

Translation Quality Indicators for Pivot-based Statistical MT

Michael Paul and Eiichiro Sumita

National Institute of Information and Communications Technology

MASTAR Project

Kyoto, Japan

michael.paul@nict.go.jp

Abstract

Recent research on multilingual statistical machine translation focuses on the usage of *pivot languages* in order to overcome resource limitations for certain language pairs. This paper provides new insights into what factors make a good pivot language and investigates the impact of these factors on the overall pivot translation performance. Pivot-based SMT experiments translating between 22 Indo-European and Asian languages were used to analyze the impact of eight factors (*language family, vocabulary, sentence length, language perplexity, translation model entropy, reordering, monotonicity, engine performance*) on pivot translation performance. The results showed that 81% of system performance variations can be explained by these factors.

1 Introduction

The translation quality of statistical machine translation (SMT) approaches heavily depends on the amount and coverage of bilingual language resources available to train the statistical models. There exist several data collection initiatives¹ amassing and distributing large amounts of textual data. For frequently used language pairs like *French-English*, large text data sets are readily available. However, for less frequently used language pairs only a limited amount of bilingual resources are available, if any at all.

In order to overcome language resource limitations, recent research on SMT focuses on the usage of *pivot languages* (de Gispert and Marino, 2006; Utiyama and Isahara, 2007; Wu and Wang, 2007; Bertoldi et al., 2008). Instead of a direct translation between two languages where only a

limited amount of bilingual resources is available, the *pivot translation* approach makes use of a third language that is more appropriate due to the availability of more bilingual corpora and/or its relatedness towards either the source or the target language. For most recent research efforts, *English* is the pivot language of choice due to the richness of available language resources. However, recent research on pivot translation has shown that the usage of non-English pivot languages can improve translation quality for certain language pairs (Paul et al., 2009; Leusch et al., 2010).

Concerning the contribution of aspects of different language pairs on the quality of machine translation, (Birch et al., 2008) identified three features (*morphological complexity, amount of reordering, historical relatedness*) for predicting success of MT in translations between the official languages of the European Union. Moreover, (Koehn et al., 2009) investigated an additional feature (*translation model complexity*) using the JRC-Aquis corpus covering not only Indo-European languages, but also one semitic and three Finno-Ugric languages.

This paper differs from previous research in the following aspects: we focus on the framework of *pivot translation*, where a target language translation of a source language input is obtained through an intermediate (*pivot*) language, investigate what factors make a good pivot language and what impact these factors have on the overall translation quality of language pairs not only including Indo-European languages, but also a large variety of Asian languages. In Section 2, we report on pivot-based SMT experiments translating between 22 Indo-European as well as Asian languages in order to provide new insights into how much language diversity affects the translation performance of pivot translation approaches. In Section 3, eight factors (*language family, vocabulary, sentence length, language perplexity,*

¹LDC: <http://www ldc.upenn.edu>, ELRA: <http://www.elra.info>

translation model entropy, reordering, monotonicity, engine performance) are investigated to determine the significance of each factor in predicting translation quality using linear regression analysis.

2 Pivot Translation

Pivot translation is a translation from a source language (SRC) to a target language (TRG) through an intermediate *pivot* (or *bridging*) language (PVT). Within the SMT framework, various coupling strategies like *cascading*, *phrase-table composition*, and *pseudo-corpus generation* have been proposed. For the experiments reported in this paper, we utilized the *cascading* approach because it is computational less expensive, but still performs comparably well compared to the other, more sophisticated pivot translation approaches. Pivot translation using the *cascading* approach requires two translation engines where the first engine translates the source language input into the pivot language and the second engine takes the obtained pivot language output as its input and translates it into the target language. Given N languages, a total of $2*N*(N-1)$ SMT engines have to be built in order to cover all $N*(N-1)*(N-2)$ SRC-PVT-TRG language pair combinations.

The importance of translation quality factors in pivot translation are investigated using the multilingual *Basic Travel Expressions Corpus* (BTEC), which is a collection of sentences that bilingual travel experts consider useful for people going to or coming from another country (Kikui et al., 2006). The sentence-aligned corpus consists of 160k sentences pairs covering 22 Indo-European and Asian languages which belong to a variety of language families including *Germanic* (da,de,en,nl), *Romance* (es,fr,it,pt,ptb), *Slavic* (pl,ru), *Indo-Iranian* (hi), *Semitic* (ar), *Austronesian* (id,ms,tl), *Tai* (th), *Mon-khmer* (vi), and *Sinitic* (zh,zht) languages. The corpus statistics are summarized in Table 1, where *Voc* specifies the vocabulary size and *Len* the average sentence length of the respective data sets. These languages differ largely in word order (*Order*: subject-object-verb (SOV), subject-verb-object (SVO), verb-subject-object (VSO)), segmentation unit (*Unit*: phrase, word, none), and degree of inflection (*Inflection*: high, moderate, light). Very similar characteristics can be seen for *Indo-European* languages and for certain subsets of *Asian languages* (ja, ko; id, ms). In addition,

Table 1: Language Resources

(European Languages)

