the left- and right-hand sides of the context-free grammar rule that is used at the i th step in the derivation of the tree; the probabilities of the rules make up the parameters of the model, and maximum likelihood estimates are easy to obtain from a corpus of trees.

Collins lexicalizes the PCFG by adding every possible word and part-of-speech to each non-terminal in the grammar, greatly increasing the number of parameters. The maximum-likelihood estimates can be obtained as before (technically, this is still a PCFG), but data sparsity quickly becomes an issue, so the generative story is broken down further. Each LHS now generates a head¹, and the head generates *subcategorization frames* for the left and the right (just a set of required complements on each side of the head phrase); each terminal of the rule to the left and right of the head is then generated depending only on the head and the constraint that the final rule match the subcategorization frame. Model 2 also adds a flag which indicates

¹The head itself is not easy to define and is the subject of some debate. Collins provides hand-designed head rules in his thesis which seem to work well for the Wall Street Journal. In [4], Bikel provides similar rules for use in other languages.

whether each nonterminal is a complement or an adjunct.

I.2.2 Machine Translation

Let us state the machine translation problem: our goal is to translate French² sentences to English³ sentences. For simplicity, we will consider sentences in isolation, ignoring the impact of context. We will always denote sentences in the source language by \mathbf{f} and sentences in the target language by \mathbf{e} . m is the number of words in the source sentence \mathbf{f} and ℓ is the number of words in the target sentence \mathbf{e} . Now, in Bayesian terms, given a French sentence \mathbf{f} , we wish to find the English sentence \mathbf{e} that maximizes $\Pr(\mathbf{e}|\mathbf{f})$, in effect imagining that French sentences are generated by some unknown transformation on English sentences. Using Bayes' Law, we write:

$$\begin{aligned}\arg \max_{\mathbf{e}} \Pr(\mathbf{e}|\mathbf{f}) &= \arg \max_{\mathbf{e}} \frac{\Pr(\mathbf{f}|\mathbf{e})\Pr(\mathbf{e})}{\Pr(\mathbf{f})} \\ &= \arg \max_{\mathbf{e}} \Pr(\mathbf{f}|\mathbf{e})\Pr(\mathbf{e}).\end{aligned}$$

The $\Pr(\mathbf{f})$ term can be ignored since \mathbf{f} is constant. The first term, $\Pr(\mathbf{f}|\mathbf{e})$, is called the “translation model” and the second, $\Pr(\mathbf{e})$, is called the “language model.” It is just good to know that our translation model does not have to worry about assigning low probabilities to English sentences that look like they could be translations of \mathbf{f} but don't really look like they could be English sentences; a good language model can make up for some deficiencies in the translation model. One other advantage of factoring the model in this way is that the language model can be trained on very large unlabeled (i.e., untranslated) data sets in the target language.

In this work, we use a trigram language model; that is, we assume that the procedure that produces English sentences is a Markov process with a history of two

²or in the general case, Foreign. The original papers translated from French to English, and it has become a tradition.

³or in the general case, Native.

words: the probability of an English sentence \mathbf{e} is broken down like this:

$$\Pr(\mathbf{e}) = \Pr(e_1, e_2, \dots, e_\ell) = \prod_{i=1}^{\ell+1} \Pr(e_i | e_{i-1}, e_{i-2}),$$

where $e_{\ell+1}$ is a stop symbol implicitly included at the end of every English sentence. This is a raw trigram model. To learn the parameters, we can simply count the number of times each trigram appears in a corpus:

$$\Pr(e | e_2, e_1) = \frac{\text{count}(e_1, e_2, e)}{\text{count}(e_1, e_2)}$$

where $\text{count}(\dots)$ denotes the number of times the words \dots appear together in the corpus in the given order. One problem with this model is that it will assign a zero probability to any sentence that has a trigram that was never seen in the corpus. To fix this, one uses a *smoothed* trigram model [22]:

$$\Pr(e | e_2, e_1) = \alpha_t \Pr_t(e | e_2, e_1) + \alpha_b \Pr_b(e | e_2) + \alpha_m \Pr_m(e)$$

where $\alpha_t + \alpha_b + \alpha_m = 1$, the α s are nonnegative, and $\Pr_t(\cdot)$, $\Pr_b(\cdot)$, and $\Pr_m(\cdot)$ denote trigram, bigram, and unigram probabilities, respectively. This is called an *affine* combination of the three distributions.

I.2.3 Five IBM Models

In a seminal 1993 paper, Brown *et al* introduced a set of five machine translation systems based on fairly simple statistical models (under the “noisy-channel” framework described above) and large parallel corpora. The later models are significantly more complex than the earlier models, with each subsequent model corresponding to an increase in complexity and improved accuracy (but with diminishing returns in the later models) [7].

We assume that we have at our disposal a corpus of N pairs of sentences $(\mathbf{e}^{(1)}, \mathbf{f}^{(1)})$,

$(\mathbf{e}^{(2)}, \mathbf{f}^{(2)}), \dots, (\mathbf{e}^{(N)}, \mathbf{f}^{(N)})$.

In this document, our analysis will not extend beyond the first and second IBM Models, so we will limit our discussion of the later models to a brief overview.

Model 1

We begin by describing an idea fundamental to both Model 1 and Model 2: an *alignment* between a pair of sentences \mathbf{f} and \mathbf{e} is an ordered set $a_1, a_2, \dots, a_m \in \{0, 1, \dots, \ell\}$. The French word f_j is said to be *aligned* to the English word e_{a_j} (where e_0 denotes a fake word “NULL” that is used to explain function words that may not have an obvious analog in the English sentence). Notice that we don’t demand that English words are aligned to French words; a single English word could be used to explain an entire French sentence (a developed model would declare such an alignment as very improbable, however).

Model 1 makes the following assumptions/approximations:

- All alignments are equally likely.
- All French sentence lengths m are equally likely (we will ignore the obvious problem that there are infinitely many French sentence lengths⁴; if it bothers you, you can assume that someone gives you the length or that there are only finitely many possible French sentence lengths, which is true in practice, anyway). We will generally omit this term.
- Each word is translated independently of the other words.

These assumptions are outrageously simplifying, but it is important to start with a tractable model. It is also important to remember that our language model will clean up output problems: we can expect short-range alignment problems and, to some extent, poor grammar to be dealt with there.

⁴There really are. As proof, I present a regular expression that matches an infinite number of grammatical French sentences: *Je suis un très* grand singe*. Obvious analogs exist in other languages.

Ultimately, we obtain the following formulation⁵:

$$\begin{aligned}
 \Pr(\mathbf{f} | \mathbf{e}) &= \frac{1}{(\ell + 1)^m} \sum_{\mathbf{a}} \prod_{j=1}^m \Pr(f_j | e_{a_j}) \\
 &= \frac{1}{(\ell + 1)^m} \sum_{a_1=0}^{\ell} \sum_{a_2=0}^{\ell} \cdots \sum_{a_m=0}^{\ell} \prod_{j=1}^m \Pr(f_j | e_{a_j}) \\
 &= \frac{1}{(\ell + 1)^m} \prod_{j=1}^m \sum_{a_j=0}^{\ell} \Pr(f_j | e_{a_j}).
 \end{aligned}$$

Here, a_j is the index of the English word that is *aligned* to the j th French word; Model 1’s parameters are the translation probabilities $\Pr(f | e)$ (the probability that the French word f was generated by the English word e). Model 1 is an excellent candidate for optimization by EM; it is convex and has only one local maximum (outside of saddle points due to symmetry) so given a random starting point, it will always converge to the same, optimum translation table. Some more math gives us the following update rules:

$$T'(f | e) = \frac{1}{Z_T(e)} \sum_{\substack{i,j,k \\ \mathbf{e}_i^{(k)}=e, \mathbf{f}_j^{(k)}=f}} \frac{T(\mathbf{f}_j^{(k)} | \mathbf{e}_i^{(k)})}{\sum_{i'=0}^{\ell} T(\mathbf{f}_j^{(k)}, \mathbf{e}_{i'}^{(k)})},$$

where Z_T is a normalization constant.

Model 2

In Model 2, we wish to relax Model 1’s assumption that all alignments are equally likely. However, we will assume for simplicity that the words all “move” independently; that is, which English word a French word is aligned to is independent of the

⁵The reversal of the product and sum is an important trick, since it makes it easy to optimize the terms independently.

alignment of the remaining words. Here is the formulation of Model 2:

$$\begin{aligned} \Pr(\mathbf{f} | \mathbf{e}) &= \sum_{\mathbf{a}} \prod_{j=1}^m \Pr(f_j | e_{a_j}) \Pr(a_j | j, \ell, m) \\ &= \prod_{j=1}^m \sum_{a_j=0}^{\ell} \Pr(f_j | e_{a_j}) \Pr(a_j | j, \ell, m). \end{aligned}$$

Here, in addition to the translation probabilities $\Pr(f | e)$ Model 2 inherits from Model 1, we find alignment probabilities, $\Pr(a_j | j, \ell, m)$ (the probability that, given a French sentence of length m that is the translation of an English sentence of length ℓ , the j th French word was generated by the a_j th English word). Model 2 is not as good a candidate for EM as Model 1 was; it is riddled with saddle points and local maxima. Typically, one initializes the translation parameters by training Model 1 before training Model 2, whence we use the following update rules:

$$\begin{aligned} T'(f | e) &= \frac{1}{Z_T(e)} \sum_{\substack{i, j, k \\ \mathbf{e}_i^{(k)} = e, \mathbf{f}_j^{(k)} = f}} \frac{D(a_j = i | j, \ell, m) T(\mathbf{f}_j^{(k)} | \mathbf{e}_i^{(k)})}{\sum_{i'=0}^{\ell} D(a_j = i' | j, \ell, m) T(\mathbf{f}_j^{(k)}, \mathbf{e}_{i'}^{(k)})} \\ D'(a_j = i | j, \ell, m) &= \frac{1}{Z_D(j, \ell, m)} \sum_{\substack{k \\ |\mathbf{e}^{(k)}| = \ell, |\mathbf{f}^{(k)}| = m}} \frac{D(a_j = i | j, \ell, m) T(\mathbf{f}_j^{(k)} | \mathbf{e}_i^{(k)})}{\sum_{i'=0}^{\ell} D(a_j = i' | j, \ell, m) T(\mathbf{f}_j^{(k)}, \mathbf{e}_{i'}^{(k)})}, \end{aligned}$$

where Z_T and Z_D are normalization constants.

The pseudocode to train IBM Models 1 and 2 is given in figure I-1.

Models 3, 4, and 5

IBM Model 3 introduces fertility parameters, modeling the number of French words a single English word generates. IBM Model 4 introduces distortion parameters to the alignment models as a way of encouraging words to move in groups. Both of these models are formulated “deficiently”: that is, they assign probability mass to impossible French sentences (four-word French sentences without a third word, for instance); Model 5 is the non-deficient version of Model 4. Since this makes little

```

Initialize  $t(f|e)$  and  $D(i|j, \ell, m)$ 
do:
· zero  $t'(f|e)$  and  $D'(i|j, \ell, m)$ 
· for  $(\mathbf{e}, \mathbf{f})$  in corpus:
· ·  $m = |\mathbf{f}|, \ell = |\mathbf{e}|$ 
· · for  $j=1 \dots m$ :
· · · for  $i=0 \dots \ell$ :
· · · ·  $a_i = t(\mathbf{f}_j | \mathbf{e}_i) \cdot D(i|j, \ell, m)$ 
· · · ·  $a_i = a_i / (\sum_{i'} a_{i'})$ 
· · · · for  $i=0 \dots \ell$ :
· · · · ·  $t'(\mathbf{f}_j | \mathbf{e}_i) = t(\mathbf{f}_j | \mathbf{e}_i) + a_i$ 
· · · · ·  $D'(i|j, \ell, m) = D(i|j, \ell, m) + a_i$ 
·  $t'(f|e) = t(f|e) / \left( \sum_{f'} t'(f'|e) \right)$ 
·  $D'(i|j, \ell, m) = D(i|j, \ell, m) / \left( \sum_{i'} D'(i'|j, \ell, m) \right)$  (Set to 1 for Model 1)
·  $t = t', D = D'$ 
until convergence

```

Figure I-1: IBM Models 1 and 2 pseudocode.

empirical difference and is a great computational burden, Model 5 is rarely used in practice [38].

I.2.4 Phrase-Based Models

The primary unit of information in all of the systems we have described up to this point is the word; in phrase-based systems, the primary unit of information is the phrase, a collection of (lexically) consecutive words and the lexical entry in a phrase-based system is a triple containing a source phrase $f_1 \dots f_n$, a target phrase $e_1 \dots e_m$, and a score $s \in [0, 1]$. That is, instead of considering probabilities of word-to-word translations and word-movement, a phrase-based system will deal with probabilities of phrase-to-phrase translations and phrase-movement. There is a great deal of evidence to suggest that machine translation systems generally experience a performance boost by making this change.

Some phrase-based models, such as Marcu’s [33], simply introduce mechanisms for phrase-to-phrase translations and invent policies to assign probability mass to phrase-to-phrase translations. Others, such as Köhn’s [31], build a dictionary of phrases from other information sources. Our experiments are centered on the Köhn system, as it achieves state-of-the-art performance.

There are a number of ways one can build phrase dictionaries depending on the data available. Phrases can be built from word-based alignments (such as those generated by the IBM Models). If syntactic information is available, it can be used to restrict our attention to syntactic phrases; although it seems that syntactic phrases may be more useful, experiment suggests that phrases that are not syntactically motivated are, generally, just as useful⁶. Furthermore, even weighting syntactic phrases produces virtually no improvement at best and is sometimes harmful. Phrase dictionaries can also be built from phrase-aligned data generated from phrase-based systems. Again, we can place more confidence in these phrases if we wish, but gen-

⁶For example, the German phrase “es gibt” corresponds nicely to the English phrase “there is”, even though they cross a syntax boundary in both languages.

erally, the lesson from the experiments with syntactic phrases is applicable: it is better in practice simply to consider as many phrases as possible than to restrict our knowledge to satisfy any bias we may have [29].

Experiments by Köhn *et al* show that simple heuristic methods based on word-based alignments from the IBM models generate state-of-the-art translations. To generate a phrase dictionary, he begins by observing that the IBM models are not symmetric; the alignments generated by a model trained to translate from French to English can be different from alignments generated by a model trained to translate from English to French (in fact, it is often impossible for alignments generated in one direction to match those generated in the other direction, due to inherent restrictions of the IBM models). Köhn's method begins by considering the intersection of the two alignments as a starting point for the phrases it must generate; that is, it begins by suggesting that words that are aligned in both models are probably related. Next, Köhn uses a growing technique to induce the phrase dictionary; phrase dictionaries generated in this fashion tend to be very large because of the generality of this technique; however, the method naturally also generates scores and generates many very low-scoring phrase pairs.

Decoding is done using an algorithm described in [24]. The output sentence is generated left-to-right in the form of partial translations which are formed using a beam-search algorithm including scores for the phrase translations and for the language model, rewards for increased sentence length, and penalties for phrase reordering. Typical states are depicted in Figure I-2.

The two factors that govern the quality of translations generated using this technique are the quantity and quality of the alignments that are given to the system during training. In practice, translation quality from this method is significantly better than any of the IBM Models. Surprisingly, it does almost as well with IBM Model 2 alignments as with IBM Model 4 alignments.

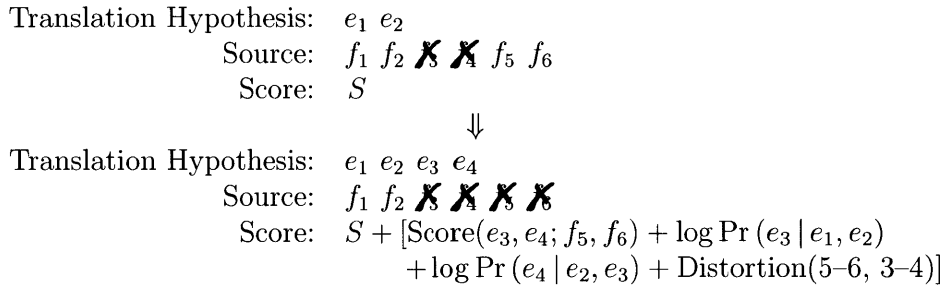


Figure I-2: Typical states in the Köhn decoder. In the original state, depicted above, the decoder has hypothesised that the French phrase $f_3 f_4$ corresponds with the English phrase $e_1 e_2$ with score S . This state is updated by adding that the phrase $f_5 f_6$ corresponds to the English phrase $e_3 e_4$ with score terms corresponding to this assignment ($\text{Score}(e_3, e_4; f_5, f_6)$), log likelihood terms corresponding to the language model, and a score for the phrase reordering.

I.3 Answering Kay

It is now, armed with a metric and a modern system design substantially different from those that Kay described in his statement [27] above, that we can respond to his demands to point to the lessons that have been learned and which summarize the mainstream agenda of the field.

1. Noisy-channel methods, like those described above, are vastly superior to rule-based systems, and require substantially less effort to design.
2. Linguistic formalisms can be incorporated into statistical translation models, but the effort required to do so is substantial, and the effect on the scores does not justify the effort.
3. n -gram language models are easy to train and tough to beat.

This second point is now a byword in the field; the following quote, by Fred Jelinek, famously conveys this:

Every time I fire a linguist, the performance goes up!

We will review these and other lessons in Chapter 3.

I.4 Outline

Chapter 1 presents details on our work in using human assistance for parsing, with an eye toward but without application to machine translation. Chapter 2 presents our work on word-sense disambiguation for machine translation using source-language human assistance. Chapter 3 revisits Kay's question, presents a critique of automatic metrics, and presents CLAIRE, our alternative. We close with conclusions and a few useful appendices.



Chapter 1

Human Assisted Parsing

What makes language processing tasks difficult? Early practitioners believed that better-than-human performance was just around the corner for many problems that remain enticingly unsolved, despite decades of active inquiry.

Perhaps the difficulties can be summarized as two broad issues: linguistic formalism and language ambiguity.

The former corresponds roughly to choice of representation, which is indeed a major concern. We do not wish to restrict ourselves to a synthetic or controlled language; we wish to permit all of the ambiguity and flexibility of open-domain natural language. Logical forms are too fragile and disconnected from the text, and traditional parse trees do not supply much of the information needed to make confident use of the text, and yet often contain a great deal of information that corresponds to linguistic hair-splitting and not to any valuable distinction in emphasis, structure, or meaning. This problem corresponds to a major share of the scope of the field of linguistics, so it comes as no surprise that we find it challenging.

However, in fully automatic systems, genuine ambiguity is such a great concern that the linguistic formalism is often regarded as a minor issue with little impact on performance. This is exemplified by the fact that the techniques that are known to consistently improve the quality of modern systems are those that reduce ambiguity

in the problem (most notably, constraining the domain of the text in question).

Consider the classic example:

I saw the man with the telescope.

There are several ambiguities in this relatively short, simple sentence, including:

- the intended sense of the word *saw* (either past tense of the verb *to see* or present tense of the verb *to saw*);
- the attachment of the prepositional phrase *with the telescope*, either to the verb *saw* or the object *the man*.
- the many possible meanings of the preposition *with* (which may be reduced given a particular attachment choice):
 - I used the telescope to see the man.
 - I saw the man when I saw the telescope.
 - I saw the man adjacent to the telescope.
 - I saw the man who owns a telescope.

and so on. Clearly there is ambiguity at many levels (lexical, syntactic, semantic), resulting in a combinatorial explosion in the number of readings of a sentence. In spite of this, language users generally have such little difficulty winnowing them that the ambiguities are often not even evident without explanation. Most English speakers would never consider the possibility that the author of this sentence wished to convey that he was using a telescope to cut a man in two, yet it is not difficult to imagine a context (perhaps a crime novel) where this *would* be the intended meaning of the sentence. In any case, two of the readings (that he used the telescope to see the man and that the man had the telescope when he saw him) are entirely plausible, and, without context, no system should entirely eliminate either of these readings.

Our goal is a system that enables a non-expert user to perform expert parsing with performance that exceeds that of fully automatic systems. We propose a transformation-based approach: users input a sentence for processing, and candidate readings of the sentence are transformed and presented to the user for correction.

Why would such a system be valuable? Some processing tasks currently must be performed by experts because the quality of fully automatic systems is still too poor for many uses. Researchers have developed tools to assist the experts in tasks such as computer assisted translation, but expert labeling is still prohibitively expensive. On the other hand, non-experts are in bountiful supply. Furthermore, as long as the performance of human-assisted systems surpasses that of fully automatic systems, their outputs could be used to augment existing corpora.

It may even be argued that, in lieu of strong artificial intelligence, fully-automatic systems will never be able to resolve world-knowledge or common-sense ambiguities. Certainly, without some sort of reasoning subsystem, it is inconceivable that we will be able to properly respond to context, even if we were to solve syntax and reverse-engineer the language center. Without the mind's superior pattern-matching machinery and world model, language understanding would (probably) still be unsolved.

What difficulties might we face? First of all, an improvement in quality may require prohibitively great effort on the part of the user. This is particularly true of a system that cannot generate a reasonable partial understanding automatically, in which case it may be difficult to form meaningful questions.

In spite of these possibilities, there may yet be hope. Ambiguities that cannot be expressed with transformation (i.e., by *rewriting the sentence*) do not exist for all practical purposes, so an ideal system should be able to communicate any difficulty it has. Anyone is qualified to respond: virtually every human being has a first-class world-model and pattern-recognizer waiting to be used. The labels he gives for one sentence, or even for part of one sentence, may be useful for subsequent ambiguities, so we can hope that the number of questions we have to ask will increase more slowly

than the amount of information we are interested in (that is, we have to know a great deal to ask an intelligent question, but much of that knowledge will be helpful for the following question). Finally, we can hope that the number of questions we need to ask overall will gradually diminish as we collect enough data to better estimate the answers.

1.1 Previous Work

Many recent papers have targeted specific attachment types directly with automatic methods, by reranking or self-training. These have achieved significant gains, but the authors are quick to note that these gains are specific to the corpus being used (i.e., the Wall Street Journal corpus), and clearly do not represent gains in broad-domain text.

Disambiguating by transformation is not a radical thing; the Treebank parsing guidelines include criteria for especially tricky examples by syntactically transforming difficult sentences.[3]

However, perhaps the earliest paper to suggest disambiguation by asking in any formal context is Kay’s description of the MIND system [26], a multilingual translation engine. Later proposals for the use of monolingual consultants in translation are described in Chapter 2.

In parsing, Carter [13] presents TreeBanker, a system similar to an earlier system by Tomita for translation [46], but targeted at expert consultants and with an eye toward corpus-building.

1.2 Our Approach

To reiterate, our goal is to improve quality of parsing using transformations, with an eye toward applications that would make use of argument structure.

Our approach differs primarily from the previous approaches in three respects: first, we incorporate interaction in a statistical (not rule-based) parser; second, our mode of interaction with the user is by way of transformations of the original text using a probabilistic tree transducer; third, instead of trying to design tests for ambiguity patterns, we propose using linguistically meaningful transformations *without* an eye toward particular ambiguity classes. The advantage of favoring statistical techniques is that they are less labor-intensive and less fragile than their rule-based counterparts, at the cost of being mechanistically opaque to system designers. Considering that we are expecting human attention, using probabilistic techniques confers another serious advantage: the system can be designed to continuously adapt its behavior to new data.

Ideally, the machine learning practitioner would have sentences in some canonical form, along with corresponding versions authoritatively transformed in some way. In lieu of such a corpus, we are consigned to manual design methods. We therefore produced a list of transformation types and designed the corresponding transduction rules. We envision that statistics would play its part by incorporating failure probabilities and in lexicalizing the rules.

Machine learning concepts also motivate how we question the user. We envision a system in which the user may choose to stop answering questions at any time, so we assume every question we ask may be the last. Consequently, we will always ask the question that gives us the most information in expectation. The amount of information a question and answer pair gives us is given by the change in the entropy of the belief distribution, which is based on probabilities given by our baseline parser.