Language		Voc	Len	Order	Unit	Inflection
Danish	da	26.5k	7.2	SVO	word	high
German	de	25.7k	7.1	SVO	word	high
English	en	15.4k	7.5	SVO	word	moderate
Spanish	es	20.8k	7.4	SVO	word	high
French	fr	19.3k	7.6	SVO	word	high
Hindi	hi	33.6k	7.8	SOV	word	high
Italian	it	23.8k	6.7	SVO	word	high
Dutch	nl	22.3k	7.2	SVO	word	high
Polish	pl	36.4k	6.5	SVO	word	high
Portuguese	pt	20.8k	7.0	SVO	word	high
Brazilian Portuguese	ptb	20.5k	7.0	SVO	word	high
Russian	ru	36.2k	6.4	SVO	word	high

(Asian Languages)

Language		Voc	Len	Order	Unit	Inflection
Arabic	ar	47.8k	6.4	VSO	word	high
Indonesian	id	18.6k	6.8	SVO	word	high
Japanese	ja	17.2k	8.5	SOV	none	moderate
Korean	ko	17.2k	8.1	SOV	phrase	moderate
Malay	ms	19.3k	6.8	SVO	word	high
Thai	th	7.4k	7.8	SVO	none	light
Tagalog	tl	28.7k	7.4	VSO	word	high
Vietnamese	vi	9.9k	9.0	SVO	phrase	light
Chinese	zh	13.3k	6.8	SVO	none	light
Taiwanese	zht	39.5k	5.9	SVO	none	light

tion, *Indo-European* languages have, in general, a higher degree of inflection compared to Asian languages. Concerning word segmentation, the corpora were preprocessed using language-specific word-segmentation tools for languages that do not use white-space to separate word/phrase tokens (ja,ko,th,zh,zht). For all other languages, simple tokenization tools were applied. All data sets were case-sensitive with punctuation marks preserved.

The language resources were randomly split into three subsets for the evaluation of translation quality (*eval*, 1000 sentences), the tuning of the SMT model weights (*dev*, 1000 sentences) and the training of the statistical models (*train*). However, in a real-world application, identical language resources covering three or more languages are not necessarily to be expected. In order to avoid a trilingual scenario for the pivot translation experiments, the *train* corpus was randomly split into two subsets of 80k sentences each, whereby the first set of sentence pairs was used to train the SRC-PVT translation models and the second subset of sentence pairs was used to train the PVT-TRG translation models. In total, 924 SMT translation engines were built to cover all 9,240 language pair combinations.

For the training of the SMT models, standard

Table 3: Oracle Pivot Translation Quality (BLEU)
(European Languages) (Asian Languages)

TRG → ↓ SRC	da	de	en	es	fr	hi	it	nl	pl	pt	ptb	ru	ar	id	ja	ko	ms	th	tl	vi	zh	zht
da	–	53.9 (en)	60.3 (nl)	59.1 (en)	57.6 (en)	45.3 (en)	53.4 (en)	57.6 (en)	49.8 (en)	57.8 (ptb)	57.8 (en)	49.5 (en)	48.8 (en)	52.5 (ms)	37.5 (ko)	36.9 (en)	51.9 (id)	51.6 (en)	47.7 (en)	52.6 (en)	34.2 (en)	39.9 (en)
de	57.2 (en)	–	61.3 (nl)	59.3 (en)	57.3 (en)	45.6 (en)	53.6 (en)	58.5 (en)	49.7 (en)	59.2 (ptb)	58.3 (pt)	49.1 (en)	47.8 (en)	52.1 (ms)	37.8 (en)	36.8 (en)	51.5 (en)	51.8 (en)	48.3 (en)	52.2 (en)	33.3 (en)	41.1 (en)
en	59.8 (es)	55.5 (nl)	–	62.7 (pt)	60.7 (es)	45.8 (es)	56.9 (es)	60.1 (es)	49.9 (es)	65.7 (ptb)	65.5 (pt)	50.7 (es)	50.1 (es)	57.2 (ms)	39.4 (ko)	38.0 (ja)	56.8 (id)	51.3 (es)	49.4 (es)	53.6 (es)	33.6 (es)	40.4 (es)
es	59.0 (en)	54.4 (en)	63.3 (pt)	–	59.4 (en)	45.6 (en)	55.7 (en)	58.6 (en)	51.7 (en)	64.7 (ptb)	64.6 (pt)	50.5 (en)	50.1 (en)	55.3 (ms)	38.5 (ko)	37.7 (en)	54.4 (id)	52.5 (en)	49.6 (en)	54.0 (en)	34.1 (en)	40.4 (en)
fr	56.4 (en)	50.9 (en)	58.8 (es)	58.2 (en)	–	43.2 (en)	52.4 (en)	54.8 (en)	47.3 (en)	58.7 (ptb)	57.9 (pt)	47.3 (en)	48.2 (en)	52.5 (ms)	37.8 (en)	37.6 (ja)	51.1 (id)	49.5 (en)	46.6 (en)	50.5 (en)	33.4 (es)	39.9 (en)
hi	50.3 (en)	47.4 (en)	50.5 (ptb)	51.8 (en)	50.8 (en)	–	47.9 (en)	50.2 (en)	44.4 (en)	51.5 (ptb)	51.6 (pt)	44.6 (en)	44.7 (en)	50.3 (ms)	35.7 (ko)	34.6 (en)	50.8 (id)	48.1 (en)	43.6 (en)	48.2 (en)	30.8 (en)	36.9 (en)
it	56.7 (en)	52.8 (en)	60.6 (pt)	59.5 (en)	58.1 (en)	44.8 (en)	–	55.7 (en)	48.8 (en)	60.5 (ptb)	60.2 (pt)	48.1 (en)	47.1 (en)	52.5 (ms)	38.1 (en)	36.8 (en)	52.1 (id)	50.6 (en)	47.3 (en)	51.6 (en)	32.3 (es)	40.5 (en)
nl	60.3 (en)	55.8 (en)	60.9 (es)	61.5 (en)	59.6 (en)	46.3 (en)	55.0 (en)	–	51.1 (en)	60.0 (ptb)	59.5 (en)	50.3 (en)	49.7 (en)	52.6 (ms)	37.7 (en)	36.8 (en)	51.9 (id)	52.0 (en)	49.0 (en)	53.3 (en)	33.3 (en)	39.9 (en)
pl	54.7 (en)	51.1 (en)	56.1 (pt)	56.2 (en)	54.0 (en)	44.2 (en)	51.2 (en)	53.5 (en)	–	56.1 (ptb)	56.6 (pt)	48.7 (en)	46.4 (en)	52.3 (ms)	37.4 (ko)	37.6 (en)	51.6 (id)	50.1 (en)	47.3 (en)	50.5 (en)	32.7 (en)	39.7 (en)
pt	60.6 (ptb)	55.8 (ptb)	68.7 (ptb)	67.0 (ptb)	63.6 (ptb)	47.3 (ptb)	58.8 (ptb)	60.1 (ptb)	51.8 (ptb)	–	67.8 (es)	52.2 (ptb)	52.4 (ptb)	54.8 (ms)	38.1 (ko)	37.3 (en)	53.6 (id)	53.5 (ptb)	50.1 (ptb)	54.8 (ptb)	34.1 (ptb)	42.6 (ptb)
ptb	60.4 (pt)	56.5 (pt)	68.9 (pt)	66.9 (pt)	62.8 (pt)	47.9 (pt)	59.1 (pt)	60.0 (pt)	52.8 (pt)	70.0 (es)	–	51.5 (pt)	52.2 (pt)	54.9 (pt)	38.7 (ko)	37.2 (en)	54.2 (pt)	52.9 (pt)	50.5 (pt)	54.8 (pt)	34.8 (pt)	42.2 (pt)
ru	51.6 (en)	47.6 (en)	53.6 (en)	53.8 (en)	51.5 (en)	42.2 (en)	47.5 (en)	51.2 (en)	46.8 (en)	53.2 (ptb)	53.8 (pt)	–	44.8 (en)	50.3 (ms)	36.7 (en)	35.8 (en)	50.3 (id)	47.4 (en)	44.2 (en)	49.1 (en)	32.0 (en)	37.0 (en)
ar	54.7 (en)	51.3 (en)	56.5 (pt)	57.1 (en)	56.1 (en)	44.9 (en)	51.7 (en)	54.4 (en)	47.2 (en)	55.7 (ptb)	55.6 (pt)	47.9 (en)	–	52.0 (ms)	36.4 (en)	36.1 (en)	52.0 (en)	49.2 (en)	45.6 (en)	51.8 (en)	32.4 (en)	38.7 (en)
id	52.0 (ms)	48.5 (ms)	56.7 (ms)	54.0 (ms)	51.9 (ms)	46.1 (ms)	49.2 (ms)	51.7 (ms)	48.4 (ms)	51.3 (ptb)	51.3 (pt)	46.9 (ms)	47.8 (ms)	–	39.1 (ms)	37.5 (ja)	59.6 (en)	51.7 (ms)	47.8 (ms)	52.7 (ms)	34.6 (ms)	41.5 (ms)
ja	33.5 (en)	31.9 (en)	38.8 (ko)	37.9 (ko)	38.6 (en)	29.3 (ko)	33.0 (en)	34.1 (en)	31.1 (en)	35.8 (ptb)	36.3 (pt)	30.7 (en)	29.8 (ko)	35.5 (ko)	–	46.7 (zh)	33.9 (id)	37.9 (ko)	33.7 (ko)	35.5 (ko)	46.9 (ko)	33.1 (ko)
ko	33.2 (ja)	31.8 (ja)	38.7 (ja)	37.1 (ja)	38.8 (ja)	28.8 (ja)	32.4 (ja)	32.7 (ja)	30.7 (ja)	34.5 (ja)	36.3 (ja)	29.5 (ja)	29.7 (ja)	35.2 (ja)	45.8 (zh)	–	34.2 (id)	38.1 (ja)	32.4 (ja)	33.7 (ja)	47.2 (ja)	32.9 (ja)
ms	53.1 (id)	50.5 (id)	57.8 (id)	55.1 (id)	53.8 (id)	47.3 (id)	49.5 (id)	53.0 (id)	49.3 (id)	52.3 (id)	53.2 (id)	48.7 (id)	48.5 (id)	60.2 (en)	39.8 (id)	37.0 (id)	–	53.2 (id)	48.7 (id)	53.4 (id)	35.1 (id)	42.9 (id)
th	49.5 (en)	45.0 (en)	50.1 (ptb)	49.2 (en)	48.5 (en)	40.6 (en)	45.1 (en)	47.9 (en)	43.2 (en)	50.1 (ptb)	49.8 (pt)	40.8 (en)	41.4 (en)	48.5 (ms)	36.1 (ko)	36.8 (ja)	47.9 (id)	–	41.7 (en)	47.5 (en)	31.4 (id)	37.0 (en)
tl	53.6 (en)	50.2 (en)	54.5 (pt)	55.6 (en)	53.7 (en)	43.7 (en)	49.7 (en)	51.8 (en)	47.0 (en)	53.5 (ptb)	53.1 (pt)	46.0 (en)	45.4 (en)	52.5 (ms)	37.5 (ko)	36.7 (en)	50.8 (id)	50.2 (en)	–	51.3 (en)	33.0 (en)	39.3 (en)
vi	53.2 (en)	48.8 (en)	53.7 (pt)	53.8 (en)	52.6 (en)	42.8 (en)	48.8 (en)	51.8 (en)	46.4 (en)	52.2 (ptb)	53.3 (pt)	45.4 (en)	46.0 (en)	53.0 (ms)	36.8 (ko)	35.4 (ja)	52.8 (id)	49.3 (en)	45.2 (en)	–	31.6 (ms)	38.3 (en)
zh	31.7 (en)	31.2 (nl)	35.4 (zht)	35.8 (en)	34.9 (en)	27.9 (ja)	31.5 (en)	32.1 (en)	29.1 (en)	33.1 (ptb)	33.4 (ja)	27.7 (ms)	27.0 (en)	34.3 (ms)	47.4 (ko)	47.6 (ja)	32.3 (id)	36.3 (en)	30.9 (en)	33.2 (en)	–	33.8 (ja)
zht	44.1 (en)	41.4 (en)	44.5 (pt)	45.4 (en)	44.5 (en)	36.8 (en)	40.8 (en)	43.5 (en)	39.4 (en)	44.6 (ptb)	44.3 (pt)	39.5 (en)	38.2 (en)	44.5 (ms)	40.8 (zh)	38.5 (zh)	44.0 (id)	43.0 (en)	39.4 (en)	42.8 (en)	35.0 (ja)	–