1.2.1 Collins Parser

Our baseline parser is the Model 2 Collins parser (described briefly in the introduction and in full in [16]) modified to recover the full Penn extended tag-set ([20] and Appendix B). This parser is fully automatic and is a standard baseline, because of

its relative simplicity and high performance. It recovers labeled constituents with $\sim 88\%$ precision/recall. Certain recurring attachment decisions are highly ambiguous because constituents can feasibly be attached to many previous points in a sentence; consequently, these correspond to a substantial portion of the error that Collins reports. For instance, in the sentence:

She announced a program to promote safety in trucks and vans.

the prepositional phrase *in trucks and vans* can attach to several points in the sentence:

- Safety in trucks and vans is what the program promotes.
- The program promotes safety, and it does so in trucks and vans.
- The announcement was made in trucks and vans.

and the conjunct *and vans* can attach to several points in the sentence as well:

- Safety in trucks and vans is what the program promotes.
- The program promotes safety and also promotes vans.
- She announced a program and she announced vans.

So it comes as no surprise that the scores for these attachment types are lower: prepositional phrase attachments are recovered with $\sim 82\%$ precision/recall, and coordinating conjunction attachments with $\sim 62\%$ precision/recall.

1.2.2 Algorithm

We make use of any available syntactic transformations that are meaningful linguistically (ones we can expect a non-expert consultant to understand). We select which transformation to use based on our belief, which is initially just the probability scores

given by the parser, and which is updated by any information we obtain from the user(s).

Our intended algorithm, then, is as follows:

Given a sentence s , obtain top N output from the Collins parser

Transform this sentence using the available rules for each candidate parse

Group the outputs of the transformation by the word span of the constituents

do

- For each collection of word spans, compute the entropy of the partition
- Ask the user to label the outputs corresponding to the maximum entropy
- Adjust the scores of the parses

until the user quits or we run out of questions

return the highest scoring parse

Now we approach the issue of the transformations themselves.

1.2.3 Tree Transduction Grammar

For the purposes of this chapter, we will only be transforming clauses from active to passive voice or clefting sentences. However, since we expect to use other transformations in later work, and since changing voice in particular is a fairly complex procedure, we will sketch the formal machinery and attach our grammars for both transformations (see the end of chapter).

A transducer is a finite state automaton with two tapes. We use a transducer that operates on the internal nodes of a tree, and call this machine a tree transducer. The match or input portion of the transducer amounts to little more than executing an exhaustive regular expression matcher at each level of the tree for each rule, and then making the results of the match available to parent nodes for recursive matching. The output portion is slightly more complicated: since changing voice can involve long-range movement of constituents, output rules have to be able to “pass” trees as

Original Sentence	John saw <u>the man</u> <u>with the telescope</u> .
Passive	<u>The man</u> <u>with the telescope</u> was seen by John. <u>The man</u> was seen <u>with the telescope</u> by John.
Cleft	It is John that saw <u>the man</u> <u>with the telescope</u> . It is <u>the man</u> <u>with the telescope</u> that John saw. It is <u>the man</u> that John saw <u>with the telescope</u> . It is <u>with the telescope</u> that John saw <u>the man</u> .

Table 1.1: Examples of voice-changing and clefting.

arguments to children outputs. To leave this stage as flexible as possible, the output rules are simply lambda expressions. Stylized outputs from the active-to-passive rule and the cleft rule are given in table 1.2.3, illustrating how they may be used to disambiguate prepositional phrase attachments.

Once again, the input to the transducer is a *parse forest*—the dynamic programming table that is generated when the parser analyzes a sentence—and the output is a collection of all of the possible outputs of the original sentence and the amount of information we would gain by knowing whether or not each output is correct (e.g., if it is semantically equivalent).

We define a tree transduction grammar to be a collection of tree transduction rules, and we define a tree transduction rule to be a triple consisting of a label, an input rule, and an output rule. Input rules map a particular node in a candidate tree to a set (possibly empty) of “matches” (generally, input rules will correspond to regular expressions) and output rules will map a match to an output for that node. A particular node label is designated to be the root of the grammar; the output of the transducer is the set all of the possible outputs labeled with the root of the grammar. The input rules can require that a child match another input rule in the grammar, so the transduction grammar is a graph as well. We have constructed two tree transduction grammars: one to change sentences from active to passive voice,

and another to cleft a sentence’s arguments.

The transducer does the work of applying every rule to every node in the parse forest without repeating work; in general, neither the transduction grammar nor the parse forest will be acyclic (a parse tree containing a “cycle” is possible because the Collins grammar permits unary productions; such a tree will never be the maximum likelihood tree, however).

An early version of the transducer operated on single parse trees and was applied to top N output of the Collins parser, but the transducer was modified after our early experiments to operate on parse forests (i.e., on the parser’s dynamic programming chart) instead. This posed a significant technical challenge, but ultimately (somewhat surprisingly) resulted in a much more efficient system (since there are generally only minor differences between the top N trees, much of the work the early system was doing was repeated). The modifications to the code are described in appendix C, and we have made an implementation of the transducer (in python) available online for future work.

The transducer is a very powerful tool, but because it’s so flexible, it’s difficult to give performance guarantees. In practice, even for the types of grammars that we gave, the performance is decent; on a modestly powered computer, it can produce outputs for a 40-word sentence in about two seconds, comparable to the time it takes to parse the sentence.

The grammars for the active-to-passive and cleft transformations are given at the end of the chapter.

For each output that is generated by the system, we can compute the probability that the nodes that produced the output would appear in the correct parse tree. If we were completely confident that each grammar produced perfect output, we simply choose the output that partitions the space of possibilities most evenly. This does not change if we don’t completely trust the grammars, as long as we trust them all equally, which we do for simplicity.

FOREST	:=	List of NODES
NODE	:=	label headLabel headWord spanStart, spanEnd (sentence span) List of CHILDRENS
CHILDREN	:=	score List of NODES
MATCH	:=	content NODE CHILDREN spanStart, spanEnd (subspan of the CHILDREN, for partial matches)
TGRAMMAR	:=	List of TRULES rootLabel
TRULE	:=	label INPUTRULE OUTPUTRULE
INPUTRULE	:=	RERULE NAMEDRULE NOTRULE KLEENERULE ANDRULE ORRULE CONCATRULE mapping from (NODE, CHILDREN) \mapsto MATCH
RERULE	:=	regularExpression denoted by "<regexp>"
NAMEDRULE	:=	label denoted by "#<label>"
NOTRULE	:=	INPUTRULE denoted by "!<rule>"
KLEENERULE	:=	INPUTRULE denoted by "<rule>*"
ANDRULE	:=	List of INPUTRULES denoted by "<rule> & <rule> & ... & <rule>"
ORRULE	:=	List of INPUTRULES denoted by "<rule> <rule> ... <rule>"
CONCATRULE	:=	List of INPUTRULES denoted by "<rule> @ <rule> @ ... @ <rule>"
OUTPUTRULE	:=	mapping which accepts MATCH

Figure 1-1: Types used in the Transducer

```

toExamine = {Leaf Rules} × {(Nodes, Children)}
toExamine' = {}
partialMatches = empty hash with default value of {}
matches = empty hash with default value of {}
while |toExamine| > 0:
  · for each rule, node, children in toExamine:
  ·   ·  $n = |\text{children}|$ 
  ·   · preMatchCount = |matches[rule, node, children]|
  ·   ·   + |partialMatches[rule, node, children]|
  ·   · for each subspan of children:
  ·   ·   · update partialMatches[rule, node, children, subspan]
  ·   ·   · update matches[rule, node, children, subspan]
  ·   ·   · postMatchCount = |matches[rule, node, children]|
  ·   ·   ·   + |partialMatches[rule, node, children]|
  ·   ·   · if preMatchCount  $\geq$  postMatchCount:
  ·   ·   ·   · continue
  ·   ·   · for each parent rule pRule of rule:
  ·   ·   ·   · for each parent node pNode of node:
  ·   ·   ·   ·   · for each pChildren in of pNode.children:
  ·   ·   ·   ·   ·   · if node  $\in$  pChildren:
  ·   ·   ·   ·   ·   ·   · toExamine'.add((pRule, pNode, pChildren))
  · toExamine = toExamine'
  · toExamine' = {}
return matches

```

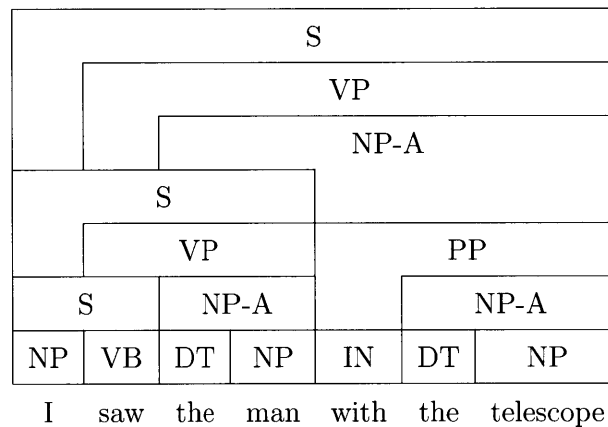
Figure 1-2: Transducer Pseudocode

1.2.4 Sample Output

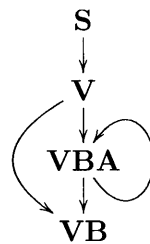
We will briefly run our sample sentence “I saw the man with the telescope” through the active-to-passive transducer to make its operation more clear. Let us suppose that our parser produced the following set of possible subtrees:

{(NP, 0-1), (VB, 1-2), (DT, 2-3), (NP, 3-4), (IN, 4-5), (DT, 5-6), (NP, 6-7),
 (S, 0-2), (NP-A, 2-4), (NP-A, 5-7), (VP, 1-4), (PP, 4-7), (S, 0-4), (NP-A, 2-7),
 (VP, 1-7), (S, 0-7)}

This corresponds to the following chart:



The active-to-passive transduction grammar (given at the end of this chapter) has the following dependency structure:



We begin by applying the input rule of the leaves of this dependency structure on the nodes in the parse forest in reverse tree-order (so that a node will only be processed after the children are processed). The algorithm begins by setting the

variable toExamine accordingly (we exclude the children of each node):

$$\text{toExamine} \leftarrow \{(\mathbf{VB}, (\text{NP}, 0-1)), (\mathbf{VB}, (\text{VB}, 1-2)), \dots, (\mathbf{VB}, (\text{S}, 0-7))\}$$

After the first iteration of the transducer, **VB** will match the node (VB, 1-2), so the parents of **VB** must be examined:

$$\begin{aligned} \text{toExamine} \leftarrow \{(\mathbf{V}, (\text{S}, 0-2)), (\mathbf{V}, (\text{VP}, 1-4)), (\mathbf{V}, (\text{VP}, 1-7)), \\ (\mathbf{VBA}, (\text{S}, 0-2)), (\mathbf{VBA}, (\text{VP}, 1-4)), (\mathbf{VBA}, (\text{VP}, 1-7))\} \end{aligned}$$

After the second iteration of the transducer, **V** will match the nodes (VP, 1-4) and (VP, 1-7) with outputs ($\dagger \mapsto$ was seen by \dagger , $\dagger \mapsto$ the man) and ($\dagger \mapsto$ was seen by \dagger , $\dagger \mapsto$ the man with the telescope), respectively. Now the parents of **V** must be examined (this time, we include the children):

$$\begin{aligned} \text{toExamine} \leftarrow \{(\mathbf{S}, \text{S} \rightarrow \text{NP}(0-1) \text{VP}(1-7)), \\ (\mathbf{S}, \text{S} \rightarrow \text{NP}(0-1) \text{VP}(1-4) \text{PP}(4-7))\} \end{aligned}$$

After the third and final iteration of the transducer, **S** matches both nodes and yields the following outputs:

The man was seen by me with the telescope.

The man with the telescope was seen by me.

1.3 Further Augmenting the Tree

The trees in the Penn treebank have a great deal of information beyond flat sentences in the Wall Street Journal; they include basic part-of-speech information at the word level, constituent dependency information, and information about the argument structure. This information is so rich that the vanilla Collins parser does

not attempt to recover it. However, there is an absence of many syntactic features on internal nodes even on these trees. This is true even of agreement features, such as number and gender, which we find ourselves in need of to confidently write new sentences with the noun phrases in the corpus. For a moment, let us attempt to recover these features given the information that is present in the corpus.

Why are these features absent in the first place? Perhaps it is imagined that number, for example, can be reliably extracted deterministically, particularly considering that number is given at the leaves. An algorithm that immediately comes to mind is to simply use the head of a tree to recover these properties at each internal node.

The first difficulty we encounter with the “head” approach is that there is some disagreement about what exactly the head should be! Conjunction phrases, for example, are notorious sticky points: should the head of the first conjoined phrase be the head of the whole phrase? Perhaps the conjunction itself should be the head? When we consider these options with our problem, the answer becomes even less clear: if we form a phrase from two singular phrases using *and*, certainly the entire phrase should not be singular. Yet, we have some well-founded discomfort with the notion that *and* is plural. (Clearly there is a problem with the representation that is used in the treebank; this is addressed in [19].)

We would run into even more trouble with other languages, like Arabic, where verbs are conjugated based on gender as well. If we merge (via *and*) a masculine noun phrase and a feminine noun phrase, the resulting noun phrase is considered masculine; however, even if we get over our discomfort at assigning *and* noun-like features, there is no correct answer here: no fixed decision on constituent head assignment works for any combination of masculine and feminine subphrases. Other conjunctions, like *or*, are even more difficult; it is not clear whether a phrase like *a man or the boys* should be singular or plural, so we should not be surprised that a general rule is elusive.

Even if we adopt some set of rules, say the hand-written head rules used to parametrize the Collins parser, we run into difficulty with fairly common phrases like *a number of men*, where *number*, a singular noun, is the head. Perhaps the rea-

son that such rules would fail is that number is really a semantic thing with some syntactic manifestations, and not simply an arbitrary syntactic property. If English speakers were to begin using a new word, say *joople*, they may assume that it refers to some singular thing based on its ending and would conjugate verbs accordingly; however, if they were told that *joople* referred to a plural thing, they would at the very least consider treating it as such (analogous phenomena were described in [32]).

Furthermore, when we are considering two such features simultaneously, it is not clear that both properties for an entire noun phrase would even come from a single head. These difficulties pushed the NLP community to use statistical methods in the first place.

1.3.1 Method

Statistical parsers are built using hand annotated parse trees, such as the treebank. Unfortunately, the treebank does not label entire noun phrases as singular or plural, without which we cannot build a supervised model.

However, the treebank does have rich argument labels, so we can extract subject-verb pairs from sentences in the treebank. Since the verbs and their subjects must agree in number in English, wherever we can identify the number of the verb, we can identify the number its subject. Hence, the same agreement property of the language that motivated this problem in the first place is the source of the data we will use to solve it. This method can be used to recover any syntactic property that is needed for agreement in any language. We apply our approach to number in English and gender in Arabic.

1.3.2 Data Collection

In both treebanks, we make use of the extended set of internal node tags; in particular, the NP-SBJ tag helps us handle sentence constructions that do not have ordinary agreement, such as expletive and cleft constructions. In English, we search for the

first verb leaf of the predicate; if it is labeled VBZ, we take the subject to be singular; if it is labeled VBP, we take the subject to be plural; otherwise, we discard the subject. Note that this heuristic will work for sentences with auxilliary verbs as well, since these verbs carry the number agreement, and its verbal arguments do not carry number (e.g., *John (and Bob) eats (eat) the apple* become *John (and Bob) has (have) eaten/been eating/etc. the apple*).

In Arabic, circumstances are slightly more complicated. Formally, Arabic has strict gender agreement. However, there are two major word orders in Arabic with differing number agreement rules. In particular, sentences with the VSO (Verb Subject Object) ordering (جُمْلَةٌ فِعْلِيَّةٌ or verbal sentences) have gender agreement between verb and subject, but do not have number agreement (the verb is usually conjugated for the singular). On the other hand, sentences with the SVO order (جُمْلَةٌ اِسْمِيَّةٌ or nominal sentences) have agreement in both gender and number. In Classical Arabic, the VSO order is dominant; however, in Modern Standard Arabic (MSA) the SVO order is quite common.¹

In Arabic, we collect sentences with the same criterion as in English and search for the verb with the same heuristic (the modal verb certainly must agree with the subject; in Arabic, the verbal arguments will also agree, but this information is not used in our experiments). Also, since Arabic is morphologically rich, the leaf tags take on a more “factored” appearance; on the other hand, the internal nodes are more impoverished (there are only 17 distinct internal node labels, compared with 240 distinct internal node labels for English, but there are 284 distinct leaf labels in Arabic, compared with only 43 for English). The verb tags that are taken as singular/dual/plural and masculine/feminine are given in table 1.4.

sentence type	verb gender	verb number
VSO جُمْلَةٌ فِعْلِيَّةٌ	agreement	singular
SVO جُمْلَةٌ اِسْمِيَّةٌ	agreement	agreement

¹It may seem surprising that Arabic would preserve any verb-subject number agreement considering that the classically dominant order lacks this type of agreement; however, Arabic is pro-drop and sentences without explicit subjects have subject-verb agreement in both gender and number.

VSO — (Verbal Sentence / جُمْلَةٌ فِعْلِيَّةٌ)						
Masculine / مُذَكَّرٌ			Feminine / مُؤَنَّثٌ			
Singular / مُفْرَدٌ	أَلْفَاةٌ <i>the apple</i>	أَوْلَادٌ <i>the boy</i>	أَكَلَ <i>ate(MS)</i>	أَلْفَاةٌ <i>the apple</i>	أَبْنَاتٌ <i>the girl</i>	أَكَلَتْ <i>ate(FS)</i>
Dual / مُتَنِيٌّ	أَلْفَاةٌ <i>the apple</i>	أَوْلَادَانِ <i>the two boys</i>	أَكَلَا <i>ate(MS)</i>	أَلْفَاةٌ <i>the apple</i>	أَبْنَاتَانِ <i>the two girls</i>	أَكَلْتَا <i>ate(FS)</i>
Plural / جَمِيعٌ	أَلْفَاةٌ <i>the apple</i>	أَوْلَادٌ <i>the boys</i>	أَكَلُوا <i>ate(MS)</i>	أَلْفَاةٌ <i>the apple</i>	أَبْنَاتٌ <i>the girls</i>	أَكَلْنَ <i>ate(FS)</i>

SVO — (Nominal Sentence / جُمْلَةٌ إِسْمِيَّةٌ)						
Masculine / مُذَكَّرٌ			Feminine / مُؤَنَّثٌ			
Singular / مُفْرَدٌ	أَلْفَاةٌ <i>the apple</i>	أَكَلَ <i>ate(MS)</i>	أَوْلَادٌ <i>the boy</i>	أَلْفَاةٌ <i>the apple</i>	أَكَلَتْ <i>ate(FS)</i>	أَبْنَاتٌ <i>the girl</i>
Dual / مُتَنِيٌّ	أَلْفَاةٌ <i>the apple</i>	أَكَلَا <i>ate(MD)</i>	أَوْلَادَانِ <i>the two boys</i>	أَلْفَاةٌ <i>the apple</i>	أَكَلْتَا <i>ate(FD)</i>	أَبْنَاتَانِ <i>the two girls</i>
Plural / جَمِيعٌ	أَلْفَاةٌ <i>the apple</i>	أَكَلُوا <i>ate(MP)</i>	أَوْلَادٌ <i>the boys</i>	أَلْفَاةٌ <i>the apple</i>	أَكَلْنَ <i>ate(MP)</i>	أَبْنَاتٌ <i>the girls</i>

Table 1.2: To distinguish between the different conjugated forms in Arabic, we append labels to the English verb in the gloss; M and F denote masculine and feminine, respectively, and S, D, and P denote singular, dual, and plural, respectively.

1.3.3 The Model

Our setting is a standard supervised classification task: the input is a noun phrase tree and the output is the number/gender of the noun phrase. Our model is that the number/gender is propagated upward from the leaves to the root according to feature-specific head rules; the parameters of our model dictate the probability at each stage that the feature-head is a particular child. For each noun phrase, we are given the feature value for the root (by virtue of subject-verb agreement) and for all of the leaves (which are hand labeled). We treat the labels of the remaining nodes as hidden variables and use expectation-maximization (EM) to optimize the model.

With some notation, we can make this more explicit: let T denote a parse tree, T_ϕ the root node of T , T_i the i th subtree of the root, and T_* the root rule (i.e., $T_\phi \rightarrow T_{1\phi} T_{2\phi} \dots T_{N\phi}$). Our goal is to model the propagation of some feature F ; we do so as follows:

$$\Pr(F(T_\phi) = f | T) = \sum_i \Pr(F(T_{i\phi}) = f | T_i) \Pr(F\text{-head} = i | T_*)$$

Hence, the parameters of the model are the probabilities $\Pr(F\text{-head} = i | T_*)$ for each rule in the grammar. The distributions of the feature-values for the leaf nodes is fixed

	Masculine	Feminine
Singular	PVSUFF_SUBJ:3MS	PVSUFF_SUBJ:3FS
	PVSUFF_SUBJ:2MS	PVSUFF_SUBJ:2FS
	IV3MS	IV3FS
	IV2MS	IV2FS
Dual	CVSUFF_SUBJ:2MS	CVSUFF_SUBJ:2FS
	PVSUFF_SUBJ:3MD	PVSUFF_SUBJ:3FD
	IV3MD	IV3FD
Plural	PVSUFF_SUBJ:3MP	PVSUFF_SUBJ:3FP
	PVSUFF_SUBJ:2MP	PVSUFF_SUBJ:2FP
	IV3MP	IV3FP
	IV2MP	IV2FP
	CVSUFF_SUBJ:2MP	CVSUFF_SUBJ:2FP

Table 1.3: The agreement tags of conjugated verbs in the Arabic treebank. PV, IV, and CV are perfect, imperfect, and imperative (command) aspects, respectively.

for each language and is given in table 3.

We optimize the parameters of the model using EM; the resulting algorithm is given in table 1.5.

After the model has been optimized, we can evaluate new trees; in our algorithm, in case of a tie, we back-off to a simple majority model. That is, for example, if singular noun phrases are more common and the probability that a given noun phrase is plural according to our model is exactly 0.5, we will return singular.

We compare our method to a simple majority model in Arabic and to the one suggested in the introduction (drawing number from the head of the noun phrase) in English, using the head rules given in Collins 1999. We trained our models on a random selection of 80% of the noun phrases obtained from each corpus; we ran the algorithm for 10 iterations (we found that increasing the number of iterations had little effect on the performance).

English	
Singular	NN, NNP
Plural	NNPS, NNPS

Arabic		
	Masculine	Feminine
Singular	CVSUFF.SUBJ:2MS	NSUFF.FEM.SG
	DEM.PRON.MS	CVSUFF.SUBJ:2FS
	IVSUFF.DO:3MS	DEM.PRON.FS
	PRON.2MS	IVSUFF.DO:3FS
	PRON.3MS	PRON.3FS
	PVSUFF.DO:3MS	DEM.PRON.F
	IV1S	IV1S
	PRON.1S	PRON.1S
	PVSUFF.DO:1S	PVSUFF.DO:1S
		PVSUFF.DO:3FS
Dual	NSUFF.MASC.DU.ACC	NSUFF.FEM.DU.ACC
	NSUFF.MASC.DU.ACC.POSS	NSUFF.FEM.DU.ACC.POSS
	NSUFF.MASC.DU.GEN	NSUFF.FEM.DU.GEN
	NSUFF.MASC.DU.GEN.POSS	NSUFF.FEM.DU.GEN.POSS
	NSUFF.MASC.DU.NOM	NSUFF.FEM.DU.NOM
	NSUFF.MASC.DU.NOM.POSS	NSUFF.FEM.DU.NOM.POSS
	CVSUFF.SUBJ:2MP	NSUFF.FEM.DU.NOM.POSS
	DEM.PRON.MP	DEM.PRON.FD
	IV1P	IV1P
	PRON.1P	PRON.1P
	PRON.2MP	DEM.PRON.F
	IV2MP	PRON.3D
	IVSUFF.DO:3MP	PVSUFF.DO:1P
	PRON.3D	PVSUFF.DO:3D
PVSUFF.DO:1P		
PVSUFF.DO:3D		
Plural	NSUFF.MASC.PL.ACC	
	NSUFF.MASC.PL.ACC.POSS	
	NSUFF.MASC.PL.GEN	
	NSUFF.MASC.PL.GEN.POSS	
	NSUFF.MASC.PL.NOM	
	NSUFF.MASC.PL.NOM.POSS	NSUFF.FEM.PL
	PRON.3MP	IV1P
	CVSUFF.SUBJ:2MP	PRON.1P
	DEM.PRON.MP	DEM.PRON.F
	IV1P	PVSUFF.DO:1P
	PRON.1P	
	PRON.2MP	
	IV2MP	
	IVSUFF.DO:3MP	
PVSUFF.DO:1P		
PVSUFF.DO:3MP		

Table 1.4: Leaf feature values in Arabic and English. Tags not given here are evenly distributed between all possible labels. Tags placed in multiple categories are evenly distributed between those categories. Also, note that we have included verbal categories for Arabic; this is in anticipation of future work with statistically inducing these features on verb phrases.

Given:

- ℓ possible feature values
- An array of N trees, T
- An array of feature values, F_T , for the roots
- An array of feature values, f , for the leaves

Initialize Θ :

- **foreach** rule T_* :
 - $n \leftarrow$ the number of children in rule T_*
 - $\Theta[T_*] \leftarrow n$ random values from $[0, 1]$
 - normalize $\Theta[T_*]$

do

- *E-Step:*
 - **foreach** tree $T[i]$:
 - $\text{estep}(T[i], \Theta)$
 - $T[i].F \leftarrow F_T[i]$
- *M-Step:*
 - zero Θ
 - **foreach** tree $T[i]$:
 - $\text{mstep}(T[i], \Theta)$
 - normalize Θ

until convergence

function $\text{estep}(\text{tree } \tau, \text{parameters } \Theta)$

- **if** τ is a leaf:
 - $\tau.F \leftarrow f(\tau.\text{tag})$
- **else:**
 - $p \leftarrow$ array of ℓ zeros
 - **foreach** subtree τ_i :
 - $\text{estep}(\tau_i, \Theta)$
 - $p \leftarrow p + \tau_{i\phi}.F * \Theta[\tau_*][i]$
 - $\tau.F \leftarrow p$

function $\text{mstep}(\text{tree } \tau, \text{parameters } \Theta)$

- **if** τ is not a leaf:
 - **foreach** subtree τ_i :
 - $\text{mstep}(\tau_i, \Theta)$
 - $\Theta[\tau_*][i] \leftarrow \Theta[\tau_*][i] + \tau.F \cdot \tau_i.F$

Table 1.5: The model optimization algorithm.

1.4 Results

Let’s start by looking at the results of the tree-augmenter. We performed experiments on the English and Arabic Penn Treebank corpora. For English, we start with 35,752 sentences with a total of 910K tokens which yielded 28K subject noun phrases for which numbers could be extracted. For Arabic, we start with 12,412 sentences with a total of 430K tokens which yielded 24K subject noun phrases for which gender could be extracted. Unfortunately, the SVO order is relatively rare in the Arabic corpus; only 30 noun phrases could be marked for number. Hence the results for Arabic number were not statistically significant.

Language	Feature	Baseline	Head	EM
Arabic	Gender	66%	–	69%
English	Number	61%	75%	82%

We ran our transformation experiments on the standard evaluation section (section 0) of the Penn Treebank II (the Wall Street Journal corpus). We ran the modified Collins parser on this data. We ran the active-to-passive and cleft tree-transducers on this data to produce questionnaires for each sentence; 14% of the sentences processed produced multiple outputs; another 4% produced a single output. We gave the questionnaires for 300 randomly selected sentences to two users for labeling. For each user, the scores for the parses were adjusted for the labels (parses corresponding to images that were marked as incorrect were simply penalized by 10 log probability points). The resulting top-scoring parses were compared to the gold-standard parses using the evalb metric described in [16]; the results are as follows:

Baseline (Collins Parser)	0.871	
Collins + Oracle	0.893	+17.0%
Collins + User 1	0.877	+4.7%
Collins + User 2	0.874	+2.2%
Collins + Users 1 and 2	0.875	+2.4%

User 2 reported that he labelled the data very quickly and may have made some mistakes, possibly explaining his diminished performance. Below are the agreement statistics between the two users.

		User 1		Total
		Correct	Incorrect	
User 2	Correct	65%	5%	70%
	Incorrect	12%	18%	30%
Total		77%	23%	100%

We expect that the reason for the high ($\sim 20\%$) false image rate is due to the rough implementation of the transformations. Despite this, the method shows significant improvement to a baseline parser with relatively little effort on the part of the user (one or two questions per sentence). We expect that better statistics would be achieved once the transformation was improved to incorporate lexical information. The labels from both users correspond to a significant improvement in ambiguous, semantically meaningful attachment types (particularly PP).

1.5 Future Work

The first and most obvious step is to further refine the transduction grammar by hand. Improvements could be made simply by modifying the parser to output traces. Several transformations can be added to help other attachment problems, and also to attempt anaphora resolution and other language understanding tasks.

Naturally, more data would improve our performance and allow us to advance our model. In particular, an attractive next step is to obtain enough data to lexicalize the rule scores, so that we can estimate when a rule will fail. An even more interesting idea (and therefore one which requires yet more data) would be to predict correlated failures; that is, we may find that the sentence will break when one constituent is placed before the other for some linguistic reason (e.g., one refers to the other), so all transformations that rearrange these constituents will fail.

Beyond a specialized corpus for this task, we are interested in using this system to develop treebanks in new domains and languages and to refine existing treebanks, as in [13].

In lieu of more labeled data, we can make use of unlabeled (i.e., parsed, but not annotated) data to estimate the lexicalized rule scores via self-training. For instance, if we find that a verb does not appear in passive voice in our data set, we can confidently hypothesize that it cannot be passivized.

The Transduction Grammars

Bold-faced capital letters denote internal nodes of the transduction grammar, normal capital letters denote internal nodes of the parse trees (often accompanied with regular expression markers), the asterisk denotes any number of children of any type, subscripts are indices for output rules that return multiple items (tuples), † and ‡ when present denote the arguments to output rules, and *c* and *pp* are auxiliary functions which output the copula and the passive participle, respectively. Patterns that appear in the output rules are intended to refer to the region in the input that was matched; regions matched by the asterisks are intended to be replaced in the order they appeared in the input.



S	←	S(-A)?	→	* NPB?-A *	V	*
	⇒	"	→	* V ₂ (NPB?-A) * V ₁ (NPB?-A) *		
V	←	VP	→	VB * VBA		*
	⇒		→ ₁ (†)	VB * VBA ₁ (†, VB)		*
			→ ₂ (†)	VBA ₂ (†, VB)		
V	←	VP	→	VB	* NPB?-A *	
	⇒		→ ₁ (†)	<i>c</i> (†) <i>pp</i> (VB) * by †		*
			→ ₂ (†)	NPB?-A		
VBA	←	VP-A	→	VB * VBA		*
	⇒		→ ₁ (†, ‡)	VB * VBA ₁ (†, VB)		*
			→ ₂ (†, ‡)	VBA ₂ (†, VB)		
VBA	←	VP-A	→	VB	* NPB?-A *	
	⇒		→ ₁ (†, ‡)	<i>c</i> (‡, NPB?-A) <i>p</i> (VB) * by †		*
			→ ₂ (†, ‡)	NPB?-A		
VB	←	VB[^N]	→	<i>not a form of to be or to have</i>		
	⇒		→	itself		

Table 1.6: The active-to-passive transduction grammar

CS	←	S(-A)?	→	(CC INTJ)* (PREARG)* CSUB * CVB *
	⇒		→	It <i>c</i> (CVB) CSUB that * CVB *
CS	←	S(-A)?	→	(CC INTJ)* (PREARG)* CSUB * CVBARG *
	⇒		→	It <i>c</i> (CVBARG) CSUB that * CVBARG *
CVBARG	←	VP	→	* CSUB *
	⇒		→ ₁	CSUB
			→ ₂	* *
CVBARG	←	VP	→	* CVBARG *
	⇒		→ ₁	CVBARG ₁
			→ ₂	* CVBARG ₂ *
CVB	←	VP	→	<i>head is not a form of to, said, add, contend</i>
	⇒		→	itself
PREARG	←	!(S.* CC INTJ)	→	*
	⇒		→	itself
CSUB	←	[^{VCMSI}].*	→	*
	⇒		→	itself

Table 1.7: The cleft transduction grammar

Chapter 2

Human Assisted Word-Sense Disambiguation for Translation

The goal of word-sense disambiguation systems is to identify the meaning of a word from context. Word-sense disambiguation, like parsing, is considered a fundamental problem in natural language processing, which is broadly applicable and, like parsing, word-sense disambiguation systems are rarely used outside of natural language research labs. In spite of this, word-sense disambiguation remains an active area of research.

Formally, for a given word or phrase in context, word-sense disambiguators are asked to identify which *synset* the word or phrase belongs to. The difficulty depends on the granularity of distinctions that are being made; in a recent evaluation which used synsets from WordNet, interannotator agreement on word senses was as low as 85%: this is a practical upper bound on an automatic systems performance [1].

Shallow methods yield surprisingly good results on word-sense disambiguation; the simplest method that achieves state-of-the-art performance is one which incorporates a handful of shallow methods by voting.

One of the most-cited potential uses of word-sense disambiguation is automatic machine translation. In [11], Carpuat and Wu incorporated a state-of-the-art Chinese

word sense disambiguation model into a state-of-the-art statistical machine translation system and found that their system performed worse than an uninformed system! This was (to the best of our knowledge) the only negative result published in ACL that year, which suggests just how unexpected the result is. Indeed, several later papers (notably [12] and [14]) seek to reverse it, with some success. Nevertheless, the fact that the original paper found that it did not help and that the later papers (despite careful design surrounding the inclusion of the system) found limited improvement is suggestive.

We will modify a current state-of-the-art system to obtain word-senses from a monolingual consultant instead of from a word-sense disambiguation system to more directly evaluate whether eliminating this type of ambiguity improves translation. Notably, our system is the first that uses source-language consultants to improve a statistical translation system.

2.1 Previous Work

The idea of disambiguating by asking is an old one; perhaps the earliest paper to suggest disambiguation by asking is Kay’s description of the MIND system [26], a multilingual translation engine. Kay points out that use of the term “fully automatic” when describing translation systems is misleading, because users of fully automatic translation systems will edit the output if they are familiar with the target language. He therefore suggests the use of a monolingual consultant to resolve ambiguities in the source language, but he does not propose algorithms for producing or processing the interaction, nor does he indicate what form the interaction should take.

Tomita [46] later refines the problem definition by disallowing assumptions about the consultant—in particular, that the consultant has any specialized background (in linguistics or computer science)—so the interaction cannot include parse trees or phrase structure rules. Tomita presents a rule-based system for parsing by augmenting a context-free grammar with annotation rules to explain each attachment and solicit

correction. For the sentence:

Mary saw a man in the park with the telescope.

Tomita's system generates the following questions:

- 1) (a man) is (in the park)
- 2) The action (Mary saw a man) takes place (in the park)

- 1) (the park) is (with a telescope)
- 2) (a man) is (with a telescope)
- 3) The action (Mary saw a man) takes place (with a telescope)

Ben-Ari *et al* propose embedding a similarly designed parsing engine in a rule-based translation system [2], with an emphasis on transfer (i.e., preferring to preserve ambiguity in translation whenever the source and target languages make that possible). They also point out that more sophisticated queries than those of Tomita may be used for certain types of ambiguity.

Maruyama [35] presents a system for ambiguity resolution in a Japanese-to-English translation system by interactively displaying dependencies. When the user selects a phrase (underlined), the system displays the phrases that the chosen phrase might modify (the automatic choice is in reverse video; all other choices are highlighted in red). In the screenshots below (adapted from [35]), the user first selects the first phrase, then the second phrase.



A recent paper [10] presents an elegant variation of the target language rewriting technique: display a chart of possible phrasal translations and permit the user to select a path through the phrase chart, but which reports negative results in spite of

the human assistance.

2.2 Methods

`moses`[30] is a free (both *gratis* and *libre*) state-of-the-art statistical translation system. The `moses` decoder can be asked to produce the set of all possible outputs (in its search space) for a given input as a directed acyclic graph. Each edge corresponds to a phrase in the output, and each path from the start node to an end node corresponds to a candidate translation.

Formally, `moses` generates a collection of nodes \mathcal{V} and edges \mathcal{E} . Each node represents a possible state of the decoder (as described in the introduction), and each edge consists of a source phrase that is removed from the source, a target phrase that is generated and added to the candidate translation, and a score that is added to the total score at that point. Each path from the empty start node to a final node (one for which the entire source sentence has been absorbed) corresponds to a candidate translation.

$$\mathcal{V} \subseteq \bigcup_{\ell \in \mathbb{N} \cup \{0\}} E^\ell \times \mathcal{P}(\{1, \dots, m\}) \times [0, 1]$$

where E denotes the set of all words in English and m is the length of the input sentence. Members of \mathcal{V} are triples:

1. an English hypothesis prefix,
2. the indices of the French words that have already been translated, and
3. the score so far.

$$\mathcal{E} \subseteq \bigcup_{\ell \in \mathbb{N} \cup \{0\}} E^\ell \times \bigcup_{m \in \mathbb{N} \cup \{0\}} F^m \times [0, 1]$$

Members of \mathcal{E} are triples:

1. an English hypothesis fragment,
2. a French hypothesis fragment, and
3. the partial score.

If we had access to a bilingual consultant, we could ask them whether or not a particular edge in the graph would be visited in an ideal translation. Which nodes in a given graph would we ask about? As always, we would ask the consultant the question that would give us the most information. Assuming (for simplicity) that there are no redundant paths in the search graph (i.e., that every possible translation corresponds to at most one path), the edge that would give us the most information is the one which we are most uncertain about; that is, we are interested in edges that we will traverse with as close to a 50% probability as possible (or more generally, the edges that have the highest *entropy*).

for each node n **in** G (traverse in topological order):

$$\cdot \quad P[n] := \sum_{n'} P[n'] P[n' \rightarrow n]$$

Asking about particular edges is effectively asking a bilingual consultant whether or not a phrase in the target language is an appropriate translation of a particular phrase in the source sentence. Unfortunately, we only have access to consultants that speak the source language; hence, instead of asking questions involving both languages, we can only ask questions involving the source language. Our model provides translation tables of the form $\Pr(f | e)$, so given a particular phrase in the target language, we can produce a number of candidate phrases in the source language (including the original source phrase) as well as scores for each of these phrases. Notice that the flipped translation model that came out of the noisy channel approach is in exactly the right form for our use!

Now our questions are of the form: which of the following two phrases is the more suitable synonym of this phrase in the original sentence? How much information is

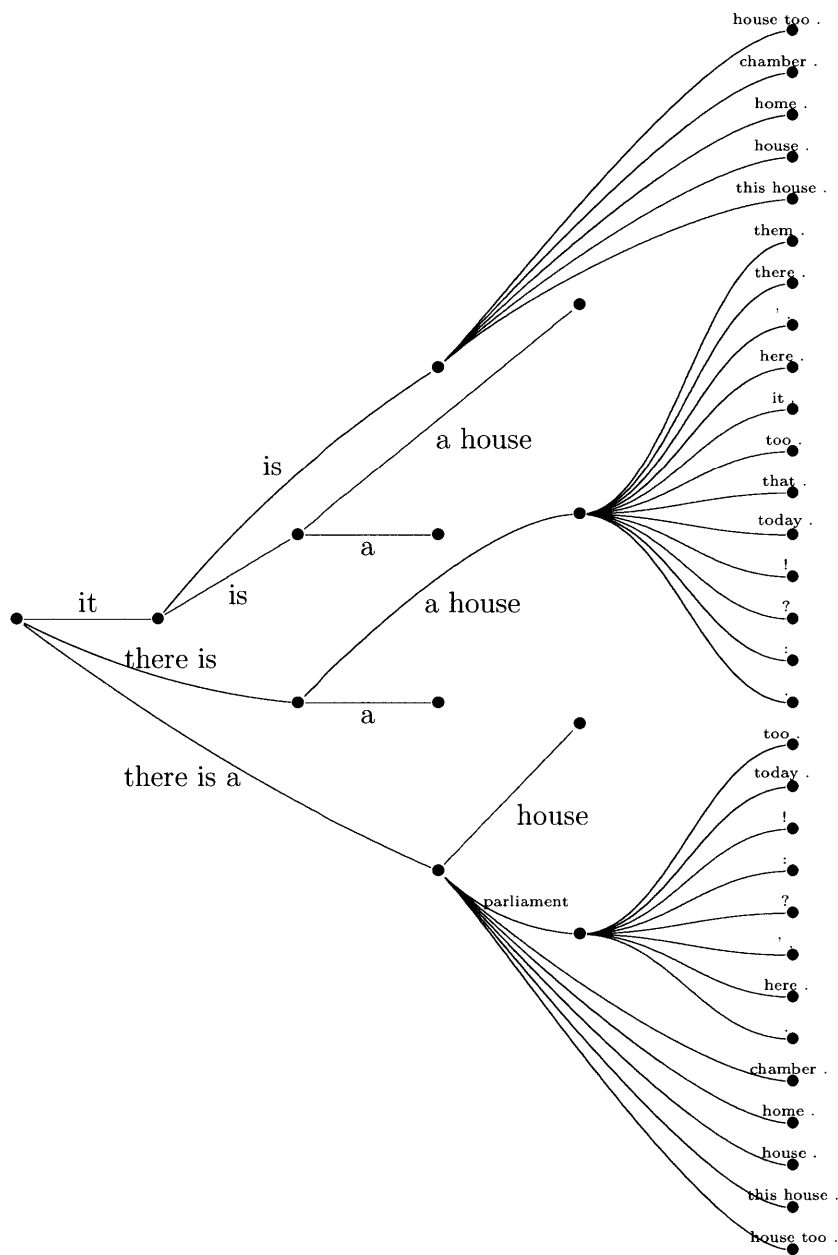


Figure 2-1: A search graph produced by *moses* for the German sentence *Es gibt ein Haus*.

gained by answering such a question?

$$\begin{aligned}
\Pr(\text{translation} \mid \text{answer}) &= \prod \Pr(\text{node} \mid \text{answer}) \\
&= \prod \Pr(\text{answer} \mid \text{node}) \cdot \frac{\Pr(\text{node})}{\Pr(\text{answer})} \\
&= \prod \Pr(f_{\text{answer}} \mid e_{\text{node}}) \cdot \frac{\Pr(\text{node})}{\Pr(\text{answer})}.
\end{aligned}$$

Without knowing the answer, we can tell what information will be gained in expectation by asking this question. We simply ask the question that will maximize this quantity.

In practice, this idea is complicated by the fact that $\Pr(f \mid e)$ is not given and is very often overtrained. `moses` will, by default, learn five models for $\Pr(f \mid e)$; the final model is optimized with a held-out subset of the corpus as an affine combination of the five models. These other models can be used to estimate the model’s uncertainty without going back to training. We simply assume that the underlying distribution of $\Pr(f \mid e)$ is drawn uniformly from the simplex of all affine combinations of the five models.

We use a Monte Carlo algorithm to estimate the entropy of the resulting model.

$$\begin{aligned}
&\arg \max_q E_{\text{ans}} [H[\mathbf{e} \mid \text{ans}] \mid q] \\
&= \arg \max_q \sum_{\text{ans}} \Pr(\text{ans}) H[\mathbf{e} \mid \text{ans}] \\
&= \arg \max_q \sum_{\text{ans}} \sum_{\mathbf{e}} \Pr(\text{ans} \mid \mathbf{e}) \Pr(\mathbf{e}) \log \left(\frac{\Pr(\text{ans} \mid \mathbf{e}) \Pr(\mathbf{e})}{\Pr(\text{ans})} \right) \\
&= \arg \max_q \sum_{\text{ans}} \sum_{\mathbf{e}} \Pr(\mathbf{e}) \Pr(\text{ans} \mid \mathbf{e}) \log(\Pr(\text{ans} \mid \mathbf{e})) \\
&\quad - \sum_{\text{ans}} \Pr(\text{ans}) \log \Pr(\text{ans})
\end{aligned}$$

To incorporate uncertainty in the values of $\Pr(\mathbf{e})$ into the model, we will sample from the space of possible models which `moses`’s training script generates. The

complete algorithm is in figure 2-2.

2.3 Mechanical Turk

The source for the human judgments needed by our algorithm is Mechanical Turk. Mechanical Turk is a service offered by Amazon for outsourcing tasks to be completed by human agents, and for doing so at a large scale. Amazon calls these tasks “Human Intelligence Tasks” (HITs), corresponding to our use-case for crowdsourcing AI complete tasks. We offer \$.10 per task. The tasks are easy to complete so this corresponds to a reasonable hourly wage. Unfortunately, a fair percentage of Mechanical Turk Workers simply spam the questions, so it is important to include a test question with a known answer. An example question including the test question is shown in Figure 2.3.

2.4 Results

We performed experiments on the same corpus used for the 2007 NIST evaluation. Our results show a significant improvement on the sentences for which the decisions had an effect. However, the system only had an impact on a small portion ($\sim 4\%$) of sentences; hence, the overall improvement is slight (from 27.7 to 28.1) and is just below the threshold of statistical significance.



```

for each phrase  $f_p$  in  $\mathbf{f}$ :
·   for each pair of english translations  $e_1$  and  $e_2$  of  $f_p$ :
·   ·    $Z' = 0$ 
·   ·   for node  $n$  in  $\mathcal{G}$  corresponding to  $f_p$ :
·   ·   ·    $Z' = Z' + n.p$ 
·   ·    $s = 0$ 
·   ·   repeat  $N$  times:
·   ·   ·   draw  $\vec{\alpha}$  from the five-dimensional simplex
·   ·   ·    $p'_1 = 0$ 
·   ·   ·    $p'_2 = 0$ 
·   ·   ·   for node  $n$  in  $\mathcal{G}$  corresponding to  $f_p$ :
·   ·   ·   ·    $p = n.p/Z'$ 
·   ·   ·   ·    $p_1 = \sum_i \alpha_i \Pr(e_1 | f_p)$ 
·   ·   ·   ·    $p_2 = \sum_i \alpha_i \Pr(e_2 | f_p)$ 
·   ·   ·   ·    $Z = p_1 + p_2$ 
·   ·   ·   ·    $p_1 = p_1/Z$ 
·   ·   ·   ·    $p_2 = p_2/Z$ 
·   ·   ·   ·    $s = s + Z'p_1 \log(-\log AnsGi) + Z'p_1 \log(-\log AnsGi)$ 
·   ·   ·   ·    $p'_1 = p'_1 + p_1$ 
·   ·   ·   ·    $p'_2 = p'_2 + p_2$ 
·   ·   ·    $H = Z_1 \log Z_1 + Z_2 \log Z_2$ 
·   ·   ·    $S = S + s - H$ 
·   ·    $S = S/N$ 

```

Figure 2-2: Algorithm for incorporating human assistance in phrase-based translation system.

Welche Wörter kann man in diesem Satz benutzen?

Wer hat _____ das neue Haus verzichtet?

- auf
- vor
- an
- für

Welche Ausdruck ist der beste Ersatz für die Kursivschrift Phrase?

Das Problem ist zu *schwer* für mich .

- schwierig
- dickleibig

Figure 2-3: Sample Mechanical Turk Questionnaire including a test question.