Table 2: Pivot Language Dependency

PVT	BLEU (%)	
	low	high
da	24.2	~ 60.2
de	25.1	~ 61.8
en	26.2	~ 67.7
es	24.9	~ 70.0
fr	24.5	~ 62.3
hi	23.7	~ 53.1
it	23.7	~ 65.8
nl	26.2	~ 61.6
pl	25.2	~ 56.8
pt	25.8	~ 68.9
ptb	24.9	~ 68.8
ru	22.6	~ 56.9

PVT	BLEU (%)	
	low	high
ar	23.3	~ 57.1
id	25.6	~ 57.9
ja	26.8	~ 59.4
ko	25.7	~ 58.3
ms	25.5	~ 57.3
th	23.3	~ 52.3
zht	23.7	~ 44.7
vi	23.6	~ 55.3

word alignment (Och and Ney, 2003) and language modeling (Stolcke, 2002) tools were used. Minimum error rate training (MERT) was used to tune the decoder’s parameters, and was performed on the *dev* set using the technique proposed in (Och and Ney, 2003). For the translation, an in-house multi-stack phrase-based de-

coder was used. For the evaluation of translation quality, we applied the standard automatic evaluation metric BLEU which calculates the geometric mean of n-gram precision by the system output with respect to reference translations multiplied by a brevity penalty to prevent very short candidates from receiving too high a score. Scores range between 0 (worst) and 1 (best) (Papineni et al., 2002). For our experiments, single translation references were used.

Table 2 summarizes the BLEU score ranges of all pivot translation experiments obtained for a given pivot language. The results show a large variation in BLEU scores for all pivot languages indicating that the best pivot choice largely depends on the respective source and target language. For European pivot languages, the best language combination scores are in general much higher than the ones obtained for Asian pivot languages.

Table 3 lists the highest BLEU scores of the pivot translation experiments obtained for all language pair combinations. The pivot language achieving

Table 4: Changes in Pivot Selection for Non-English European and Asian Language Pairs (BLEU)