Chapter 3

Human Assisted Translation Evaluation

We return to the lessons of machine translation (in response to Kay’s question) that we touched on in the introduction. The lessons are established as the modern conventional wisdom of the statistical machine translation practitioner and are intended as a review of the major results of the field since Kay first posed this question (roughly the past twenty years); after the review, we will discuss the methodology that was used to obtain these results.

3.1 Answering Kay

3.1.1 Rule-Based MT is Bad

Noisy-channel methods, those described in the introduction, are vastly superior to rule-based systems, and require substantially less effort to design. In our own experiments in 2006 with German to English translation on the EUROPARL corpus [29], Systran (a commercial rule-based translation system) achieved a BLEU score of 0.11, whereas Pharaoh [28] (the predecessor of `moses`[30]) scored a BLEU score of 0.20, a

massive difference.

3.1.2 Hybrid Systems Are OK

In the same series of experiments in 2006, we created a hybrid model where a statistical system was trained on Systran’s output of 100,000 EUROPARL sentences. The resulting BLEU score was 0.20, suggesting that statistical systems are able to fix the errors of the rule-based system. This result was reproduced by [17] and is the basis of Systran’s current hybrid systems.

3.1.3 Minimum Error-Rate Training

In [37], Och suggested adding a number of high-level parameters to statistical language models and optimizing the BLEU score over these parameters as a meta-step in training a translation system. For example, he proposed optimizing:

$$\arg \max_{\mathbf{e}} \Pr(\mathbf{e} | \mathbf{f}) \Pr(\mathbf{e})^\alpha$$

to generalize the original:

$$\arg \max_{\mathbf{e}} \Pr(\mathbf{e} | \mathbf{f}) \Pr(\mathbf{e}).$$

This has nothing to do with the Bayesian formalism that we started with for any value of α other than 1; however, minimum error-rate training typically finds $\alpha \in [2, 2.5]$ to be optimal.

3.1.4 Linguistics Doesn’t Work

It is appealing to rest an engineering effort on a scientific one; the scientific counterpart of natural language processing is linguistics (and the closely related field of psycholinguistics). A number of intricate linguistic models have been in development for decades. These formalisms can be incorporated into statistical translation models,

Average information content per word	BLEU score
3.5	0.28
5.5	0.52
7.6	0.82
10.0	1.00

Table 3.1: Amount of information content per word in nits (according to a trigram model) needed to obtain a certain BLEU score on the EUROPARL corpus. A BLEU score of 0.28 is state-of-the-art on this corpus, and corresponds to the translation model providing 3.5 nits of information per word, compared to the 10.0 nits required to obtain a perfect BLEU score.

but incorporating anything more than the most coarse, general observations (such as the idea that sentences have hierarchical structure) into the design calls for substantial effort which are not rewarded by increased BLEU scores. Hence, in spite of their appeal, formal linguistics does not often find its way into modern machine translation.

3.1.5 *n*-gram Language Models Are Best

In particular, the idea that language is a simple Markov chain is frustrating; [6] describes this model as “almost moronic... [capturing] local tactic constraints by sheer force of numbers, but the more well-protected bastions of semantic, pragmatic and discourse constraint and even morphological and global syntactic constraint remain unscathed, in fact unnoticed”. In spite of this, *n*-gram language models are used universally in speech recognition and in machine translation. In fact, *n*-gram models are easy to train and hard to beat[23, 44].

In our own experiments, we found that it is possible to remove a substantial amount of information-heavy content (as judged by a trigram model) and to obtain a high BLEU score, showing that *n*-gram models are closely tied to the BLEU score (see Table 3.1).

Perhaps it is because data remains sparse and more complex language models would be successful if only more data were available. On the other hand, more (albeit unlabeled) data is readily available for training English language models than

for almost any other artificial intelligence task.

3.1.6 More n -grams Help... Forever

No group is better placed to take advantage of this data than Google, and the Google translation systems are all based on massive language models harvested from websites and newspapers, and have achieved record performance on open-domain machine translation [5]. Beginning with a state-of-the-art Arabic to English translation system, they found that simply increasing the amount of data available to their n -gram model, even data from sources other than the original target text, consistently increases the BLEU score. This is true in spite of simplifications that had to be made to the back-off model to deal with such massive datasets. They report in their conclusions:

Significantly, we found that translation quality as indicated by BLEU continues to improve with increasing language model size, at even the largest sizes considered. This finding underscores the value of being able to train and apply very large language models, and suggests that further performance gains may be had by pursuing this direction further.

Their results are shown in figure 3-1.

Google provides researchers with an excellent opportunity for exploration: they offer the use of their translation system *gratis* on their website¹. The effect that the emphasis on the language model (based on Och's work, described above) combined with the largest language model in the world offers some interesting examples in table 3.2 to reflect on; in order to optimize the language model term $\Pr(\mathbf{e})$, the system will sacrifice $\Pr(\mathbf{f}|\mathbf{e})$.

¹<http://translate.google.com/>

fr → en	La pomme mange le garçon.	
⇒	The apple eats the boy.	Reference
⇒	The apple eats the boy.	Systran
⇒	The boy eats the apple.	Google
fr → en	Un hippopotame me veut pour Noël.	
⇒	A hippopotamus wants me for Christmas.	Reference
⇒	An hippopotamus wants me for Christmas.	Systran
⇒	I want a hippopotamus for Christmas.	Google
de → en	Leute stahlen mein Weißes Auto.	
⇒	People stole my white car.	Reference
⇒	People stole my white car.	Systran
⇒	White people stole my car.	Google
it → en	George Bush non è un idiota.	
⇒	George Bush is not an idiot.	Reference
⇒	George Bush is not an idiot.	Systran
⇒	George Bush is an idiot.	Google
it → en	Ali Mohammad non è un idiota.	
⇒	Ali Mohammad is not an idiot.	Reference
⇒	Ali Mohammad is not an idiot.	Systran
⇒	Mohammad Ali is not an idiot.	Google

Table 3.2: Examples from Systran and Google translate (at the time of this writing), elucidating some possible issues with BLEU.

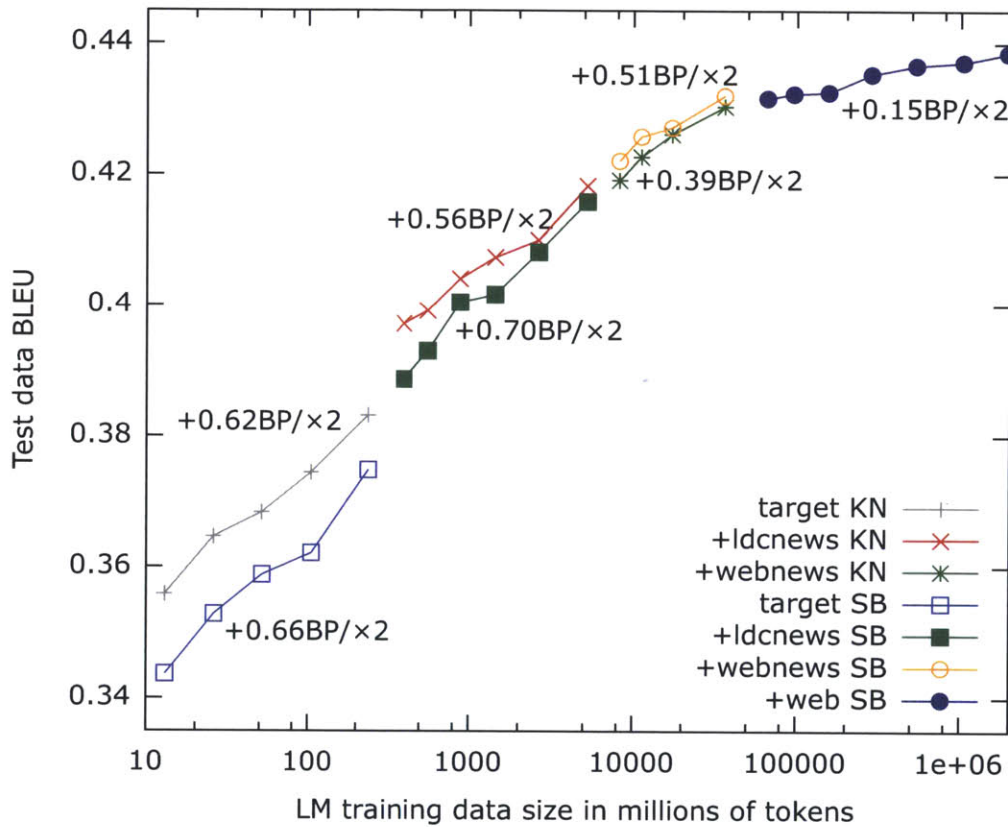


Figure 3-1: BLEU scores for a translation system with varying amounts of data using Kneser-Ney (KN) and Stupid Backoff (SB). BP/ $\times 2$ indicates the increase in the BLEU score (in 0.01 BLEU points, or BP) for every doubling of the size of the language model (in tokens). Differences greater than 0.51 BP are significant at the 0.05 level. Adapted from [5].

3.1.7 Shrinking the Translation Model is OK

Adding complexity to the language model seems to help quite a bit; equivalently, reducing the complexity of the language model will significantly degrade the quality of the translation system. This does not seem to be the case for the translation model. In fact two experiments on state-of-the-art systems causing a massive reduction in the number of parameters of the translation model have shown no statistically significant effect on the BLEU score (one is our own experiment, described in Appendix D; another is the work of [25]).

3.1.8 Word-Sense Disambiguation Does Not Help Machine Translation

The previous chapter details an attempt by Carpuat and Wu and our own attempt to incorporate a Word-Sense Disambiguation model into machine translation, with disappointing results: even with the assistance of monolingual consultants to perform disambiguation, the scores do not improve significantly. Fully automatic word-sense disambiguation systems attenuate the scores.

3.1.9 Having Monolingual Consultants Revise Translations is Bad

[8] puts forward Kay's idea [26] that monolingual consultants are able to improve the output of a modern translation system by introducing a controlled revision phase. Unfortunately, the results are disappointingly poor: the BLEU score plummeted from 0.51 to 0.43.

3.2 BLEU

Under the force of these results, many of which are surprising but not unbelievable, let's take a moment to understand the methodology that yielded them.

Evaluation is a serious problem for Machine Translation. In other supervised learning tasks, such as hand-written digit recognition, evaluation is a simple matter of comparing the output of a system to the output of a reference (usually a human capable of performing the given task). For hand-written digit recognition, a perfect match is given a score of 1, and a mismatch is given a score of 0: guessing "7" when the reference is "8" is just as bad as guessing "0", even though "7" is closer.

In automatic speech recognition, researchers can't afford to be quite as heavy handed since an all-or-nothing approach is too coarse and (for many applications) a hypothesis which has an error in one word is substantially better than random text. The typical method of evaluation in this context is Word Error Rate (WER): it is the average number of edit operations (substitutions, deletions, and insertions) that would have to be performed per word to change the output to the reference. WER is closely related to the Levenshtein distance between two strings; fast (linear-time) algorithms exist for computing WER.

WER can be used to evaluate machine translation as well; however, there are some important differences between what is acceptable for a machine translation system and what is acceptable for an automatic speech recognition system that are not captured by WER: in particular, for automatic speech recognition, a single reference is seen as authoritative, whereas in machine translation, two translators are unlikely to ever agree completely.

Indeed, human evaluations consider two aspects of translation quality separately: fluency and adequacy. Fluency is a measure of the quality of text in the target language, whereas adequacy is a measure of the accuracy of the translation. Fluency and adequacy are competing qualities; a translation which is more faithful to the source text, particularly for languages that are unrelated, tends to feel alien in the

target language. For example, the Arabic sentence:

عليّ موجودٌ

could be translated as “Ali is present” or as “Ali is a findable one.” Technically, “a findable one” is a more precise rendering of

موجودٌ

than “present”, but clearly “Ali is present” is a much more fluid English sentence.

The BLEU metric[40] is a standard method for evaluating machine translation system performance by comparing translations to one or many human translations. The translations are compared by the precision of n -grams of successively greater length; the BLEU score typically refers to a smoothed 4-gram comparison; mathematically, it can be described by the following formula:

$$\begin{aligned} \text{BLEU} &= e^{I_{c \leq r} \cdot (1-r/c)} \cdot \sqrt[4]{p_1 p_2 p_3 p_4} \\ \log \text{BLEU} &= I_{c \leq r} \left(1 - \frac{r}{c}\right) + \sum_{i=1}^4 \frac{1}{4} p_i \end{aligned}$$

where c is the total length of the candidate translation produced by the system being evaluated, r is sum of the lengths of the reference sentences that most closely match the lengths of the candidate sentences, $I_{c \leq r}$ is 1 if $c \leq r$ and 0 otherwise, and p_j refers to the j -gram precision of the test set.

Evaluating translation systems is a difficult task—BLEU has a number of useful properties that make it a popular choice: it is fast, it is cheap, and it correlates well to judgments of translation fluency and adequacy by bilingual judges. Hence, BLEU is the de facto standard for evaluating automatic machine translation systems and has been for over a decade. It is largely on the authority of BLEU that the conclusions of the previous section rest; over 90% of papers in machine translation published after BLEU use BLEU as their sole method of evaluation. BLEU has had an impressive impact on the agenda of the field.

3.2.1 BLEU Correlates Well With Human Judgments?

BLEU is used as a general metric for translation quality, but considering its definition, it is more precise to call it an n -gram precision model. The papers that introduce BLEU [39, 40] report then that an n -gram precision model correlates strongly with human judgments of translation *quality*. Of course, an improvement in a correlated value does not imply an improvement in the value of interest; however, it is the strength of the correlation that is so promising. In [39], the authors report that the correlation between BLEU and human judgments of adequacy and fluency for French-English translation systems are **0.94** and **1.00**, respectively.

These figures are incredibly impressive. A correlation of 1.0 implies that BLEU is a linear function of the human evaluation of fluency, which implies that BLEU can be used to predict the human evaluation of fluency *perfectly*.

Another interpretation of this result, however, is that humans, when asked to evaluate the fluency of a sentence *against reference sentences*, simply report values that are proportional to the n -gram precision of the sentence versus the references. This interpretation is a statement on the methods that were used to collect the human judgments; that is, under certain conditions, when humans are asked to perform semantic or syntactic comparisons, they will resort to simple string comparisons, and these papers established a set of conditions under which this is true.

It is alarming that this seemingly minor issue of methodology calls into question the research agenda that mainstream researchers have been laboring under for decades.

3.2.2 Re-evaluating BLEU

When we restate the experimental results we cited above as statements about the n -gram precision, the results are far less surprising and the conclusions are less convincing.

The problems with BLEU are becoming better known; Callison-Burch *et al* followed the surprising results of their experiment with human evaluations that showed lower correlations with BLEU than were reported by Papineni *et al* [10]. However, they conclude that, although BLEU clearly cannot be used as a substitute for human evaluations in general, it can be used to compare models within the same class.

Another common belief is that even if BLEU has its problems, it may be possible to account for them with some other automatic metric. It seems that an evaluation method for translation is hardly useful if it cannot be used to compare two arbitrary systems; any performance difference can simply be attributed to the systems being overly different. I propose that any automatic metric can be fooled—building an automatic evaluator is just as hard as building an automatic translation system. Indeed, recent experiments show that, for most language pairs, BLEU correlates more strongly to human judgments than any other available automatic metric does [9].

Reconsidering the results of table 3.1, BLEU is effectively a detector for an n -gram language model. Hence, we propose the use of monolingual consultants for evaluation. With humans in the loop, we have a good chance of mimicking a bilingual evaluator.

According to earlier work [39], monolingual evaluators favor translations that are more fluid to translations that are more adequate. On the other hand, bilingual evaluators tend to be more forgiving to sentences that favor adequacy over fluency. It seems that the gold standard should be the bilingual evaluator; an evaluator that is able to judge the source text as well as the target text will obviously have an advantage. At the same time, there is some reason to pay attention to the monolingual evaluator’s perspective, as it is generally people that have little to no familiarity with the source language that will be using a translation system.

3.3 Methods

The psychometrics literature tells us that when faced with difficult judgments, human evaluators experience cognitive fatigue and will turn to heuristics to perform the evaluation. This fatigue can be measured (and lessened) by asking questions which can be answered quickly and consistently [18]. Forced-choice binary comparisons are a premiere method of obtaining information from human evaluators without causing the kind of fatigue that we described.

We will withdraw from the machine translation setting for the remainder of this section to analyze the mathematical aspect of this question. This yields a new evaluation metric which we call CLAIRES².

You are the head judge of a baking competition and you are required to announce a full ranking of the cakes that were submitted to you. You are democratic, so you'd like to give a ranking that corresponds to the rankings that would be given by the average cake-taster. Tasters are able to compare exactly two cakes at a time, and will indicate which of the two cakes is better. With an unlimited queue of tasters at your disposal, how can you arrive at the correct ranking of these cakes?

If the tasters were absolutely consistent, this would be a sorting problem, and it could be solved in $N \log N$ comparisons in the worst case. However, we can't guarantee that the tasters are in total agreement. Let us instead model the tasting as a probabilistic process.

Let's assume that each of the cakes $i = 1, 2, \dots, N$ has a secret score $x_i \in \mathbb{R}$, such that when we ask a taster to try cakes i and j , he draws a value X which represents his momentary preference for cake i versus j ; if $X + x_i > x_j$, he will report that cake i is better than cake j , which we will denote $i \succ j$; otherwise, he reports $j \succ i$. This will happen with some probability, $F(x_i - x_j)$:

$$F(x_i - x_j) := P[i \succ j \mid \vec{x}].$$

²CLAIRES stands for **CLAIRES Lets Anyone Infer Ranks Easily**.

There are a few nice properties that will be required of this function for it to be useful; obtaining these properties will help us impose restrictions on the distribution that the preference variable X is drawn from.

1. $F : \mathbb{R} \rightarrow [0, 1]$, making it produce meaningful probabilities.
2. $F - 1/2$ is odd; also required since the choices are forced. That is, when a taster is given cakes i and j , they must either select $i \succ j$ or $j \succ i$. Consequently, F' is even. This implies that X must be drawn from a distribution that is symmetric about $X = 0$ (i.e., one whose density function is even).
3. $F(x) \rightarrow 1$ as $x \rightarrow \infty$. That is, our judges will correctly identify differences that are far apart. If we added a random selection phase to our process (i.e., after a decision is made, with some probability, a taster may decide to flip a coin instead of tasting the cakes), it would violate this condition. This will happen if X is finite in expectation.
4. $\log F$ is concave, so when the tasters report $i \succ j$ more often than they report that $\ell \succ m$, it indicates that it is quite likely that $x_i - x_j > x_\ell - x_m$. This will happen if the density function of X decreases somewhat rapidly for positive X .
5. F has a continuous first derivative (that is, F is C^1). This requires X to be drawn from a distribution which has a continuous density.

Now that we've put together a model, given some comparisons C from our tasters (where C_{ij} is the number of times that a taster reported $i \succ j$), we'd like to find the set of scores $\vec{x} \in \mathbb{R}^N$ that is most likely given our data. So we'd like to find the value of \vec{x} that maximizes the following quantity:

$$\begin{aligned} P(\vec{x} | C) &\propto P(C | \vec{x}) \\ &= \prod_{i=1}^N \prod_{j=i+1}^N \binom{c_{ij} + c_{ji}}{c_{ij}} F^{c_{ij}}(x_i - x_j) F^{c_{ji}}(x_j - x_i). \end{aligned}$$

It's more convenient to work with the logarithm of this quantity; since log monotonically increases, we can maximize that instead:

$$\begin{aligned}\mathcal{L}(C | \vec{x}) &= \log P(C | \vec{x}) \\ &= \sum_{i=1}^N \sum_{j=i+1}^N \left[\log \binom{c_{ij} + c_{ji}}{c_{ij}} + c_{ij} \log F(x_i - x_j) + c_{ji} \log F(x_j - x_i) \right].\end{aligned}$$

Since C is constant with respect to \vec{x} , we can ignore the first term and optimize this:

$$\sum_{i=1}^N \sum_{j=i+1}^N [c_{ij} \log F(x_i - x_j) + c_{ji} \log F(x_j - x_i)].$$

This is a mapping from \mathbb{R}^N to \mathbb{R} , so we'll compute the gradient of this quantity with respect to \vec{x} and find the critical points.

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial x_i} &= \sum_{j \neq i} \left[\frac{c_{ij} F'(x_i - x_j)}{F(x_i - x_j)} - \frac{c_{ji} F'(x_j - x_i)}{F(x_j - x_i)} \right] \\ &= \sum_{j \neq i} F'(x_i - x_j) \left[\frac{c_{ij}}{F(x_i - x_j)} - \frac{c_{ji}}{1 - F(x_i - x_j)} \right].\end{aligned}$$

Assuming the c_{ij} are all nonzero, we can make a few observations about this quantity:

1. If we hold x_j constant for $j \neq i$ and increase i , $\partial \mathcal{L} / \partial x_i$ will decrease (since $\log F$ is concave).
2. If we increase x_j for some $j \neq i$ and hold the remaining components constant, $\partial \mathcal{L} / \partial x_i$ will increase.
3. If we select a subset $A \subset \{1, 2, \dots, N\}$ and increase x_i for $i \in A$ such that $x_i - x_j$ remains constant for $i, j \in A$, and the remaining components of \vec{x} also remain constant, $\partial \mathcal{L} / \partial x_i$ will decrease for $i \in A$ and increase for $i \notin A$.

Assume that we have two distinct critical points \vec{x} and \vec{y} . Without loss of generality,

order the indices by the difference $y_i - x_i$, so that:

$$y_1 - x_1 \leq y_2 - x_2 \leq \dots \leq y_N - x_N.$$

We construct a sequence of points $\vec{x}^{(1)}, \vec{x}^{(2)}, \dots, \vec{x}^{(N)}$ in the following fashion:

$$\begin{aligned} \vec{x}^{(0)} &= \vec{x} \\ x_i^{(t+1)} &= \begin{cases} x_i^{(t)} + (y_t - x_i^{(t)}) & \text{for } i \geq t \\ x_i^{(t)} & \text{otherwise} \end{cases} \end{aligned}$$

Hence, $\vec{x}^{(N)} = \vec{y}$. Observe that $\vec{\nabla} \mathcal{L}(\vec{x}^{(1)}) = \vec{\nabla} \mathcal{L}(\vec{x})$ since the two vectors only differ by a constant offset in each coordinate, and that :

$$0 = \frac{\partial \mathcal{L}(\vec{x})}{\partial x_N} = \frac{\partial \mathcal{L}(\vec{x}^{(1)})}{\partial x_N} \leq \frac{\partial \mathcal{L}(\vec{x}^{(2)})}{\partial x_N} \leq \dots \leq \frac{\partial \mathcal{L}(\vec{x}^{(N)})}{\partial x_N} = \frac{\partial \mathcal{L}(\vec{y})}{\partial x_N} = 0.$$

By the sandwich theorem, all of the quantities are 0; hence, the change from $\vec{x}^{(N-1)}$ to $\vec{x}^{(N)}$ implies that $x_N - y_N = x_{N-1} - y_{N-1}$, etc. That is, $x_i - y_i = x_j - y_j$ for all i, j ; thus $\mathcal{L}(\vec{x}) = \mathcal{L}(\vec{y})$. Thus, every critical point is unique up to an offset. Ignoring such offsets, there is at most one global maximum likelihood.

If we close \mathbb{R}^N to include infinite values and extend \mathcal{L} to take its limiting values in that space and remain continuous, it must achieve a maximum value; hence there is a unique global maximum up to offset.

Now, we can optimize \mathcal{L} using a nonlinear optimization method (such as conjugate gradient with Newton-Raphson and Fletcher-Reeves in the case where F is twice differentiable); this is guaranteed to converge rapidly to a local maximum from a random starting point; since every local maximum is a global maximum, **it is easy to optimize \mathcal{L} for nice F 's** (those that satisfy the properties we outlined above).