(Non-English European Language Pairs)											(Asian Language Pairs)												
TRG → ↓ SRC	da	de	es	fr	hi	it	nl	pl	pt	ptb	ru	TRG → ↓ SRC	ar	id	ja	ko	ms	th	tl	vi	zh	zht	
da	-	51.7 (nl)	56.0 (nl)	56.2 (es)	43.1 (es)	50.6 (es)	55.2 (es)	46.7 (pt)	57.8 (ptb)	57.6 (pt)	47.5 (es)	ar	-	52.0 (ms)	35.6 (id)	33.9 (id)	52.0 (id)	46.1 (ms)	41.3 (id)	46.7 (ms)	31.1 (id)	36.4 (id)	
de	55.4 (nl)	-	57.1 (ptb)	55.4 (nl)	43.3 (ptb)	51.9 (es)	54.8 (nl)	47.3 (nl)	59.2 (ptb)	58.3 (pt)	47.8 (nl)	id	47.8 (ms)	-	39.1 (ms)	37.5 (ja)	54.9 (vi)	51.7 (ms)	47.8 (ms)	52.7 (ms)	34.6 (ms)	41.5 (ms)	
es	57.0 (pt)	52.9 (pt)	-	58.4 (pt)	43.9 (pt)	54.7 (ptb)	56.0 (ptb)	48.8 (pt)	64.7 (ptb)	64.6 (pt)	48.4 (ptb)	ja	29.8 (ko)	35.5 (ko)	-	46.7 (zh)	33.9 (id)	37.9 (ko)	33.7 (ko)	35.5 (ko)	46.9 (ko)	33.1 (ko)	
fr	53.2 (pt)	50.1 (nl)	57.2 (pt)	-	41.6 (es)	50.7 (es)	53.4 (es)	45.0 (es)	58.7 (ptb)	57.9 (pt)	45.8 (es)	ko	29.7 (ja)	35.2 (ja)	45.8 (zh)	-	34.2 (id)	38.1 (ja)	32.4 (ja)	33.7 (ja)	47.2 (ja)	32.9 (ja)	
hi	47.3 (ptb)	45.6 (nl)	49.6 (ptb)	48.4 (ptb)	-	45.2 (de)	47.6 (es)	41.4 (es)	51.5 (ptb)	51.6 (pt)	41.7 (es)	ms	48.5 (id)	53.8 (ar)	39.8 (id)	37.0 (id)	-	53.2 (id)	48.7 (id)	53.4 (id)	35.1 (id)	42.9 (id)	
it	53.8 (pt)	50.4 (nl)	58.5 (pt)	56.4 (pt)	42.4 (pt)	-	53.8 (es)	46.8 (pt)	60.5 (ptb)	60.2 (pt)	47.3 (es)	th	39.4 (ms)	48.5 (ms)	36.1 (ko)	36.8 (ja)	47.9 (id)	-	40.5 (id)	44.3 (ms)	31.4 (id)	34.7 (ms)	
nl	55.5 (es)	51.6 (da)	57.7 (ptb)	56.7 (es)	43.9 (es)	52.0 (es)	-	47.7 (pt)	60.0 (ptb)	59.5 (pt)	47.9 (es)	tl	40.8 (id)	52.5 (ms)	37.5 (ko)	36.7 (ja)	50.8 (id)	46.5 (ms)	-	47.0 (ms)	32.3 (id)	36.5 (ms)	
pl	51.8 (pt)	47.9 (pt)	53.9 (pt)	51.8 (pt)	41.9 (pt)	49.6 (es)	51.3 (es)	-	56.1 (ptb)	56.6 (pt)	45.6 (es)	vi	42.5 (ms)	53.0 (ms)	36.8 (ko)	35.4 (ja)	52.8 (ms)	48.6 (ms)	43.6 (ms)	-	31.6 (ms)	37.0 (ms)	
pt	60.6 (ptb)	55.8 (ptb)	67.0 (ptb)	63.6 (ptb)	47.3 (ptb)	58.8 (ptb)	60.1 (ptb)	51.8 (ptb)	-	67.8 (es)	52.2 (ptb)	zh	26.9 (zht)	34.3 (ms)	47.4 (ko)	47.6 (ja)	32.3 (ko)	35.9 (ja)	30.8 (ko)	32.6 (zht)	-	33.8 (ja)	
ptb	60.4 (pt)	56.5 (pt)	66.9 (pt)	62.8 (pt)	47.9 (pt)	59.1 (pt)	60.0 (pt)	52.8 (pt)	70.0 (es)	-	51.5 (pt)	zht	35.9 (id)	44.5 (ms)	40.8 (zh)	38.5 (zh)	44.0 (id)	40.6 (id)	36.8 (id)	40.5 (ms)	35.0 (ja)	-	
ru	50.0 (pt)	46.8 (nl)	52.5 (pt)	50.6 (pt)	40.7 (es)	46.9 (ptb)	49.7 (es)	44.2 (pt)	53.2 (ptb)	53.8 (pt)	-												

the highest scores (*oracle pivot*) for translating the source (*S*) language into the target (*T*) language are given in parantheses. Non-English oracle pivot languages are highlighted in boldface. The figures show that the *English* pivot approach still achieves the highest scores for the majority of the examined language pairs. However, in 49.8% (230 out of 462) of the cases, a non-English pivot language, mainly *Portuguese*, *Brazilian Portuguese*, *Malay*, *Indonesian*, *Japanese*, *Korean*, is preferable. For languages that are closely related like Portuguese vs. Brazilian Portuguese and Malay vs. Indonesian, the related language should be chosen as the pivot language when either translating from or into the respective language for 88.7% (71 out of 80) and 85.0% (68 out of 80) of the pivot translation experiments, respectively. Moreover, *Japanese* is the dominant pivot language when translating from Korean into an other language (95.0%, 19 out of 20), but not for the translation into Korean (30.0%, 6 out of 20). These results suggest that in general pivot languages closely related to the source language have a larger impact on the overall pivot translation quality than pivot languages related to the target language.

Interestingly, for European-only language pairs, only European languages are the oracle pivot language, the majority of which is English. In addition, Spanish is the pivot language of choice when translating from English into another European language and the Dutch pivot achieved the highest BLEU scores for Germanic-only language pairs. On the other hand, when translating between Asian languages, 65.6% (59 out of 90) of the oracle pivot languages are Asian lan-

guages. The Spanish (Chinese) oracle pivot languages for translations between Portuguese and Brazilian Portuguese (Japanese and Korean) also stresses the importance of language relatedness.