The psychometrics literature suggests the logistic function as a realistic choice for F (obtained by drawing X from a logistic distribution with mean 0, scale 1/2). This

is especially convenient when it comes to computing the gradient:

$$\begin{aligned} F(x_i - x_j) &= \frac{1}{2}(\tanh(x_i - x_j) + 1) \\ &= \frac{e^{x_i - x_j}}{e^{x_i - x_j} + e^{x_j - x_i}}. \end{aligned}$$

This helpfully satisfies the differential equation:

$$F' = F - F^2,$$

hence,

$$\begin{aligned} \frac{\partial \mathcal{L}(\vec{x})}{\partial x_i} &= \sum_{j \neq i} \frac{F'(x_i - x_j)}{F(x_i - x_j)(1 - F(x_i - x_j))} [c_{ij}(1 - F(x_i - x_j)) - c_{ji}F(x_i - x_j)] \\ &= \sum_{j \neq i} \frac{F'(x_i - x_j)}{F(x_i - x_j)(1 - F(x_i - x_j))} [c_{ij} - (c_{ij} + c_{ji})F(x_i - x_j)] \\ &= \sum_{j \neq i} [c_{ij} - (c_{ij} + c_{ji})F(x_i - x_j)]. \end{aligned}$$

We compute the Hessian H for use in nonlinear optimization:

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}(\vec{x})}{\partial x_i \partial x_j} \\ &= \begin{cases} (c_{ij} + c_{ji})F'(x_i - x_j) & \text{for } i \neq j \\ -\sum_{k \neq i} (c_{ik} + c_{ki})F'(x_i - x_k) & \text{otherwise.} \end{cases} \end{aligned}$$

3.3.1 Active Ranking by Forced Comparisons

Suppose now that having a person taste the cakes is expensive. What can we do to minimize the number of comparisons that must be made? We would like to maximize the amount of information that is gained by each comparison (in expectation), so let's try to minimize the entropy of the distribution $P(\vec{x} | C)$.

Let's begin by considering the effect of including an additional comparison $i \succ j$

on the distribution of \vec{x} :

$$\begin{aligned} P(\vec{x} | C \cup \{i \succ j\}) &= \frac{P(i \succ j | \vec{x}, C)P(\vec{x} | C)}{P(i \succ j | C)} \\ &= \frac{P(i \succ j | \vec{x})P(\vec{x} | C)}{P(i \succ j | C)}. \end{aligned}$$

We would like to ask the user to compare the items i and j that will have the greatest impact on the entropy in expectation:

$$\begin{aligned} &E_{i \succ j} [H(\vec{x} | C \cup \{i \succ j\}) | C] \\ &= P(i \succ j | C)H(\vec{x} | C \cup \{i \succ j\}) + P(j \succ i | C)H(\vec{x} | C \cup \{j \succ i\}) \\ &= P(i \succ j | C) \int_{\mathbb{R}^n} P(\vec{x} | C \cup \{i \succ j\}) \log P(\vec{x} | C \cup \{i \succ j\}) d\vec{x} \\ &\quad + P(j \succ i | C) \int_{\mathbb{R}^n} P(\vec{x} | C \cup \{j \succ i\}) \log P(\vec{x} | C \cup \{j \succ i\}) d\vec{x} \\ &= \int_{\mathbb{R}^n} P(i \succ j | \vec{x})P(\vec{x} | C) \log P(\vec{x} | C \cup \{i \succ j\}) d\vec{x} \\ &\quad + \int_{\mathbb{R}^n} P(j \succ i | \vec{x})P(\vec{x} | C) \log P(\vec{x} | C \cup \{j \succ i\}) d\vec{x} \\ &= \int_{\mathbb{R}^n} P(i \succ j | \vec{x})P(\vec{x} | C) \log P(i \succ j | \vec{x}) d\vec{x} \\ &\quad + \int_{\mathbb{R}^n} P(j \succ i | \vec{x})P(\vec{x} | C) \log P(j \succ i | \vec{x}) d\vec{x} \\ &\quad + \int_{\mathbb{R}^n} P(i \succ j | \vec{x})P(\vec{x} | C) \log P(\vec{x} | C) d\vec{x} \\ &\quad + \int_{\mathbb{R}^n} P(j \succ i | \vec{x})P(\vec{x} | C) \log P(\vec{x} | C) d\vec{x} \\ &\quad - \int_{\mathbb{R}^n} P(i \succ j | \vec{x})P(\vec{x} | C) \log P(i \succ j | C) d\vec{x} \\ &\quad - \int_{\mathbb{R}^n} P(j \succ i | \vec{x})P(\vec{x} | C) \log P(j \succ i | C) d\vec{x} \\ &= E_{\vec{x}} [H(i \succ j | \vec{x}) | C] + H(\vec{x} | C) - H(i \succ j | C). \end{aligned}$$

Since $H(\vec{x} | C)$ is independent of i and j , we simply choose the pair which minimizes

$$E_{\vec{x}} [H(i \succ j | \vec{x}) | C] - H(i \succ j | C).$$

This is difficult to compute analytically; however, it can be estimated to an arbitrary degree of accuracy via importance sampling, using a multivariate Gaussian centered on the maximum likelihood value of \vec{x} and with a covariance matrix such that the Gaussian has the same Hessian at the mean as \mathcal{L} (i.e., with $\Sigma = H^{-1}$ —since the nullity of H is one, eliminating any row will make H invertible). The full algorithm is given in figures 3-2 and 3-3.

3.3.2 Previous Work

Forced choice methods were introduced in the late 20s by L. L. Thurstone, who used them to measure social attitudes, and have been extensively studied in the psychometric community, where they are used to measure human sensitivity to various stimuli. They have also been used to develop guided decision-making processes based on both objective and subjective criteria. Unexpectedly, the type of analysis that we have gone about here does not (to the best of our knowledge) exist in the literature, perhaps because the prevailing scenarios in these fields are directed at minimizing computation instead of minimizing the effort of those surveyed.

Analogous models are also used to establish rankings for games, notably ELO for chess and (much more recently) TrueSkill for multiplayer games on Xbox Live [21]. It is rare to be able to select pairs for comparison in these scenarios, and these models are generally restricted in that only the scores of the items compared should be affected by a comparison, whereas in our case, a single comparison can affect all of the scores.

We believe our model has wide application. A python version of library that uses Amazon Mechanical Turk for data collection is available.

3.3.3 Consequences of Choosing a Bad Activation Function

Although scaling X does not functionally change F , changing $G(p) = F(2F^{-1}(p))$ does. Given a Lipschitz continuous bijection $G : [0, 1] \rightarrow [0, 1]$, there exists a unique F

function FINDOPTIMALWEIGHTS(c):

- $\vec{x} =$ random vector in $[0, 1]^N$
- $i = 0$
- $k = 0$
- $\vec{r} = \nabla \mathcal{L}(\vec{x})$
- $\vec{d} = \vec{r}$
- $\delta' = \|\vec{r}\|^2$
- $\delta_0 = \delta'$
- **while** $i < i_{max}$ **and** $\delta' > \varepsilon^2 \delta_0$:
 - $j = 0$
 - $\delta_D = \|\vec{d}\|^2$
 - $\alpha = (\varepsilon^2 / \delta_D) / 2 + 1$
 - **while** $j < j_{max}$ **and** $\alpha^2 \delta_D > \varepsilon^2$:
 - $\alpha = (\nabla \mathcal{L}(\vec{x}) \cdot \vec{d}) / (\vec{d} \cdot (\mathcal{H}(\vec{x}) \vec{d}))$
 - $\vec{x} = \vec{x} + \alpha \vec{d}$
 - $j = j + 1$
 - $\vec{r} = \nabla \mathcal{L}(\vec{x})$
 - $\delta_o = \delta'$
 - $\delta' = \|\vec{r}\|^2$
 - $\beta = \delta' / \delta_o$
 - $\vec{d} = \vec{r} + \beta \vec{d}$
 - $k = k + 1$
 - **if** $k == N$ **or** $\vec{r} \cdot \vec{d} < 0$:
 - $\vec{d} = \vec{r}$
 - $k = 0$
 - $i = i + 1$
- **return** \vec{x}

Figure 3-2: Find the optimal scores based on counts; this is a nonlinear gradient descent algorithm. \mathcal{H} denotes the Hessian of \mathcal{L} .

do:

- $\vec{x} = \text{FINDOPTIMALWEIGHTS}(c)$
- $\vec{x}_0 = \vec{x}$ with the first element removed
- $\mathcal{H} = \text{Hessian of } \mathcal{L} \text{ at } \vec{x} \text{ with the first row and column removed}$
- Diagonalize \mathcal{H} , yielding $\mathcal{H} = QDQ'$
- $A = QD^{-1/2}$ (so that $\Sigma = AA' = \mathcal{H}^{-1}$)
- initialize samples, an empty list of (sample, weight) tuples
- **repeat many times:**
 - · $\vec{z}_0 = \text{an } (N - 1)\text{-dimensional vector of unit normal samples}$
 - · $\vec{z}_0 = A\vec{z}_0 + \vec{x}_0$
 - · $\vec{z} = [\|\vec{x}\|_1 - \|\vec{x}_0\|_1 \mid \vec{z}_0]$ (augment the vector)
 - · $\text{weight} = \exp(\mathcal{L}(\vec{z})) - \exp(-((\vec{z}_0 - \vec{x}_0) \cdot A(\vec{z}_0 - \vec{x}_0))/2) / \sqrt{(2\pi)^{(N-1)}|\mathcal{H}|}$
 - · $\text{samples.append}((\vec{z}, \text{weight}))$
- $\text{bestScore} = -\infty$; $\text{bestPair} = \text{None}$
- **for** i **in** $[0, \dots, N - 2]$:
 - · **for** j **in** $[i + 1, \dots, N - 1]$:
 - · · $Z = 0$; $\hat{h} = 0$; $\hat{f} = 0$
 - · · **for** \vec{z}, weight **in** samples :
 - · · · $s_{ij} = F(z_i - z_j)$
 - · · · $\hat{h} = \hat{h} + (s_{ij} \log s_{ij} + (1 - s_{ij}) \log(1 - s_{ij})) \cdot \text{weight}$
 - · · · $\hat{f} = \hat{f} + s_{ij} \cdot \text{weight}$
 - · · · $Z = Z + \text{weight}$
 - · · · $\hat{h} = \hat{h} / Z$; $\hat{f} = \hat{f} / Z$
 - · · · $\text{score} = \hat{h} - (\hat{f} \log \hat{f} + (1 - \hat{f}) \log(1 - \hat{f}))$
 - · · · **if** $\text{bestScore} < \text{score}$:
 - · · · · $\text{bestScore} = \text{score}$; $\text{bestPair} = (i, j)$
 - ask the user to compare the items in bestPair and update c

until convergence

Figure 3-3: Active ranker pseudocode to rank N items.

(up to scaling) that satisfies this relation; this gives a way that one could (in practice) argue for a particular choice of F .

Suppose that for some set of items, there is a true activation function G with true scores $\hat{x}_1, \dots, \hat{x}_N$ and we are training against an activation function F . Since $\log F$ is concave, there is a path from the true values $\hat{x}_1, \dots, \hat{x}_N$ to the values optimized for F $\vec{x}_1, \dots, \vec{x}_N$. If the order of these values differs, then there will be some cross-over point where all of the items are ordered properly except for one pair. However, in the limit of infinitely many observations, this will make the likelihood decrease. Hence, in the limit, the ranking of the items will be correct regardless of the choice of activation function.

3.4 Mechanical Turk

As in the previous chapter, we make use of Mechanical Turk for our judgments. Once again, the tasks are easy to complete (average time is ~ 30 s), showing that they do not impart a significant psychometric load suggesting that the results are reliable. Once again we rely on a preposition test as proof that the workers have reasonable English fluency and are giving the questionnaire some attention. An example question including the test question is shown in Figure 3.4.

3.5 Results

When applying this technique to the machine translation task, we compared our results to those obtained in the 2007 Meta-evaluation task, which obtained results from human evaluators and a number of automatic metrics [9]. See table 3.3 for the results.

Fill in the blank (mark all the words that would result in a grammatical sentence):

I _____ the supplier to divert the shipment.

- wanted
- know
- saw
- seemed

Please select the sentence that is closest in meaning to the italicized sentence.

It should be noted that this right to be different is nothing more than a demonstration of the principle of equality, which also requires different treatment for anything that is different.

- This right of differentiation is nothing more than a result of the principle of equality, which is also to deal with what is different.
- This right of differentiation is nothing further as a result of the principle of equality, therefore also the conditions to deal with what is different.

Figure 3-4: Sample Mechanical Turk Questionnaire including a test question.

	CLAIRE	(sigma)	BLEU	Human Evaluation
systran	0.23015	0.027	0.154	1
uedin	0.21232	0.026	0.277	2
liu	0.11797	0.024	0.263	4
nrc	0.08926	0.024	0.254	5
saar	0.06191	0.024	0.198	6
cmu-uka	0.01335	0.024	0.247	7
upc	0.00000	0.021	0.250	3

Table 3.3: Our proposed metric, CLAIRE, captures the results of the shared evaluation and correlates strongly across different types of translation systems.

3.6 Conclusions

We propose an alternative to BLEU and other automatic metrics, which we call CLAIRE. CLAIRE is a method for performing evaluations of translation systems which judges them without undue preference toward a single architecture. CLAIRE scores can be obtained quickly, they are comparable against each other, and correlate strongly with more traditional full-scale human judgments of quality. The benefits of a clear and consistent design carry over from the automatic metrics, and the costs of performing incremental evaluations are minimal ($< \$20$ for a full evaluation).



Chapter 4

Conclusions and Future Work

The question answered by a typical publication in NLP is: “does modelling *⟨phenomenon⟩* improve *⟨metric⟩*?” Consequently, we have learned how to improve certain popular metrics. Over time, differences between the metric and genuine measures of performance yield designs that optimize the former at the expense of the latter.

In particular, we analyze BLEU, the dominant metric in statistical machine translation. We show differences between BLEU and human measures of performance, issues with the papers that originally presented BLEU, and show how systems have optimized BLEU scores to eventually yield diminished performance in human evaluations. We argue that automatic metrics are just as difficult to design as automatic systems. Consequently, we contribute a new metric based on inexpensive human evaluations in the cloud guided by the psychometrics literature; our metric correlates strongly with more expensive conventional human evaluations and is sensitive to minor differences in performance.

We proposed a promising research agenda that can gradually draw us closer to this goal, putting us hand-in-hand with modern linguistics efforts, and more realistically promising systems that are broadly usable. We show ways of applying this technique to parsing, word-sense disambiguation, and machine translation.

Invoking Fred Jelinek’s quote from the introduction yields a second interpretation:

Every time I fire a linguist, the performance goes up!

The fruits of linguistics as a field extend beyond linguistic theories for the language process: the many researchers devoted to evaluating the theories scientifically by providing counterexamples that falsify them. Language systems are a marriage between linguistic theory and statistical model: these are scientific elements that warrant scientific rigor. It is difficult to imagine an established scientific enterprise directed at optimizing an automatic metric such as BLEU: physical models so motivated would ignore bodies moving near the speed of light, and the strange behavior of very small particles passing through narrow slits as rare events that are infrequently experienced. How could we have progressed beyond Newtonian mechanics by ignoring the unusual phenomena that lie at the fringe of our models?

As Fred Jelinek himself says, researchers have a right to devote themselves to the solution of intrinsically interesting questions even during an era of senseless product competition. We should be willing to accept an immediate reduction in performance (even in our improved metrics) if we are building a principled system that can reasonably offer enhanced performance in the long-run, and funding agencies should be willing to accept this, too.

The members of ALPAC did not reserve all of their criticism for machine translation. ALPAC's chairman, John R. Pierce, wrote of speech recognition in 1969[41]:

Most recognizers behave, not like scientists, but like mad inventors or untrustworthy engineers. The typical recognizer gets it into his head that he can solve "the problem." The basis for this is either individual inspiration (the "mad inventor" source of knowledge) or acceptance of untested rules, schemes, or information (the untrustworthy engineer approach)... The typical recognizer... builds or programs an elaborate system that either does very little or flops in an obscure way. A lot of money and time are spent. No simple, clear, sure knowledge is gained. The work has been an experience, not an experiment.

It took decades for speech recognition to plateau, to exhaust the immediate gratification that can be had by doubling the clock-speed of a server or the size of a data-set; recent research in speech recognition begins to look like real science, including explorations on the benefits of linguistic representation, a willingness to sacrifice performance in the short-term for nuanced models that can capture the rare events as well as the common ones. It has become a respectable scientific enterprise with many applications and (with the pervasiveness of mobile phones) ubiquity. Research agenda in hand, I have high hopes that machine translation can match and surpass that success.



Appendix A

EM Algorithm Reference

A.1 Definitions

$$\begin{aligned}\mathcal{L}(\Theta) &= \sum_X \log P(X; \Theta) && \text{Likelihood} \\ Q(\Theta', \Theta) &= \sum_X E_Y (\log P(X, Y; \Theta') | X; \Theta) \\ H(\Theta', \Theta) &= - \sum_X E_Y (\log P(Y | X; \Theta') | X; \Theta) && \text{Cross-Entropy} \\ KL(\Theta', \Theta) &= H(\Theta', \Theta) - H(\Theta, \Theta) && \text{KL-Divergence} \\ \Theta^{(i)} &= \arg \max_{\Theta} Q(\Theta, \Theta^{(i-1)})\end{aligned}$$

A.2 Lemmas

Jensen's Inequality:

$$E(\log(X)) \leq \log(E(X))$$

An Obvious Identity:

$$E_X \left(\frac{f(X)}{P(X)} \right) = \int_X f(X) dX$$

Lemma:

$$H(\Theta, \Theta) \leq H(\Theta', \Theta) \quad \forall \Theta, \Theta'$$

Equivalently:

$$KL(\Theta', \Theta) \geq 0 \quad \forall \Theta, \Theta'$$

Proof:

$$\begin{aligned} \forall X \quad 0 &= \log(1) = \log \left(\int_Y P(Y|X; \Theta') dY \right) \\ &= \log \left(\int_Y P(Y|X; \Theta) \cdot \frac{P(Y|X; \Theta')}{P(Y|X; \Theta)} dY \right) \\ &= \log E_Y \left(\frac{P(Y|X; \Theta')}{P(Y|X; \Theta)} \middle| X; \Theta \right) \\ &\geq E_Y \left(\log \left[\frac{P(Y|X; \Theta')}{P(Y|X; \Theta)} \right] \middle| X; \Theta \right) \\ &= E_Y (\log P(Y|X; \Theta') | X; \Theta) - E_Y (\log P(Y|X; \Theta) | X; \Theta) \\ E_X(0) = 0 &\geq E_X (E_Y (\log P(Y|X; \Theta') | X; \Theta)) - E_X (E_Y (\log P(Y|X; \Theta) | X; \Theta)) \\ &= H(\Theta, \Theta) - H(\Theta', \Theta). \end{aligned}$$

$$\Rightarrow H(\Theta, \Theta) \leq H(\Theta', \Theta)$$

A.3 EM is Nondecreasing

Theorem:

$$\mathcal{L}(\Theta^{(i)}) \leq \mathcal{L}(\Theta^{(i+1)})$$

Proof:

$$\begin{aligned} \mathcal{L}(\Theta) &= \sum_X \log P(X; \Theta) \\ &= \sum_X E_Y (\log P(X; \Theta) | X; \Theta^{(i)}) \\ &= \sum_X E_Y (\log P(X; \Theta) | X; \Theta^{(i)}) - Q(\Theta, \Theta^{(i)}) + Q(\Theta, \Theta^{(i)}) \end{aligned}$$

$$\begin{aligned}
&= \sum_X (E_Y (\log P(X; \Theta) | X; \Theta^{(i)}) - E_Y (\log P(X, Y; \Theta) | X; \Theta^{(i)})) \\
&\quad + Q(\Theta, \Theta^{(i)}) \\
&= - \sum_X E_Y \left(\log \left(\frac{P(X, Y; \Theta)}{P(X; \Theta)} \right) \middle| X; \Theta^{(i)} \right) + Q(\Theta, \Theta^{(i)}) \\
&= - \sum_X E_Y (\log P(Y | X; \Theta) | X; \Theta^{(i)}) + Q(\Theta, \Theta^{(i)}) \\
&= H(\Theta, \Theta^{(i)}) + Q(\Theta, \Theta^{(i)}). \\
\mathcal{L}(\Theta^{(i+1)}) - \mathcal{L}(\Theta^{(i)}) &= [H(\Theta^{(i+1)}, \Theta^{(i)}) - H(\Theta^{(i)}, \Theta^{(i)})] \\
&\quad + [Q(\Theta^{(i+1)}, \Theta^{(i)}) - Q(\Theta^{(i+1)}, \Theta^{(i)})] \\
&\geq 0.
\end{aligned}$$

Thus, the likelihood of successive EM parameter vectors is non-decreasing. (This is a long way from convergence proof...)

Incidentally, we have also shown that, $\forall \Theta, \Theta'$:

$$\boxed{\mathcal{L}(\Theta) = H(\Theta, \Theta') + Q(\Theta, \Theta')}$$

A.4 EM on Multinomials

EM is easy in the special case when $P(X, Y | \Theta)$ is a multinomial distribution; that is, it can be written in the form:

$$P(X, Y; \Theta) = \prod_{r=1}^N \Theta_r^{\text{Count}_r(X, Y)}.$$

This is a form that occurs very often in language processing tasks.

Let's go!

$$\begin{aligned}
Q(\Theta', \Theta) &= \sum_X E_Y (\log P(X, Y; \Theta') | X; \Theta) \\
&= \sum_X \sum_{r=1}^N E_Y (\text{Count}_r(X, Y) | X; \Theta) \log \Theta'_r
\end{aligned}$$

$$= \sum_{r=1}^N \left(\sum_X E_Y (\text{Count}_r(X, Y) | X; \Theta) \right) \cdot \log \Theta'_r$$

It's easy to write an algorithm to maximize Q ; we just have to set each Θ'_r to its coefficient and normalize (in ways corresponding to inherent constraints of the parameters):

$$\Theta'_r \propto \sum_X E_Y (\text{Count}_r(X, Y) | X; \Theta).$$

Let's formalize the normalization; the indices $r \in \{1, 2, \dots, N\}$ are partitioned into disjoint subsets R_1, R_2, \dots, R_m such that $\sum_{r \in R_i} \Theta_r = 1$. Now we set:

$$\Theta'_r = \frac{\sum_X E_Y (\text{Count}_r(X, Y) | X; \Theta)}{\sum_{r' \in R_i} \sum_X E_Y (\text{Count}_{r'}(X, Y) | X; \Theta)} \quad \forall r \in R_i$$

A.5 Parametrizing Multinomials

Suppose we wish to further parametrize the parameters Θ by another set of parameters α ; that is, we define a set of events E_r and set $\Theta_r = P(E_r; \alpha)$ whilst preserving the normalization conditions on Θ_r (i.e., that $\sum_{r \in R_i} P(E_r; \alpha) = 1, \forall R_i$); thus,

$$P(X, Y; \alpha) = \prod_{r=1}^N P(E_r; \alpha)^{\text{Count}_r(X, Y)}.$$

We assume that we can easily find maximum-likelihood α given E_r data (that is, the number of occurrences of each event E_r —not necessarily integral). We proceed:

$$\begin{aligned} Q(\alpha', \alpha) &= \sum_X E_Y (\log P(X, Y; \alpha') | X; \alpha) \\ &= \sum_{r=1}^N \left(\sum_X E_Y (\text{Count}_r(X, Y) | X; \alpha) \right) \log P(E_r; \alpha') \end{aligned}$$

Then, the EM update of α is simply the maximum-likelihood α' where event E_r has

occurred this many times:

$$\sum_X E_Y(\text{Count}_r(X, Y) | X, \alpha).$$

Note that these counts are pre-normalization! When the α_i are “grouped” the same way as the Θ_r , the normalization does not enter into the picture (that is, when

$$P(E_r; \alpha') = P(E_r; \alpha'_i) \quad \forall r \in R_i,$$

where the R_i are defined as above to be sets of Θ parameters that must be normalized); consequently, the Q -maximizing Θ_r can themselves be used as counts to maximize the likelihood of α . Again, if, for any i , we were to multiply the coefficients of Θ_r for $r \in R_i$ by a constant, the maximum likelihood values do not change; thus, no normalization is necessary.

[It is easy to see that normalization can be harmful; consider, for instance, the following experiment: we repeatedly select one of two biased coins to flip and record which coin we flipped and the outcome. Then, the maximum-likelihood probability that the first coin will flip heads, for instance, is the number of heads we got from the first coin divided by the number of times we flipped the first coin. Suppose, however, that we add the constraint that the coins are identically biased. Then the number of times we flipped each coin is important; we cannot correctly estimate the probability of heads with the unconstrained maximum-likelihood probability of heads for each coin alone.]

A.6 EM on IBM2+1dG

Here we give the full derivation of the EM updates for the one-dimensional gaussian framework described in Appendix D for the sake of the mathematically skeptical.

Let’s begin by defining the model:

- Training Data

The training data consists of several triplets $(\mathbf{e}, \mathbf{f}, \vec{a})$.

1. \mathbf{e} and \mathbf{f} are English and French sentences, respectively, that are translations of each other (observed in training).
2. $\vec{a} \in \{0, \dots, \ell\}^m$ represents an alignment between the sentences, where $\ell = |\mathbf{e}|$ and $m = |\mathbf{f}|$ and $a_j = i$ implies that the i th English word corresponds to the j th French word (hidden in training).

- Parameters

1. $T(f|e)$ for all French words f and English words e (the NULL word is added to the English vocabulary).
2. $\mu_{j,\ell,m}$, $\sigma_{j,\ell,m}$ for all French sentence lengths m , English sentence lengths ℓ , and French word indices $j \in \{1, \dots, m\}$, respectively the mean and standard deviation of the index of the corresponding English word, given that it is not the NULL word.
3. $N(j, \ell, m)$ for all French sentence lengths m , English sentence lengths ℓ , and French word indices $j \in \{1, \dots, m\}$, the probability that that French word is aligned to the English NULL word.

- Model

$$P(\mathbf{f}, \vec{a} | \mathbf{e}; T, \mu, \sigma) = \prod_{j=1}^m T(\mathbf{f}_j | \mathbf{e}_{a_j}) D(a_j | j, \ell, m)$$

where

$$D(a_j | j, \ell, m) = \begin{cases} (1 - N(j, \ell, m)) \cdot f_{\mathcal{N}}(a_j | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) & j \neq 0 \\ N(j, \ell, m) & j = 0 \end{cases}$$

Note that this model is deficient; that is, we are not enforcing the proper normalization constraints on D .

- Normalization conditions

$$\sum_f T(f|e) = 1 \quad \forall e$$

- Parametrized Multinomial

$$P(\mathbf{f}, \bar{\mathbf{a}} | \mathbf{e}; T, \mu, \sigma) = \prod_{f,e} T(f|e)^{\text{Count}_{f,e}(\mathbf{e}, \mathbf{f}, \bar{\mathbf{a}})} \times \prod_{i,j,\ell,m} D(i|j, \ell, m)^{\text{Count}_{i,j,\ell,m}(\mathbf{e}, \mathbf{f}, \bar{\mathbf{a}})}$$

where

$$\begin{aligned} \text{Count}_{f,e}(\mathbf{e}, \mathbf{f}, \bar{\mathbf{a}}) &= \sum_{j=1}^m \delta(f, \mathbf{f}_j) \delta(e, \mathbf{e}_{a_j}) \\ \text{Count}_{i,j,\ell,m}(\mathbf{e}, \mathbf{f}, \bar{\mathbf{a}}) &= \delta(\ell, |\mathbf{e}|) \delta(m, |\mathbf{f}|) \delta(a_j, i), \end{aligned}$$

where $\delta(\cdot, \cdot)$ denotes the Kronecker delta function.

We parametrize T and D by T, μ, σ , and N . Thus, T is trivially parametrized by itself, whereas D is parametrized as described above by μ, σ , and N .

We wish to maximize the function $Q((T, \mu, \sigma, N), (T^i, \mu^i, \sigma^i, N^i))$ given T^i, μ^i, σ^i , and N^i . Since the model can neatly be factored into a term that depends on T alone and a term that depends on D alone, we can optimize these parameters independently. Clearly, the EM procedure for finding the optimal T is unchanged from Model 2; thus, we need only focus on finding the optimal values for μ, σ , and N . We compute:

$$\begin{aligned} &Q((\mu, \sigma, N), (T', \mu', \sigma', N')) \\ &= \sum_{\mathbf{e}, \mathbf{f}} \sum_{\bar{\mathbf{a}}} P(\bar{\mathbf{a}} | \mathbf{e}, \mathbf{f}, T', N', \mu', \sigma') \log P(\mathbf{f}, \bar{\mathbf{a}} | \mathbf{e}; N, \mu, \sigma) \\ &= \sum_{i,j,\ell,m} \left[\underbrace{\left(\sum_{\mathbf{e}, \mathbf{f}} E_{\bar{\mathbf{a}}}(\text{Count}_{i,j,\ell,m}(\mathbf{e}, \mathbf{f}, \bar{\mathbf{a}}) | \mathbf{e}, \mathbf{f}; T', \mu', \sigma', N') \right)}_{C(i, j, \ell, m)} \log D(i|j, \ell, m) \right] \\ &= \sum_{j,\ell,m} C(0, j, \ell, m) \log N(j, \ell, m) \end{aligned}$$

$$\begin{aligned}
& + \sum_{j,\ell,m} \sum_{i=1}^{\ell} C(i, j, \ell, m) \log[(1 - N(j, \ell, m)) f_{\mathcal{N}}(i | \mu_{i,j,\ell,m}, \sigma_{i,j,\ell,m})] \\
= & \sum_{j,\ell,m} C(0, j, \ell, m) \log N(j, \ell, m) + \sum_{j,\ell,m} \sum_{i=1}^{\ell} C(i, j, \ell, m) \log(1 - N(j, \ell, m)) \\
& + \sum_{j,\ell,m} \sum_{i=1}^{\ell} C(i, j, \ell, m) \log f_{\mathcal{N}}(i | \mu_{i,j,\ell,m}, \sigma_{i,j,\ell,m}).
\end{aligned}$$

Thus, the optimal value of N is given by:

$$N(j, \ell, m) = \frac{C(0, j, \ell, m)}{\sum_{i=0}^{\ell} C(i, j, \ell, m)}$$

and the μ and σ are optimized by their usual maximum-likelihood estimators:

$$\begin{aligned}
\mu_{j,\ell,m} &= \frac{\sum_{i=1}^{\ell} i C(i, j, \ell, m)}{\sum_{i=1}^{\ell} C(i, j, \ell, m)} \\
\sigma_{j,\ell,m}^2 &= \frac{\sum_{i=1}^{\ell} (i - \mu_{j,\ell,m})^2 C(i, j, \ell, m)}{\sum_{i=1}^{\ell} C(i, j, \ell, m)}
\end{aligned}$$

Appendix B

Information on the WSJ Corpus

Tagset

The annotations in the Penn treebank in general are quite difficult to interpret; the manual that is distributed with it is a mere 300 pages or so. The early chapters are extremely formal and handle “simple” sentences, but the last half of the book is special cases and unusual structures that are difficult to analyze.

We’ll begin with the nominal terminal tags, a list of all of the members of each of the closed-class parts-of-speech, a list of the nonterminal tags, a list of the extended tags and their definitions, and a list of all of the combinations of extended tags with representative examples that appear in the corpus. This final table suffers contradictions and redundancies; clearly the order of the extended tags is important in some situations, unimportant in others; the inconsistencies in the table reflect corresponding inconsistencies in the corpus and are not errors.

B.1 Terminal Tags (Parts of Speech)

CC	Coordinating Conjunction
CD	Cardinal Number

DT	Determiner
EX	Existential <i>there</i>
FW	Foreign Word
IN	Preposition/Subordinating Conjunction
JJ[RS]?	Adjective
	R Comparative
	S Superlative
LS	List
MD	Modal
NNP?S?	Noun (singular or mass)
	P Proper
	S Plural
PDT	Predeterminer
POS	Possessive Ending
PRP\$?	Personal Pronoun
	\$ Possessive Pronoun
RB[RS]?	Adverb
	R Comparative
	S Superlative
RP	Particle
SYM	Symbol
TO	<i>to</i>
UH	Interjection
VB[DGNPZ]?	Verb (base form)
	D past tense
	G gerund
	N past participle
	P present, not 3rd-person singular
	Z present, 3rd-person singular
WDT	Wh-determiner

WP\$?	Wh-pronoun
	\$ Wh-possessive Pronoun
WRB	Wh-adverb

B.2 Members of Closed-Class Tags

Punctuation

‘ ‘	‘ ‘
’ ’	’ ’
,	,
:	- - ; : ...
.	! ? .
\$	\$ A\$ C\$ HK\$ FFr M\$ NZ\$ S\$ US\$
#	#
-LRB-	-LCB- -LRB-
-RRB-	-RCB- -RRB-
SYM	= @ * ** & a b c d e f r x z

Closed Class

CC	& and and/or both but either et less minus 'n 'n' neither nor or plus so times versus vs. whether yet
DT	a all an another any both each either every half many neither no some that the them these this those
EX	there

IN aboard about above across against ago a/k/a albeit along alongside amid
among amongst around astride at atop because behind below beneath
beside besides between beyond but by despite down during except expect
fiscal for from in inside into lest like minus near nearer nearest neither
next notwithstanding of off on onto opposite out outside over par past
pending per plus post save so than that then though through throughout
'til till toward towards under underneath unless unlike until up upon v.
versus via vs. whereas whether while with within without worth
after although as because before how if once since than that though 'til
till until when where whether while¹
both... and, either... or, neither... nor, not only... but also, so... as
whether... or

LS 1 2 3 4 a b c d e f r x first second third

MD ca can could 'd dare 'll may might mighta must need ought shall should
will wo would

PDT all both half many nary quite such

POS ' 's

PRP 'em he her him his I it me mine one ours s 's she 't- t' theirs them they us
we ya y'all you herself himself itself myself ourselves themselves thyself
yourself

PRP\$ her his its my our their your

RBR about down less less-perfectly more than worse

RBS best hardest highest least most worst

RP about across ahead along apart around aside at away back before behind
by down even for forth forward in of off on open out over through together
up upon with yet

TO to na²

WDT that what whatever which whichever

¹subordinating conjunctions

²as in *gonna*

WP what who whoever whom
WP\$ whose
WRB how however when whenever where whereby wherever why

Other Things

UH ah alas amen aw bam boy damn egad heck hello howdy indeed man no
nope oh oink ok okay please quack say true uh uh-uh welcome well wham
whoopee wow yeah yes zounds

B.3 Nonterminals

ADJP Adjective Phrase
ADVP Adverb Phrase
CONJP Conjunction Phrase
FRAG Fragment
INTJ Interjection
LST List Marker
NAC Not A Constituent
NP Noun Phrase
NX NP head marker
PP Prepositional Phrase
PRN Parenthetical
PRT Particle
QP Quantifier Phrase
RRC Reduced Relative Clause
SBARQ Direct question
SBAR Subordinate clause
SINV Sentence, Inverted
SQ Sentence, Question (inverted yes/no, or the argument of a Wh)
S Sentence, Declarative
UCP Unlike Coordinated Phrase
VP Verb Phrase

WHADJP	Wh-adjective Phrase
WHADVP	Wh-adverb Phrase
WHNP	Wh-noun Phrase
WHPP	Wh-prepositional Phrase
X	Unknown, Uncertain, or Unbracketable

B.4 Extended Tags

I took the next two tables from “The Penn Treebank: Annotating Predicate Argument Structure” by Mitch Marcus *et al* [34].

Text Categories

- HLN headlines and datelines
- LST list markers
- TTL titles

Grammatical Functions

- CLF true clefts
- NOM non NPs that function as NPs
- ADV clausal and NP adverbials
- LGS logical subjects in passives
- PRD non VP predicates
- SBJ surface subject
- TPC topicalized and fronted constituents
- CLR closely related; should be part of the VP

Semantic Roles

- PUT where something is put? (not in the Marcus paper)
 - VOC vocatives
 - DIR direction and trajectory
 - LOC location
 - MNR manner
 - PRP purpose and reason
 - TMP temporal phrases
-

B.5 Pseudo-Attachment/Null Element Markers

-
- *ICH* Interpret Constituent Here
 - *PPA* Permanent Predictable Ambiguity
 - *RNR* Right Node Raising
 - *EXP* Expletive
-

B.6 Examples of the Full Extended Tags

ADJP	Adjective Phrase
ADJP	This is an <i>enviously low</i> level.
ADJP-ADV	...Sarah Lee closed <i>unchanged</i> at 60 1/8.
ADJP-CLR	falls <i>flat</i> , let <i>loose</i> , sat <i>idle</i> ended <i>mixed</i> , closed <i>higher</i> rated <i>triple-A</i> , increased <i>nearly fivefold</i>
ADJP-LOC	the site <i>adjacent to the refinery</i>
ADJP-MNR	It's better to sell <i>private</i> .
ADJP-PRD	Growth is <i>relatively small</i> .
ADJP-PRD-TPC	<i>Conspicuous by its absence</i> is California.
ADJP-SBJ	<i>Bigger</i> is better.
ADJP-TPC	<i>Not likely</i> , I think. <i>Typical</i> is this response:...
ADVP	Adverb Phrase
ADVP	<i>Again</i> , I think <i>so, too</i> . The dollar finished <i>lower</i> . No one <i>else</i> does.
ADVP-CLR	goes <i>further</i> , leave them <i>alone</i> , broke <i>loose</i>
ADVP-CLR-MNR	live <i>together</i> cooperatively
ADVP-CLR-TPC	<i>So</i> says John, <i>Along</i> comes Bobby
ADVP-DIR	shot <i>up</i> , go <i>out</i>
ADVP-DIR-CLR	bring my prices <i>down</i> , trickle <i>down</i> , head <i>off</i> in some direction
ADVP-DIR-TPC	<i>here</i> comes the sun, <i>along</i> came Bob
ADVP-EXT	It rose <i>slightly</i> , prolong <i>somewhat</i> ,
ADVP-LOC	Stop <i>there</i> , He bought <i>elsewhere</i>
ADVP-LOC-CLR	It won't get <i>there</i> , He laid <i>low</i> , want <i>out</i>
ADVP-LOC-CLR-TPC	<i>therein</i> lies the draw

ADVP-LOC-PRD	It was <i>all over the place</i> .
ADVP-LOC-PRD-TPC	<i>Here</i> is an example.
ADVP-LOC-TMP	Where and when it will occur *
ADVP-LOC-TPC	<i>Here</i> is an example.
ADVP-LOC-TPC-PRD	<i>Here</i> is an example.
ADVP-MNR	Speak <i>sincerely</i> , quit <i>suddenly</i>
ADVP-MNR-CLR	doing <i>very well</i> , do <i>better</i>
ADVP-MNR-TMP	how and when the goals would be achieved * *
ADVP-MNR-TPC	But as stock prices recovered, <i>so</i> did the US currency
ADVP-PRD	Is it <i>over</i> ? It was <i>down a little</i> , Give it <i>back</i>
ADVP-PRD-LOC	get <i>down in the dumps</i> ,
ADVP-PRD-LOC-TPC	<i>Here's</i> a look at some of the alternatives:
ADVP-PRD-TMP	It only pays when there's a catastrophe *
ADVP-PRD-TPC	<i>so</i> did he, <i>here's</i> an idea, <i>first</i> on the list of ideas is... <i>Soon to feel the glare of attention</i> are lawyers...
ADVP-PRP	Why is this happening *? Climate varies <i>due to natural causes</i> .
ADVP-PUT	That put's us <i>back in the soup</i> .
ADVP-PUT-TPC	<i>Next door</i> she put a glass house.
ADVP-TMP	I want it <i>now</i> . <i>Finally</i> , he got help.
ADVP-TMP-CLR	It starts <i>as soon as tomorrow</i> , it didn't last <i>long</i> , why it took <i>so long</i> , it begins <i>soon</i>
ADVP-TMP-PRD	Why must it be <i>so soon</i> ? That was <i>early in our history</i> , The time is <i>now</i> .
ADVP-TMP-TPC	<i>Initially</i> , the company said "hello."
ADVP-TPC	He was hungry, or <i>so</i> it seemed.
ADVP-TPC-PRD	John panicked, and <i>so</i> did Bob.
CONJP	Conjunction Phrase

CONJP	<i>as well as, rather than, not only, if not, but also, instead of</i>
FRAG	Fragment
FRAG	The game hasn't changed, <i>only the name.</i> <i>While no guarantee,</i> an increased salary might improve performance. He yelled out " <i>dolce! dolce!</i> " <i>As usual,</i> I'm going nuts.
FRAG-ADV	<i>Earthquake or not, If anything, Not only that, Who knows?</i>
FRAG-HLN	
FRAG-PRD	The answer is, " <i>Yes, of course.</i> "
FRAG-TPC	<i>Not so fast,</i> said Boris.
FRAG-TTL	His magazine, " <i>Cornhuskers,</i> " criticizes wheat farmers.
INTJ	Interjection
INTJ	<i>Yes, they are.</i> He said <i>no</i> again. It was, <i>well,</i> fake.
INTJ-CLR	He's learning to say <i>no.</i>
INTJ-HLN	
LST	List Marker
LST	2. Provide better toilet paper.
NAC	Not A Constituent
NAC	<i>MIT</i> students study the wrong books. <i>Former president</i> George W. Bush
NAC-LOC	a <i>Boston</i> manufacturer, the <i>Chicago</i> office
NAC-TMP	The <i>Oct 12</i> editorial
NAC-TTL	He engraved a " <i>#1 Dad</i> " plaque
NP	Noun Phrase
NP	

NP-ADV	costing \$280 <i>a share</i> ; a dozen cases <i>a year</i> ; will be billed <i>several weeks after the expenditure</i> ; climbed <i>a solid 47%</i>
NP-BNF	pour <i>me</i> a cup of tea, buy <i>the cat</i> a present
NP-CLR	take <i>heart</i> ; follow <i>suit</i> ; cost <i>about \$2</i> ; Thank <i>Goodness!</i> ; take <i>several steps</i> ; make <i>it</i> past sth; etc.
NP-CLR-LOC	the house, located <i>about 50 yards from here</i> , was destroyed. (the only example)
NP-CLR-TMP	spend two days working and <i>two days</i> in the yard (the only example, and 'two days' from 'two days working' is labeled as NP-TMP-CLR)
NP-DIR	foolish to look <i>the other way</i> , sit <i>here</i> and wait
NP-EXT	sales grew <i>16%</i> ; funds increased <i>\$13 billion</i>
NP-HLN	<i>DISCIPLINARY PROCEEDINGS</i> against lawyers open to public in Illinois Dow shot up 23 points, in part due to buy programs generated by <i>stock-index arbitrage</i> (<i>Jacksonville, Florida</i>) Some Headline Here
NP-LGS	bills recently passed by <i>the House and Senate</i> , earnings reduced by <i>the sale of 4 million shares</i>
NP-LOC	Petrie stores, of <i>Secaucus, N.J.</i> ; Brown and Platt, <i>Chicago</i>
NP-LOC-CLR	changes sweeping <i>the East bloc</i> ; i once lived <i>there</i>
NP-LOC-HLN	
NP-LOC-PRD	the earthquake was <i>50 miles to the south</i> ; you're <i>right there</i>
NP-LOC-TPC-PRD	<i>here</i> are some of the major components
NP-MNR	he tries to have it <i>both ways</i> on the abortion issue; the bigotry of seeing things <i>only the Japanese way</i> the law should restrict citizens <i>as little as is consistent with good manners</i> .

NP-MNR-CLR	
NP-PRD	i swam in the lake – <i>lake eerie</i> , that is. although it is good, <i>Hollywood on the Hudson</i> it isn't.
NP-PRD-TPC	<i>issues discussed</i> were a, b, and c.
NP-PRD-TTL	he calls cotton “ <i>the fabric of our lives.</i> ” titled “ <i>Comments from Students,</i> ” it focuses on the real shame of college athletics
NP-SBJ	<i>the farmer</i> leaves. <i>both</i> deny wrongdoing.
NP-SBJ-TTL	“ <i>Feelings</i> ” is a good song.
NP-TMP	she gave up running <i>three times a week</i> in favor of playing golf. company a was acquired <i>last year</i> by company b. in a meeting <i>last tuesday</i> , we discussed stocks
NP-TMP-CLR	expires <i>November 16th</i> , ended <i>yesterday</i>
NP-TMP-HLN	
NP-TMP-PRD	it was <i>a long time</i> coming, it was <i>5:00pm</i> when i started working
NP-TPC	“ <i>no wonder the competition’s green with envy,</i> ” said bob; “ <i>wonderful!</i> ” said jim
NP-TTL	the company stopped selling “ <i>the Big Earl;</i> ” I approached “ <i>Mastergate</i> ” with trepidation
NP-TTL-PRD	a movie called “ <i>Marmelade</i> ”; a speech titled “ <i>Marmelade is Tasty</i> ”
NP-TTL-SBJ	“ <i>Baker’s Boys</i> ” is both bluesy and funny
NP-TTL-TPC	“ <i>Tivoli Motel,</i> ” I read on the sign.
NP-VOC	no, <i>darling</i> ; move over, <i>pornographic phone services; Hol- lywood</i> , you slay me.
NX	NP head marker
NX	Seems that this is intended to mark the head of a noun phrase; I couldn't find anything terribly meaningful

NX-TTL	
PP	Prepositional Phrase
PP	He lived <i>with her</i> .
	I'm out <i>of bed</i> .
	He had lots <i>of them</i> .
	He became angry <i>in return</i> .
	There is hope <i>of change</i> .
	He can live <i>with little pleasures</i> .
PP-BNF	He prints ads <i>exclusively for retailers</i> .
PP-CLR	beware <i>of my dog</i> , look <i>at me</i> , feed <i>on</i> , serve <i>as emcee</i> , talked <i>of the aftermath</i>
PP-CLR-LOC	used <i>in these strategies</i> , based <i>in Houston</i> ,
PP-CLR-TMP	expires <i>on Dec 31 1990</i>
PP-CLR-TPC	<i>Out of the mouths of revolutionaries</i> are coming words of moderation. <i>As factors contributing to the slowdown</i> , he cited ...
PP-DIR	The yield rose <i>to 5.38%</i> . Fanuc gained 100 <i>to 7,580</i> . We're not rushing <i>into anything</i> . Mr. Bush returned <i>to Washington</i> Saturday night.
PP-DIR-CLR	The family moves <i>to another house</i> at night. China exported 65 million pounds of mushrooms, valued at \$47 million, <i>to the U.S.</i> This will lead <i>to increased litigation</i> .
PP-DIR-PRD	People say they swim, and that may mean they've been <i>to the beach</i> . If he wants \$70K <i>out of me</i> , they have to take everything I have.
PP-DTV	They lied <i>to me</i> , They sent it <i>to the Senate</i>
PP-EXT	Kyocera advanced 80 yen <i>to 5,440</i> .

PP-HLN

PP-LGS Buyouts may be curbed *by two rules pending legislation.*

PP-LOC Most sleep *on the floor*, Reform starts *in the Pentagon*

PP-LOC-CLR We took it *on the chin*, Politics got *in the way*, Mr. Bass is based *in Ft. Worth*

PP-LOC-CLR-TPC *At the core of all this* stands a hotel.

PP-LOC-HLN

PP-LOC-MNR It should open up channels of communications with the Tigrean rebels *through neighboring Sudan.*

PP-LOC-PRD I'm *out of bed*. You're *in good company*. Help is *on the way.*

PP-LOC-PRD-TPC *Among the leading products* is a flu shot, *At the core* is a love for plants.