In order to investigate the dependency of pivot language selection and language families further, Table 4 summarizes the BLEU scores of pivot translations between only (a) non-English European and (b) Asian language pairs. The results of the European-only language pairs in the table on the left confirm the findings of Table 3. *Portuguese* and *Brazilian Portuguese* are still the dominant pivot languages for non-English European language pairs. An increase of Spanish/Dutch oracle pivot language pairs can be seen for the translation between only Romance/Germanic languages, respectively. Similarly, Malay and Indonesian are the dominant pivot languages, followed by Japanese and Korean, for Asian-only language pairs, most of which achieve BLEU scores that are only slightly lower than the ones for the English oracle pivot language experiments reported in Table 3.

Table 5 summarizes the percentages for the language pairs where the respective pivot language achieved the highest automatic evaluation score for the pivot translation experiments summarized in Table 3 (all language pairs) and Table 4 (non-English European language pairs, Asian language pairs). The results show that English is indeed the pivot language of choice for the majority of the investigated translation directions, but for almost half of the language pairs a non-English pivot language is preferable.

In order to investigate how much improvement

Table 5: Oracle Pivot Language Distribution
(All Language Pairs)

PVT	usage (%)	PVT	usage (%)
en	232 (50.2)	ko	21 (4.5)
pt	40 (8.7)	es	19 (4.1)
ptb	38 (8.2)	nl	5 (1.1)
id	37 (8.0)	zh	4 (0.9)
ms	36 (7.8)	zht	1 (0.2)
ja	29 (6.3)		

(Non-English European Language Pairs)

PVT	usage (%)	PVT	usage (%)
pt	40 (36.3)	nl	10 (9.1)
ptb	32 (29.1)	de	1 (0.9)
es	26 (23.7)	da	1 (0.9)

(Asian Language Pairs)

PVT	usage (%)	PVT	usage (%)
id	28 (31.1)	zh	4 (4.4)
ms	27 (30.0)	zht	2 (2.2)
ja	15 (16.6)	vi	1 (1.1)
ko	12 (13.3)	ar	1 (1.1)

Table 6: Gain of non-English Pivot Languages

PVT	(oracle)	Gain in BLEU (%)		
		avg	min	max
zh	(4)	4.7	3.2	6.1
ja	(27)	2.5	0.1	13.3
id	(35)	2.4	0.6	5.4
pt	(31)	2.3	0.3	4.6
ptb	(32)	2.1	0.3	4.9
ko	(19)	1.9	0.1	11.4
ms	(34)	1.8	0.1	3.9
es	(4)	0.8	0.1	2.4
nl	(2)	0.6	0.5	0.8

in pivot translation performance can be achieved by using non-English pivot languages instead of an English pivot, we calculated the difference in BLEU scores for all 188 non-English language pairs where the non-English pivot language improved translation quality. Table 6 summarizes the average, minimal and maximal gains in BLEU scores for the respective pivot language translation experiments. The pivot languages are sorted according to the highest average increase in translation performance and the amount of improved language pairs are given in parantheses. In total, an average gain of 2.2 BLEU points were obtained for the investigated language pairs. The highest gains (13.4/11.4 BLEU points) were achieved for the Japanese/Korean pivots when translating Korean/Japanese into Chinese, respectively.

3 Indicators of Pivot Translation Quality

The diversity of the pivot language selection reported in the last section rises the question of what makes a language a good pivot language for a given language pair.

We investigated the following eight factors (comprised of a total of 45 distinct features) based on the language resources and SMT engines (SRC-PVT, PVT-TRG) used for the pivot translation experiments described in Section 2 where the total number of features of each factor is given in brackets. For SMT-engine-related features, both translation directions (SRC-PVT, PVT-TRG) are taken into account.

- *language family* [2]: a binary feature verifying whether the source and target languages of the SMT engines belong to the same family or not.
- *vocabulary* [15]: the training data vocabulary size of source and target languages, the ratio of source and target vocabulary sizes, and the overlap between source and target vocabulary.
- *sentence length* [12]: the average sentence length of source and target training sets and the ratio of source and target sentence length.
- *reordering* [6]: the amount and span of word order differences (reordering) in the training data and the *Reordering Quantity* score as proposed in (Birch et al., 2008).
- *language perplexity* [4]: perplexity of the utilized language models measured on the *dev/eval* data sets.
- *translation model entropy* [2]: amount of uncertainty involved in choosing candidate translation phrases as proposed in (Koehn et al., 2009).
- *engine performance* [2]: the BLEU scores of the respective SMT engine used for the pivot translation experiments.
- *monotonicity* [2]: the BLEU score difference of a given SMT engine for decoding with and without a reordering model.