PP-LOC-TPC *Behind the posturing* lies a dispute.

PP-LOC-TPC-PRD *Among the new issues* was Massachusetts's debt.

PP-MNR He responded *in kind.*

This being typed *in a standing position.*

This will reduce spending *in a very effective fashion.*

PP-MNR-PRD The only way to find out is *by listening.*

PP-NOM He spent *between \$5 and \$6.*

PP-PRD That is *for the future.*

It looks *like a holiday.*

PP-PRD-LOC He remains *on the board* as director.

It's ironic that David Boren should be *in the center of this.*

PP-PRD-LOC-TPC *In the corner of the room* is a desk.

PP-PRD-TPC *Among the most upbeat* was Bobby.

PP-PRP The shop is closed *for a holiday.*

PP-PRP-CLR She was jailed in a child custody case *for refusing to reveal the whereabouts of her daughter.*

PP-PRP-PRD	The rise was <i>partly because of higher demand</i> .
PP-PUT	This put Mrs. Thatcher <i>in a bind</i> .
PP-SBJ	
PP-TMP	<i>In 1985</i> , it was warm; I'm done in <i>two minutes</i> .
PP-TMP-CLR	He reset opening arguments <i>for today</i> . The sentencing is set <i>for Jan 5</i> .
PP-TMP-PRD	The change is <i>since year-end</i> .
PP-TMP-TPC	<i>Starting in September</i> , the index started to slide.
PP-TPC	<i>Of 1500 people sent a questionnaire</i> , 951 replied.
PP-TPC-CLR	<i>With that authority</i> goes an accountability.
PP-TPC-LOC-PRD	<i>Among the possible suitors</i> is Italy's Fiat.
PP-TPC-PRD	<i>Along with the exodus of shopping</i> is an exodus of jobs.
PP-TTL	Tomorrow's " <i>On Sports</i> " will look at another aspect.
PP-TTL-PRD	The name of this column is " <i>On Sports</i> ".
PRN	Parenthetical
PRN	The alternative— <i>Giuliani</i> —is ghastly. The rest, <i>as they say</i> , is history.
PRT	Particle
PRT	He cashed <i>in</i> . Dream <i>on</i> .
PRT ADVP	Help the U.S. win <i>back</i> business.
QP	Quantifier Phrase
QP	He talked <i>about 20</i> minutes. <i>Not a</i> peso is offered.
RRC	Reduced Relative Clause
RRC	Everyone <i>at the Stick that day</i> started out as a spectator and ended up as a participant. There are still some uncertainties, <i>particularly regarding possible side effects</i> .
SBARQ	Direct question
SBARQ	<i>But who knows?</i>
SBARQ-HLN	<i>WHO'S NEWS??</i>

SBARQ-NOM	now the interest is in <i>what else can i do</i>
SBARQ-PRD	the question is, <i>what is stock worth?</i>
SBARQ-TPC	<i>why be a middleman?</i> asked joe
SBARQ-TTL	the old refrain, " <i>Who am I to judge</i> "
SBAR	Subordinate clause
SBAR	
SBAR-ADV	<i>as it turns out, John loves Mary</i> ; i can run <i>if i need to run.</i>
SBAR-ADV-TPC	<i>if profits don't improve</i> , we may need to close the company.
SBAR-CLR	i feel <i>as though i'm being watched</i> ; the agreement calls <i>for you to give me a lot of money</i>
SBAR-DIR	capital flows <i>where it is needed</i> ; go <i>where the money is</i>
SBAR-DIR-TPC	i hold that <i>wherever Mary goes</i> , John will follow.
SBAR-HLN	<i>WHO'S NEWS:</i>
SBAR-LOC	seems indistinguishable from SBAR-LOC-CLR
SBAR-LOC-CLR	my parents will stay <i>where they are.</i>
SBAR-LOC-PRD	that is <i>where i first met my wife.</i>
SBAR-MNR	they didn't play the game on saturday <i>as scheduled.</i>
SBAR-NOM	he has <i>what all publishers wish for.</i>
	i hate answering questions about <i>what would happen if we went to war.</i>
SBAR-NOM-LGS	they are put off by <i>what they consider to be restrictive investment regulations.</i>
SBAR-NOM-PRD	that is <i>what we did.</i>
SBAR-NOM-SBJ	<i>what is true for sheep</i> is true for goats as well.
SBAR-NOM-TPC	<i>whatever peopl want to buy</i> , i'll sell.
SBAR-PRD	our hope is <i>that the technique could identify diseased vessels.</i>
SBAR-PRD-TPC	<i>What counts</i> is the bottom line.

SBAR-PRP	the charge didn't affect net for the quarter, <i>as it was offset by tax benefits.</i> we expect a large market in the future, <i>so the long term it will be profitable.</i>
SBAR-PRP-PRD	that is <i>because John ran up the hill.</i>
SBAR-PUT	...put our resources <i>where they could do the most</i> ; put his money <i>where his mouth is</i>
SBAR-SBJ	he has made it clear <i>that the issue is important to him personally.</i> we want to make sure <i>we hold on to our existing customers.</i> <i>where they lag behind the Japanese</i> is in turning the inventiveness into increased production.
SBAR-TMP	i will be happy <i>when terms are fixed Oct. 26.</i>
SBAR-TMP-CLR	it didn't help <i>when she was charged with public drunkenness.</i>
SBAR-TMP-PRD	that was <i>before the tax reform made things more complicated.</i>
SBAR-TPC	he jailed them for several hours <i>after they defied his order;</i>
SBAR-TTL	<i>"When Harry Met Sally"; "When Irish Eyes Are Smiling"</i>
SINV	Sentence, Inverted
SINV	<i>"I am hungry," said Bob.</i> <i>Says Joe, "I am hungry, too."</i>
SINV-ADV	protected themselves against squalls in any area, <i>be it stocks, bonds, or real estate</i>
SINV-HLN	seems same as SINV; just used when it is a headline instead of a normal sentence
SINV-TPC	<i>Offsetting the lower stake in Lyondell were high crude oil prices, among other things.</i>

SINV-TTL	the children sang “ <i>Here Comes Santa Claus</i> ”
SQ	Sentence, Question (inverted yes/no, or the argument of a Wh)
SQ	How the hell <i>can you live with yourself?</i>
SQ-PRD	What gets by me every time is <i>has the milk expired?</i> Jimmy asked, “ <i>Can I go to the store?</i> ”
SQ-TPC	<i>Is that the forecast? Is the government really not helping anybody? Would I have done all those things?</i>
SQ-TTL	“ <i>Is Science, Or Private Gain, Driving Ozone Policy?</i> ” (article title)
S	Sentence, Declarative
S	
S-ADV	<i>A piece down</i> , the computer resigned. Investment bonds ended <i>1/4 point lower</i> . The company wouldn’t elaborate, <i>citing competitive reasons</i> .
S-CLF	<i>It is the total relationship that is important.</i>
S-CLF-TPC	“ <i>It’s not very often something like this comes up,</i> ” said Ron.
S-CLR	It helps <i>to be male</i> . The farmer stands <i>to go</i> .
S-CLR-ADV	Share prices closed lower in Paris, and <i>mixed</i> in Amsterdam.
S-HLN	<i>JAMAICA FIRES BACK</i>
S-LOC	At the end of the third quarter McDonald’s had 10K units operating <i>world-wide</i> .
S-MNR	Bonuses would be paid <i>based on playing time and performance</i> .
S-MNR-CLR	He began his career <i>peddling stock to individual investors</i> .
S-NOM	He apologizes for <i>sounding pushy</i> . They don’t flinch at <i>writing them</i> .

S-NOM-LGS	The insurance provided by <i>purchasing puts</i> is worthwhile. It was followed by <i>our driving to the nearest watering hole</i> .
S-NOM-PRD	That is <i>gilding the lily</i> .
S-NOM-SBJ	<i>Avoiding failure</i> is easy.
S-PRD	That lawsuit is <i>pending</i> . There is <i>more volatility to come</i> . It is <i>the New Journalism come to television</i> .
S-PRD-TPC	<i>Still to come</i> are issues by Monsanto.
S-PRP	Mr. Gargan favors simply giving money to the SEC <i>to hire more staff</i> .
S-PRP-CLR	Lotus Notes is designed <i>to sort e-mail sent within work groups</i> .
S-PRP-PRD	The manufacturers said 14.2% of their spending is designed to improve products, 17.5% is <i>to cut costs</i> , etc.
S-PRP-TPC	<i>To provide for the restructuring's costs</i> , Trinova took an after-tax charge.
S-SBJ	<i>To watch your child die</i> is an inhuman experience. “ <i>Work hard, play hard</i> ” is advice best taken with caution.
S-TMP	<i>Going into the fourth quarter</i> the sales comparison will be more difficult.
S-TPC	<i>The market changed</i> , he adds.
S-TPC-TMP	<i>A year ago you'd spend two days working and two days in the yard</i> , he recalls.
S-TTL	He stressed the “ <i>always working</i> ” theme.
S-TTL-PRD	The theme of the conference will be “ <i>take a pension fund manager to lunch</i> .”
S-TTL-SBJ	It's possible that “ <i>Look Who's Talking</i> ” isn't as entertaining as it seems.
UCP	Unlike Coordinated Phrase

UCP	This requires <i>regulatory and shareholder</i> approval. <i>consumer and other goods</i>
UCP-ADV	<i>Third and most important</i>
UCP-CLR	Stocks closed <i>lower but above intraday lows</i>
UCP-DIR	They moved <i>away from one thing and toward another.</i>
UCP-EXT	We will stay <i>through the skiing season or until the money runs out.</i>
UCP-LOC	The cuts will be made <i>half within Germany and half abroad.</i>
UCP-LOC-PRD	
UCP-MNR	It was mentioned <i>very briefly and in passing.</i>
UCP-PRD	Long dollar bonds were <i>flat to up 3/8 point.</i>
UCP-PRD-LOC	The SEC is <i>closer to the markets and in a good position to sing.</i>
UCP-PRP	Growers bred them <i>more for looks and to satisfy demands of long-term storage.</i>
UCP-TMP	<i>The next day or even an hour later, this year and in 1990</i>
UCP-TPC	
VP	Verb Phrase
VP	He <i>dies.</i>
VP-TPC	Also <i>baking a cake</i> is his mother.
VP-TTL	This newspaper's <i>Heard on the Street</i> column
WHADJP	Wh-adjective Phrase
WHADJP	<i>How strong</i> is Mr. Mohammad?
WHADV	Wh-adverb Phrase
WHADV	<i>How, Why, How quickly, etc.</i>
WHADV-TMP	<i>When</i>
WHNP	Wh-noun Phrase
WHNP	<i>Who, What, Whom, Which</i>
WHPP	Wh-prepositional Phrase

WHPP	<i>by how much, under what weight, for whom</i>
X	Unknown, Uncertain, or Unbracketable
X	<i>the closer they got, the more</i> the price rose the stock tumbled, to end <i>at</i> the earthquake <i>was</i> the crowd shouted, " <i>viva peace, viva.</i> " <i>c-</i> list item 3 i struggled <i>to</i> to eat my sandwich. i am married, <i>no children.</i> it was a funny time, <i>what</i> with the vietnam war and all. i was hungry <i>to begin with</i>
X-ADV	<i>the more extensive</i> the voir dire, the easier you make it. the more he muzzles his colleagues, <i>the more</i> leaks will pop up.
X-CLF	
X-DIR	earnings declined by \$120 million <i>last year's robust levels.</i>
X-EXT	exports from canada jumped 11% while imports from canada rose <i>only 2.7%</i>
X-HLN	
X-PUT	mr. bush's veto power puts him <i>a commanding position</i> in the narrowly divided house
X-TTL	

Appendix C

Modifying the Bikel Parser

The Collins-Bikel parser is a very heavily optimized chart parser. The algorithm for the parser is described in the appendices of Collins' PhD thesis [16], and remains the basis of Bikel's implementation. The main modification that is done is to add an `equivalentItems` list to each element in the chart and to store every item and link that would have been pruned away either by the search or just by virtue of the dynamic program (which is looking for the top-scoring parse). In the Bikel parser, this change should occur in the `add` method of the `Chart` class.

Following this, any calculations that need to be done (to compute inside and outside probabilities, for example) can be done in `Decoder.parse` once the entire forest has been computed. This is also the appropriate point for the parse forest to be emitted.

Appendix D

A Massively Simpler Alignment Model

D.1 Motivation

Assuming that a random variable is Gaussian is a natural choice when the distribution is unknown, because Gaussians have many nice properties. In particular, they are an especially good choice for use in EM algorithms since the maximum likelihood estimates for a Gaussian can be written in closed form. Furthermore, assuming that a variable is Gaussian is often a good approximation due to the Central Limit Theorem, which states that a sum of independent and identically distributed variables with finite mean and variance tends to be Gaussian as the number of addends approaches infinity, and particularly due to the fuzzy Central Limit Theorem, which states that data influenced by many independent sources of noise are roughly normally distributed [48].

For instance, the number of words in English sentences in the EUROPARL corpus, depicted in Figure D-1, is roughly Gaussian. One could explain this based on that fuzzy Central Limit Theorem by imagining a number of independent sources of noise that would influence the length of an English sentence, including such things as

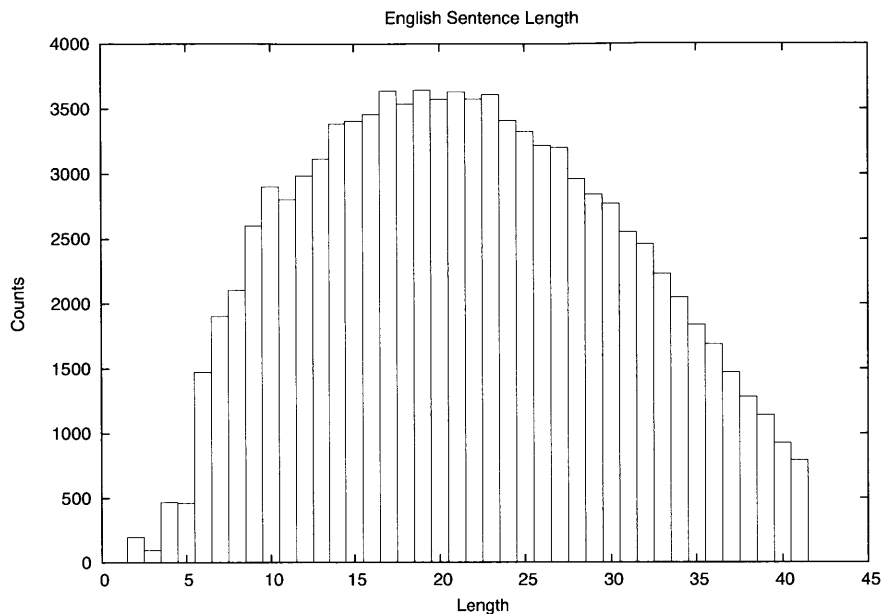
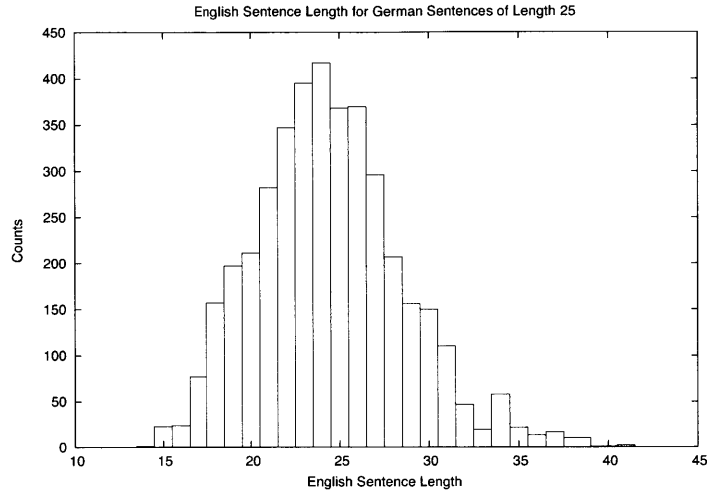
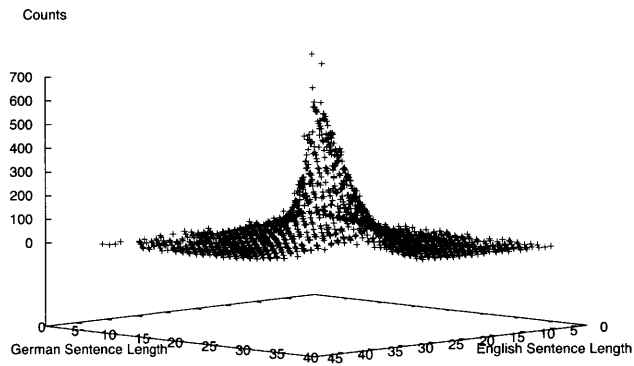


Figure D-1: The length of English sentences drawn from the English translation of the proceedings of the European Parliament appears to have a Gaussian distribution.

the connotation of certain phrases (which influences the author's choice and thereby influences the length of the sentence) and the author's desire to be precise. In fact, the distribution of English sentence lengths for a fixed German sentence length also appears to be Gaussian (see Figure D-2). Most relevant to us, however, is that when one trains an IBM2 model to translate from German to English, the distribution of the index of the word in a, say, 25-word English sentence that the 13th German word of a 25-word German sentence is aligned to also suggests the Gaussian shape, as Figure D-3 shows. This can again be argued using the fuzzy Central Limit Theorem: all else equal, we imagine that a word is most likely to remain in the same relative spot within a sentence; for each transformation that would move it to one side, we imagine there is another transformation that is likely to move it to the other side. Ultimately, the reasons for a translator's choices are innumerable and will be written off as noise here.



(a)



(b)

Figure D-2: (a) The length of English sentences whose German translations are twenty-five words long appears to have a Gaussian distribution. (b) The length of sentences in English and German appears to have a bivariate Gaussian distribution. All things are indeed Gaussian.

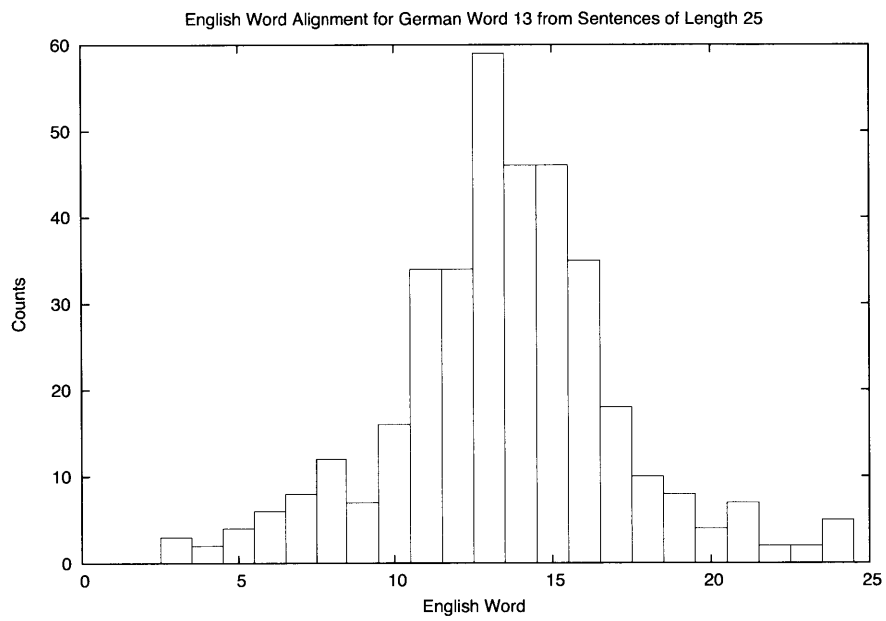


Figure D-3: Our motivation: The indices of aligned words can be approximated by a Gaussian.

One might argue that the best way to model a random variable of unknown distribution is by simply modeling a probability for each of its possible values. This technique is certainly flexible; unfortunately, it is a byword in the machine learning community that one must pay for added flexibility with more training data to avoid over-fitting. Given that even the most sophisticated machine translation models are far from perfect and that they demand massive amounts of training data, our approach is to begin by further constraining existing models instead of creating more flexible ones.

In accordance to this intuition and our guess that alignments look Gaussian, instead of modeling every possible value for each alignment as a separate probability, we will model the alignment probabilities by a single Gaussian. The mathematical formulation for IBM Model 2 is changed by the addition of the second line:

$$P(\mathbf{f}, \vec{a} | \mathbf{e}) = \prod_{j=1}^m T(\mathbf{f}_j | \mathbf{e}_{a_j}) D(a_j | j, \ell, m)$$

$$D(a_j | j, \ell, m) = \frac{f_{\mathcal{N}}(a_j | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}{\sum_{i=0}^{\ell} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m})},$$

where $f_{\mathcal{N}}(\cdot | \mu, \sigma)$ denotes the density of the gaussian with mean μ and standard deviation σ , $1/(\sqrt{2\pi}\sigma) \cdot \exp(-(\mu - \cdot)^2/(2\sigma^2))$. That is to say, we replace $D(a_j | j, \ell, m)$ for $a_j = 1, \dots, \ell$ with $\mu_{j,\ell,m}$ and $\sigma_{j,\ell,m}$ as our parameters. This typically results in fewer than 10% as many alignment parameters as IBM Model 2.

D.1.1 Algorithm

IBM Model 2 is a multinomial model; consequently, the EM updates are very easy to compute. The type of extension we are discussing (“parametrizing the parameters”) corresponds to a simple addition to the algorithm in this case. The full derivation is discussed in the appendix.

Let’s begin by casting the model in the standard form: each observation consists of a sentence pair (\mathbf{e}, \mathbf{f}) generated from a hidden alignment $\vec{a} \in \{0, \dots, \ell\}^m$. The IBM

Model 2 probability is:

$$\begin{aligned}
P(\mathbf{f}, \vec{a} | \mathbf{e}; T, D) &= \prod_{j=1}^m T(\mathbf{f}_j | \mathbf{e}_{a_j}) D(a_j | j, \ell, m) \\
&= \prod_{f,e} T(f | e)^{\text{Count}_{f,e}(\mathbf{e}, \mathbf{f}, \vec{a})} \times \prod_{i,j,\ell,m} D(i | j, \ell, m)^{\text{Count}_{i,j,\ell,m}(\mathbf{e}, \mathbf{f}, \vec{a})}
\end{aligned}$$

where

$$\begin{aligned}
\text{Count}_{f,e}(\mathbf{e}, \mathbf{f}, \vec{a}) &:= \sum_{j=1}^m \delta(f, \mathbf{f}_j) \delta(e, \mathbf{e}_{a_j}) \\
\text{Count}_{i,j,\ell,m}(\mathbf{e}, \mathbf{f}, \vec{a}) &:= \delta(\ell, |\mathbf{e}|) \delta(m, |\mathbf{f}|) \delta(a_j, i),
\end{aligned}$$

and $\delta(\cdot, \cdot)$ denotes the Kronecker delta function (1 if the two parameters are equal, and 0 otherwise). The probability model we propose is:

$$\begin{aligned}
P(\mathbf{f}, \vec{a} | \mathbf{e}; T, \mu, \sigma) &= \prod_{f,e} T(f | e)^{\text{Count}_{f,e}(\mathbf{e}, \mathbf{f}, \vec{a})} \times \prod_{j,\ell,m} N(j, \ell, m)^{\text{Count}_{0,j,\ell,m}(\mathbf{e}, \mathbf{f}, \vec{a})} \\
&\quad \times \prod_{i,j,\ell,m} ((1 - N(j, \ell, m)) \cdot f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}))^{\text{Count}_{i,j,\ell,m}(\mathbf{e}, \mathbf{f}, \vec{a})}.
\end{aligned}$$

Here, $N(j, \ell, m)$ is the probability that the j th French word aligns to the NULL English word (we do not wish this probability to be modeled by some slot in the Gaussian, because no position corresponds logically to the NULL word; thus we separate it in this fashion). Note that the $f_{\mathcal{N}}$ term is unnormalized—that is, the model is deficient under this framework for the sake of mathematical convenience. Model 2 clearly falls under the multinomial framework; consequently, this new model falls under the parametrized multinomial framework described in the appendix.