The impact of the above factors in isolation on the translation performance is measured using linear regression which models the relationship between a response variable and one or more explanatory variables. Data sets are modeled using linear functions and unknown model parameters are estimated from the data. In this paper, the response variable is defined by the BLEU metric (measuring the pivot translation performance) and the explanatory variables are given by the feature values obtained for each of the respective language pair combinations. Figure 1 gives an example for a simple linear regression using the *reordering quantity* feature as the explanatory variable for (a) all language pairs, (b) European languages only, and (c) Asian languages only. The “goodness of

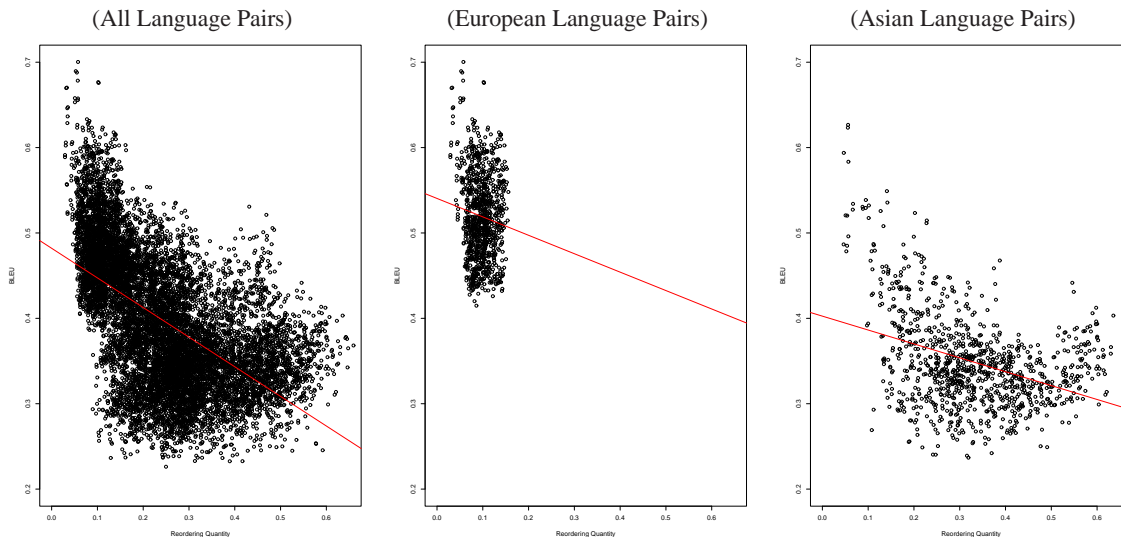


Figure 1: Linear Regression Example (*Reordering Quantity*)

fit” of the explanatory variable(s) is calculated using the R^2 coefficient of determination, which is a statistical measure of how well the regression line approximates the real data points. An R^2 of 1.0 indicates that the regression line perfectly fits the data. For the *translation model entropy* factor, for example, we obtain an R^2 of 0.4604 for all language pairs, which indicates that 46.04% of the differences in translation performance can be explained by this factor.

3.1 Predictive Power of Single Factors

Table 7 summarizes the R^2 scores of the multiple linear regression analysis of the respective investigated factors, i.e. all features of a given factor are combined and treated as multiple explanatory variables. In total, 81% of the system performance variations can be explained when all investigated factors are taken into account. For European language pairs, the impact is even larger (91%). However, for Asian language pairs, the investigated factors have much less correlation (R^2 of 0.5888) with the overall pivot translation translation quality, indicating the difficulty of selecting an appropriate pivot language for translation tasks including Asian languages.

The impact of each factor on the translation performance is also given in Table 7. The results show that *engine performance* is the most correlated factor, followed by *translation model entropy* and *reordering* when all language combinations are taken into account. *Language family* and *language perplexity* seems to have the least impact on translation performance. However, when applying linear regression on language subsets (only Euro-

Table 7: Impact on Translation Performance

Explanatory Variable	R^2		
	All	European	Asian
all factors	0.8102	0.9106	0.5880
engine performance	0.7438	0.7906	0.5151
translation model entropy	0.4604	0.3669	0.1661
reordering	0.4383	0.4593	0.1806
vocabulary	0.3112	0.3867	0.2389
monotonicity	0.2682	0.0149	0.1323
sentence length	0.1717	0.6052	0.0724
language family	0.1204	0.1280	0.0982
language perplexity	0.0826	0.1100	0.0337

pean vs. only Asian languages), the impact of factors largely differs. Similar to all language pairs, the *engine performance* factor is most relevant for both European and Asian language subsets.

For pivot translations between European languages, *sentence length*, *reordering* and *vocabulary* are more predictive than the *translation model entropy* factor. Moreover, the *monotonicity* factor obtains the lowest R^2 score indicating that word order differences between European languages occur mainly on the phrase-level (*local reordering*) and that only minor gains can be achieved when reordering successive phrases. The high R^2 score for *sentence length* also suggests that the ratio of sentence length is an important feature when selecting an appropriate pivot language for closely related languages.

On the other hand, looking at the Asian language pair regression results, the lower R^2 scores underline the large diversity between the Asian languages. Relatively high R^2 scores for *reordering* and *monotonicity* are obtained for Asian lan-

Table 8: Factor Contribution

Explanatory Variable	R^2		
	All	European	Asian
all factors	0.8102	0.9106	0.5880
w/o engine performance	0.5621	0.8755	0.3683
w/o language perplexity	0.7734	0.8895	0.5488
w/o sentence length	0.7856	0.8989	0.5501
w/o reordering	0.7958	0.8999	0.5712
w/o vocabulary	0.7961	0.8766	0.5669
w/o translation model entropy	0.8004	0.9024	0.5748
w/o monotonicity	0.8026	0.9024	0.5768
w/o language family	0.8035	0.9022	0.5793

guages, indicating that structural differences between the pivot language and the source/target language largely affects the overall pivot translation quality.