Therefore, the addition we made to the mathematical formulation in the last section results in a single, easy change to the IBM Model 2 algorithm: after each EM iteration, the old D values are replaced by their maximum-likelihood Gaussian

counterparts. That is to say, we transform the D 's obtained from Model 2 as follows:

$$\begin{aligned}\mu_{j,\ell,m} &= \frac{1}{m} \sum_{i=1}^m i \cdot D(i | j, \ell, m) \\ \sigma_{j,\ell,m} &= \sqrt{\frac{1}{m} \sum_{i=1}^m i^2 \cdot D(i | j, \ell, m) - \mu_{j,\ell,m}^2}\end{aligned}$$

Adding the constraint that the alignment variables are samples of univariate Gaussians corresponds to an altogether simple change to the algorithm; the full algorithm is shown in Figure D-4. (We let the 0th word, e_0 , of every English sentence be the NULL word to simplify the notation.)

Our argument for this Gaussian approximation is only based on the idea that the alignment variables look like one-dimensional Gaussian densities and that Gaussians are easy to deal with; this is clearly not the only approximation that satisfies these properties. In fact, there are two other obvious choices based on the Gaussian: “truncated” Gaussians and “integrated” Gaussians. Here they are in math beneath our original formulation:

$$\begin{aligned}\text{vanilla: } D(a_j | j, \ell, m) &= f_{\mathcal{N}}(a_j | \mu_{j,\ell,m}, \sigma_{j,\ell,m}), \\ \text{truncated: } D(a_j | j, \ell, m) &= \frac{f_{\mathcal{N}}(a_j | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}{\sum_i f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}, \\ \text{and integrated: } D(a_j | j, \ell, m) &= \frac{\int_{a_{j-1}}^{a_j} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) di}{\int_0^\ell f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) di}.\end{aligned}$$

In both the vanilla and truncated models, the alignment variable a_j can take on the values $1, \dots, \ell$; the difference is that, for mathematical simplicity, the vanilla version is deficient, assigning probabilities to values outside this range (i.e., that $\sum_{a_j} D(a_j | j, \ell, m) \lesssim 1$). Although we obviously cannot align words outside of the sentence, we allow the algorithm to assign non-zero probability mass to those alignments; we recover by normalizing over legal alignments afterward. We can instead make this restriction in the algorithm itself. Unfortunately, the M-step of the EM algorithm is no longer a beautiful closed form, but instead requires numerical opti-

```

Initialize  $t(f|e)$  and  $D(i|j, \ell, m)$ 
do:
· zero  $t'(f|e)$  and  $D'(i|j, \ell, m)$ 
· for  $(\mathbf{e}, \mathbf{f})$  in corpus:
· ·  $m = |\mathbf{f}|, \ell = |\mathbf{e}|$ 
· · for  $j=1 \dots m$ :
· · · for  $i=0 \dots \ell$ :
· · · ·  $a_i = t(\mathbf{f}_j | \mathbf{e}_i) \cdot D(i|j, \ell, m)$ 
· · · ·  $a_i = a_i / (\sum_{i'} a_{i'})$ 
· · · · for  $i=0 \dots \ell$ :
· · · · ·  $t'(\mathbf{f}_j | \mathbf{e}_i) = t(\mathbf{f}_j | \mathbf{e}_i) + a_i$ 
· · · · ·  $D'(i|j, \ell, m) = D(i|j, \ell, m) + a_i$ 
·  $t'(f|e) = t'(f|e) / (\sum_{f'} t'(f'|e))$ 
·  $D'(i|j, \ell, m) = D'(i|j, \ell, m) / (\sum_{i'} D'(i'|j, \ell, m))$ 
·  $t = t', D = D'$ 
· for  $\ell, m$ :
· · for  $j=1 \dots m$ :
· · ·  $\mu = \frac{1}{\ell} \left( \sum_{i=1}^{\ell} D(i|j, \ell, m) \cdot i \right) / (1 - D(0|j, \ell, m))$ 
· · ·  $\sigma = \frac{1}{\ell} \left( \sum_{i=1}^{\ell} D(i|j, \ell, m) \cdot (i - \mu)^2 \right) / (1 - D(0|j, \ell, m))$ 
· · · for  $i=1 \dots \ell$ :
· · · ·  $D(i|j, \ell, m) = \exp(-(i - \mu)^2 / 2\sigma)$ 
· · · ·  $D(i|j, \ell, m) = D(i|j, \ell, m) \cdot (1 - D(0|j, \ell, m)) / (\sum_{i'=1}^{\ell} D(i'|j, \ell, m))$ 
until convergence

```

Figure D-4: The Vanilla Gaussian Algorithm.

mization. Although this is still tractable, it is undesirable and the results are not sufficiently improved to warrant this computational burden (our experiments showed virtually identical results to the vanilla model).

Likewise, the integrated version is very attractive intuitively, but the shape of the density is so close to that of the vanilla edition that any improvements are minute and are outweighed by the added computational complexity, as the Q -function in EM must once again be optimized by a numerical optimization technique. For completeness, we explain how one would perform such optimization; in this case, just as in the vanilla algorithm, the change only amounts to “fitting” (in the maximum-likelihood sense) the appropriate density to the intermediate D -values after each EM step. We can compute the gradient of the likelihood function in both of these cases, so we optimize it using a gradient optimization technique. In our experience, the likelihood functions tend to have long, narrow valleys, so we find that optimizing using conjugate gradient descent is faster than just using steepest descent. (Note again that we do not include the NULL parameter $D(0 | \cdot, \cdot, \cdot)$ in the fit.)

We have to fit a curve to each set of D -values, so fix j , ℓ , and m . Define $\bar{D} = D/(1 - D_0)$ to reflect the fact that we are not including the NULL word position in our model and to thus further ease the notational burden. We wish to maximize the log-likelihood, so let’s begin by writing it down:

$$\begin{aligned}
\mathcal{L}_{\text{truncated}} &= \sum_{i=1}^{\ell} \bar{D}(i | j, \ell, m) \log \frac{f_{\mathcal{N}}(a_j | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}{\sum_i f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m})} \\
&= \sum_{i=1}^{\ell} \bar{D}(i | j, \ell, m) \cdot \frac{-(i - \mu)^2}{2\sigma^2} - \log \left(\sum_{i=1}^{\ell} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) \right) \\
\mathcal{L}_{\text{integrated}} &= \sum_{i=1}^{\ell} \bar{D}(i | j, \ell, m) \log \frac{\int_{a_{j-1}}^{a_j} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) di}{\int_0^{\ell} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) di} \\
&= \sum_{i=1}^{\ell} \bar{D}(i | j, \ell, m) \cdot \log(F_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - F_{\mathcal{N}}(i - 1 | \mu_{j,\ell,m}, \sigma_{j,\ell,m})) \\
&\quad - \log(F_{\mathcal{N}}(\ell | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - F_{\mathcal{N}}(0 | \mu_{j,\ell,m}, \sigma_{j,\ell,m})),
\end{aligned}$$

where $F_{\mathcal{N}}$ denotes the cumulative distribution function. We evaluate the gradient:

$$\begin{aligned}
\frac{\partial \mathcal{L}_{\text{truncated}}}{\partial \mu} &= \sum_{i=1}^{\ell} \bar{D}(i|j, \ell, m) \cdot \frac{i - \mu}{\sigma^2} - \frac{\sum_{i=1}^{\ell} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) \cdot \frac{i - \mu}{\sigma^2}}{\sum_{i=1}^{\ell} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m})} \\
\frac{\partial \mathcal{L}_{\text{truncated}}}{\partial \sigma^2} &= \sum_{i=1}^{\ell} \bar{D}(i|j, \ell, m) \cdot \frac{-(i - \mu)^2}{2\sigma^4} - \frac{\sum_{i=1}^{\ell} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) \cdot \frac{-(i - \mu)^2}{2\sigma^4}}{\sum_{i=1}^{\ell} f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m})} \\
\frac{\partial \mathcal{L}_{\text{integrated}}}{\partial \mu} &= \sum_{i=1}^{\ell} \bar{D}(i|j, \ell, m) \cdot \sqrt{2\sigma^2} \cdot \frac{f_{\mathcal{N}}(i - 1 | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}{F_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - F_{\mathcal{N}}(i - 1 | \mu_{j,\ell,m}, \sigma_{j,\ell,m})} \\
&\quad - \sqrt{2\sigma^2} \cdot \frac{f_{\mathcal{N}}(0 | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - f_{\mathcal{N}}(\ell | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}{F_{\mathcal{N}}(\ell | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - F_{\mathcal{N}}(0 | \mu_{j,\ell,m}, \sigma_{j,\ell,m})} \\
\frac{\partial \mathcal{L}_{\text{integrated}}}{\partial \sigma^2} &= \sum_{i=1}^{\ell} \bar{D}(i|j, \ell, m) \cdot \frac{1}{4\sigma^2} \cdot \frac{f_{\mathcal{N}}(i - 1 | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - f_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}{F_{\mathcal{N}}(i | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - F_{\mathcal{N}}(i - 1 | \mu_{j,\ell,m}, \sigma_{j,\ell,m})} \\
&\quad - \frac{1}{4\sigma^2} \cdot \frac{f_{\mathcal{N}}(0 | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - f_{\mathcal{N}}(\ell | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}{F_{\mathcal{N}}(\ell | \mu_{j,\ell,m}, \sigma_{j,\ell,m}) - F_{\mathcal{N}}(0 | \mu_{j,\ell,m}, \sigma_{j,\ell,m})}
\end{aligned}$$

We apply the conjugate gradient descent algorithm to the D values [43]:

Given $D \in \mathbb{R}^m$

- Select $\vec{x}_0 = (\mu, \sigma^2) \in \mathbb{R}^2$ at random
- $i = 0$, $\vec{g}_0 = \nabla \mathcal{L}(\vec{x}_0)$, $\vec{h}_0 = -\vec{g}_0$
- **do**:
- · $\lambda_i = \arg \min_{\lambda \geq 0} \mathcal{L}(\vec{x}_i + \lambda_i \vec{h}_i)$
- · $\vec{x}_{i+1} = \vec{x}_i + \lambda_i \vec{h}_i$
- · $\vec{g}_{i+1} = \nabla \mathcal{L}(\vec{x}_{i+1})$
- · $\gamma_i = (\vec{g}_{i+1} - \vec{g}_i) \cdot \vec{g}_{i+1} / \|\vec{g}_i\|^2$
- · $\vec{h}_{i+1} = -\vec{g}_{i+1} + \gamma_i \vec{h}_i$
- · $i = i + 1$
- **until** convergence

return \vec{x}_i

D.1.2 Implementation and Evaluation

The dataset used to evaluate our system was the standard EUROPARL set of Köhn *et al.* We used German-English as a representative language pair, as translation is neither especially easy nor especially difficult [29]. The data was aligned at the sentence-level using the standard tools and sentences of vastly differing lengths were removed. Finally, we trained the system on the data and produced Viterbi alignments for each sentence pair. These alignments were output to the Köhn phrase-based system, which itself produced a phrase dictionary. This dictionary was applied to the Pharoah decoder (with a language model trained on the entire available training set). Finally, we applied the system to the standard test set chosen by Köhn in [29]. The resulting translations were compared to the human-translated reference using the BLEU metric of [40].

The BLEU metric is a standard method for evaluating machine translation system performance by comparing translations to one or many human translations. The translations are compared by precision and recall on n -grams of successively greater length; the BLEU score typically refers to a smoothed 4-gram comparison; mathematically, it can be described by the following formula:

$$\text{BLEU} = e^{I_{c \leq r} \cdot (1-r/c)} \cdot \sqrt[4]{p_1 p_2 p_3 p_4},$$

where r is the length of the reference corpus, c is the total length of the candidate translation produced by the system being evaluated, r is sum of the lengths of the reference sentences that most closely match the lengths of the candidate sentences, $I_{c \leq r}$ is 1 if $c \leq r$ and 0 otherwise, and p_j refers to the j -gram precision of the test set.

We evaluated our technique using the EUROPARL corpus [29] and the applied a BLEU scorer to our model's output on the standard section 24 test set. Our results (when training is done on the full data set) are shown in Table 1. Our results are clearly very competitive with IBM Model 2.

Considering the small number of parameters in the one-dimensional gaussian

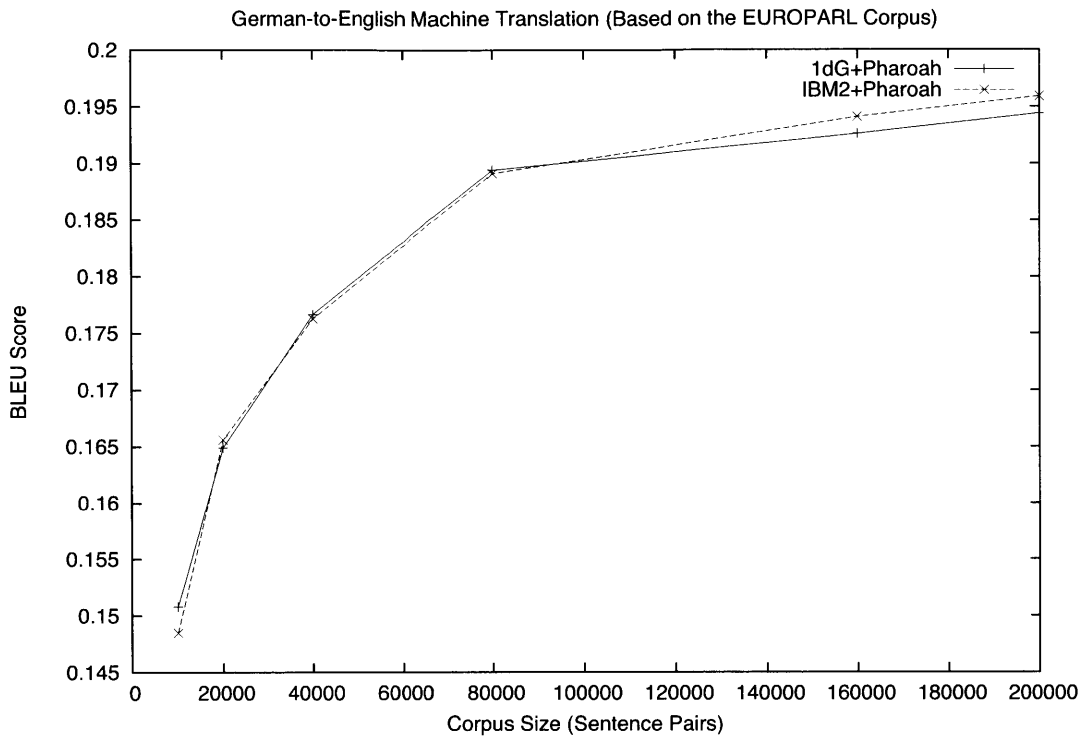


Figure D-5: The performance of the alignments induced by the one-dimensional gaussian model is similar to those induced by the IBM-2 model. This graph shows the BLEU metric of the two models when applied to the EUROPARL training data and standard test set for the German-English language pair.

model, intuition suggests that it should converge to its limiting BLEU score with less data. If this is true, it is an insignificant effect, as Figure D-5 shows; we believe that this is due to the still overwhelming number of translation parameters in the model. The performance of the one-dimensional gaussian model is, in fact, indistinguishable from that of Model 2.

Bibliography

- [1] Second international workshop on evaluating word sense disambiguation systems. <http://www.senseval.org/>, July 2001.
- [2] Danit Ben-Ari, Daniel M. Berry, and Mori Rimon. Translational ambiguity rephrased. In *Proceedings of Second International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Language*, pages 257–262, Pittsburgh, PA, 1988.
- [3] Ann Bies, Mark Ferguson, Karen Katz, and Robert MacIntyre. *Bracketing Guidelines for Treebank II Style Penn Treebank Project*. Unpublished, January 1995.
- [4] Dan Bikel. A multilingual statistical parsing engine. <http://www.cis.upenn.edu/~dbikel/software.html#stat-parser>.
- [5] Thorsten Brants, Ashok C. Popat, Peng Xu, Franz J. Och, and Jeffrey Dean. Large language models in machine translation. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 858–867, Prague, Czech Republic, June 2007. Association for Computational Linguistics.
- [6] Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, Jennifer C. Lai, and Robert L. Mercer. An estimate of an upper bound for the entropy of english. *Computational Linguistics*, 18(1), 1992.

- [7] Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. The mathematics of statistical machine translation: Parameter estimation. *Association for Computational Linguistics*, 19(2):263–311, 1993.
- [8] Chris Callison-Burch. Linear b system description for the 2005 nist mt evaluation exercise. In *Proceedings of Machine Translation Evaluation Workshop*, 2005.
- [9] Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. (meta-) evaluation of machine translation. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 136–158, Prague, Czech Republic, June 2007. Association for Computational Linguistics.
- [10] Chris Callison-Burch, Miles Osborne, and Philipp Koehn. Re-evaluating the role of bleu in machine translation research. In *Proceedings of the Eleventh Conference of the European Chapter of the Association for Computational Linguistics*, pages 249–256, Trento, Italy, 2006.
- [11] Marine Carpuat and Dekai Wu. Word sense disambiguation vs. statistical machine translation. In *Proceedings of the 43rd Annual meeting of the Association for Computational Linguistics*, pages 387–394, Ann Arbor, July 2005.
- [12] Marine Carpuat and Dekai Wu. Improving statistical machine translation using word sense disambiguation. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 61–72, Prague, Czech Republic, June 2007. Association for Computational Linguistics.
- [13] David Carter. The treebanker: a tool for supervised training of parsed corpora. In *Computational Environments for Grammar Development and Linguistic Engineering*. Association for Computational Linguistics, July 1997.
- [14] Yee Seng Chan, Hwee Tou Ng, and David Chiang. Word sense disambiguation improves statistical machine translation. In *Proceedings of the 45th Annual Meet-*

ing of the Association of Computational Linguistics, pages 33–40, Prague, Czech Republic, June 2007. Association for Computational Linguistics.

- [15] Kenneth W. Church and Eduard H. Hovy. Good applications for crummy machine translation. *Machine Translation*, 8:239–258, 1993.
- [16] Michael Collins. *Head-Driven Statistical Models for Natural Language Processing*. PhD thesis, University of Pennsylvania, 1999.
- [17] Loïc Dugast, Jean Senellart, and Philipp Koehn. Statistical post-editing on sysstrans rule-based translation system. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 220–223, Prague, June 2007. Association for Computational Linguistics.
- [18] Andreas Fink and Aljoscha C. Neubauer. Speed of information processing, psychometric intelligence and time estimation as an index of cognitive load. *Personality and Individual Differences*, 30:1009–1021, 2001.
- [19] Sandiway Fong and Robert C. Berwick. New approaches to parsing conjunctions using prolog. 1985.
- [20] Ryan Gabbard, Mitchell Marcus, and Seth Kulick. Fully parsing the penn treebank. In *Proceedings of the Human Language Technology Conference of the north American Chapter of the ACL*, pages 184–191, New York, NY, 2006.
- [21] Ralf Herbrich and Thore Graepel. Trueskill(tm): A bayesian skill rating system. 2006.
- [22] F. Jelinek, R. L. Mercer, L. Bahl, and J. Baker. Interpolated estimation of Markov source parameters from sparse data. In *Pattern Recognition in Practice: Proceedings, International Workshop*. Pattern Recognition in Practice, May 1980.
- [23] Frederick Jelinek. Up from trigrams! the struggle for improved language models. In *Proceedings of the Second European Conference on Speech and Technology*,

pages 1037–1040, Genova, September 1991. International Speech Communication Association.

- [24] Frederick Jelinek. *Statistical methods for speech recognition*. MIT Press, Cambridge, MA, USA, 1997.
- [25] Howard Johnson, Joel Martin, George Foster, and Roland Kuhn. Improving translation quality by discarding most of the phrasetable. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 967–975, Prague, Czech Republic, June 2007. Association for Computational Linguistics.
- [26] Martin Kay. *The MIND System*, pages 155–188. Algorithmics Press, New York, 1973.
- [27] Martin Kay. Machine translation will not work. In *Proceedings of the 24th Annual meeting of the Association for Computational Linguistics*, page 268, New York, July 1986. Association for Computational Linguistics.
- [28] Philipp Koehn. Pharaoh: a Beam Search Decoder for Phrase-Based Statistical Machine Translation Models. In *The 6th Conference of the Association for Machine Translation in the Americas*, pages 115–124, Washington DC, September 2004. Association for Machine Translation in the Americas.
- [29] Philipp Koehn. Europarl: A Parallel Corpus for Statistical Machine Translation. In *The Tenth Machine Translation Summit*, pages 79–86, Phuket, Thailand, September 2005. Asia-Pacific Association for Machine Translation.
- [30] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion*

Volume Proceedings of the Demo and Poster Sessions, pages 177–180, Prague, Czech Republic, June 2007. Association for Computational Linguistics.

- [31] Philipp Koehn, Franz Josef Och, and Daniel Marcu. Statistical phrase-based translation. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 48–54, Edmonton, AB Canada, May–June 2003. Association for Computational Linguistics.
- [32] Beth Levin. *English Verb Classes And Alternations: A Preliminary Investigation*. University of Chicago Press, Chicago, 1993.
- [33] Daniel Marcu and William Wong. A phrase-based, joint probability model for statistical machine translation. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing*, pages 133–139, Philadelphia, July 2002. Association for Computational Linguistics.
- [34] Mitchell Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert Macintyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. The Penn Treebank: Annotating Predicate Argument Structure. pages 114–119, March 1994.
- [35] Hiroshi Maruyama, Hideo Watanabe, and Shiho Ogino. An interactive japanese parser for machine translation. In *International Conference On Computational Linguistics*, pages 257–262, Stanford, CA, 1990.
- [36] Barack Obama. A strategy for american innovation. September 2009.
- [37] Franz Josef Och. Minimum error rate training in statistical machine translation. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 160–167, Sapporo, Japan, July 2003. Association for Computational Linguistics.
- [38] Franz Josef Och and Hermann Ney. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51, 2003.

- [39] Kishore Papineni, Salim Roukos, Todd Ward, John Henderson, and Florence Reeder. Corpus-based comprehensive and diagnostic mt evaluation: Initial arabic, chinese, french, and spanish results. In *Proceedings of HLT 2002, Second International Conference on Human Language Technology Research*, pages 132–137, San Francisco, March 2002. Association for Computational Linguistics.
- [40] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, July 2002. Association for Computational Linguistics.
- [41] John R. Pierce. Whither speech recognition? pages 1049–1051, 1949.
- [42] John R. Pierce, John B. Carroll, Eric P. Hamp, David G. Hays, Charles F. Hockett, Anthony G. Oettinger, and Alan Perlis. *Language and Machines: Computers in Translation and Linguistics*. National Academy of Sciences, Washington, D.C., 1966.
- [43] E. Polak. *Optimization: Algorithms and Consistent Approximations*. Springer, New York, 1997.
- [44] Matt Post and Daniel Gildea. Parsers as language models for statistical machine translation. In *Proceedings of the Eighth AMTA Conference*, pages 172–181, Hawaii, October 2008. Association for Computational Linguistics.
- [45] Benoit Thouin. *The METEO System*, pages 39–44. 1981.
- [46] Masaru Tomita. Disambiguating grammatically ambiguous sentences by asking. In *Proceedings of the 22nd Annual Meeting of the Association for Computational Linguistics*, pages 475–480, Stanford, CA, 1984.
- [47] Warren Weaver. *Translation*, pages 15–23. Technology Press of the Massachusetts Institute of Technology, Cambridge, MA, 1949.

[48] Eric W. Weisstein. Central limit theorem. <http://mathworld.wolfram.com/CentralLimitTheorem.html>.