3.2 Contribution of Single Factors

Besides the predictive power of each factor, we calculated the R^2 scores of all the factors besides one (*leave-one-out*) in order to investigate the contribution of each factor to the multiple linear regression analysis. In general, the smaller the R^2 score after omitting a given factor, the larger the contribution of this factor on the explanation of the overall translation performance is supposed to be.

The results summarized in Table 8 show that the largest contribution for all language pairs is obtained for the *engine performance* factor, followed by *language perplexity* and *sentence length*. Interestingly, the *vocabulary* factor contributes as much as the *engine performance* factor for European languages, but not for Asian languages. This confirms that morphological similarities between highly inflected languages are important to identify an appropriate pivot language. Moreover, for European-only and Asian-only language pairs, the omission of any of these factors led to lower R^2 scores, but the difference towards the complete factor set is much smaller. This shows the importance of all the investigated features for the task of pivot language selection, especially if languages of large diversity are to be taken into account.

3.3 Translation Direction Dependency

In order to investigate whether the selection of a pivot language depends more on its relationship towards the source language or the target language, we carried out a linear regression analysis based on all factors using (a) only source-language-related features (*SRC-PVT only*) and (b)

Table 9: Source vs. Target Language Dependency

Explanatory Variable	R^2		
	All	European	Asian
all factors	0.8102	0.9106	0.5880
SRC-PVT only	0.4923	0.3125	0.2805
PVT-TRG only	0.4732	0.6505	0.2986

only target-language-related features (*PVT-TRG only*). The results are summarized in Table 9.

In order to distinguish between languages of large diversity, the source language features seem to be more predictive than the target language features. However, for more coherent language pairs, like in the case of European languages, the impact on how much language diversity affects pivot translation performance shifts towards target-language-related features. However, the restriction to either the source or the target features leads to a large decrease in the R^2 scores for all language data sets, underlining the importance of both source-language-related and target-language-related feature sets to identify an appropriate pivot language for a given language pair.

4 Conclusion

We investigated the impact of eight translation quality indicators for the task of pivot translation between 22 languages covering a large diversity of language characteristics. A linear regression analysis showed that 81% of the variation in translation performance differences can be explained by the combination of these factors. The most informative factor in identifying the best pivot language is *engine performance*, i.e., the translation quality of the SMT engines used to translate (a) the source input into the intermediate language and (b) the intermediate language MT output into the target language. In addition, the highest correlation of the investigated factors towards pivot translation performance was obtained when both source-language-related and target-language-related features were combined. The importance of source vs. target language features largely depends on the diversity of the investigated language pairs, i.e., source language features are preferable for heterogeneous language pairs whereas the focus shifts towards target-language-related features for more coherent language pairs. In addition, the differentiation between European and Asian languages revealed that the task of identifying a pivot lan-

guage for new language pairs largely depends on the availability of structurally similar languages.

As future work, we are planning to investigate the importance of the factors analyzed in Section 3 in the selection of pivot languages for new language pairs by applying a machine learning algorithm like *Support Vector Machines* (SVM) to train discriminative models for the task of predicting a pivot language that achieves the highest translation performance for a given translation task.

References

- Nicola Bertoldi, Madalina Barbaiani, Marcello Federico, and Roldano Cattoni. 2008. Phrase-Based Statistical Machine Translation with Pivot Languages. In *Proceedings of the 5th International Workshop on Spoken Language Translation (IWSLT)*, pages 143–149, Hawaii, USA.
- Alexandra Birch, Miles Osborne, and Philipp Koehn. 2008. Predicting success in machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 745–754, Honolulu, Hawaii.
- Adria de Gispert and Jose B. Marino. 2006. Catalan-english statistical machine translation without parallel corpus: bridging through spanish. In *Proceedings of 5th International Conference on Language Resources and Evaluation (LREC)*, pages 65–68, Genoa, Italy.
- Genichiro Kikui, Seiichi Yamamoto, Toshiyuki Takezawa, and Eiichiro Sumita. 2006. Comparative study on corpora for speech translation. *IEEE Transactions on Audio, Speech and Language*, 14(5):1674–1682.
- Philipp Koehn, Alexandra Birch, and Ralf Steinberger. 2009. 462 Machine Translation Systems for Europe. In *Proceedings of the Machine Translation Summit XII*, Ottawa, Canada.
- Gregor Leusch, Aurélien Max, Josep Maria Crego, and Hermann Ney. 2010. Multi-Pivot Translation by System Combination. In *Proceedings of 7th International Workshop on Spoken Language Translation (IWSLT)*, pages 299–306, Paris, France.
- Franz Josef Och and Hermann Ney. 2003. A Systematic Comparison of Various Statistical Alignment Models. *Computational Linguistics*, 29(1):19–51.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 311–318, Philadelphia, USA.
- Michael Paul, Hirofumi Yamamoto, Eiichiro Sumita, and Satoshi Nakamura. 2009. On the Importance of Pivot Language Selection for Statistical Machine Translation. In *Proceedings of the North American Chapter of the Association for Computational Linguistics - Human Language Technologies (NAACL/HLT)*, pages 221–224, Boulder, USA.
- Andreas Stolcke. 2002. SRILM - an extensible language modeling toolkit. In *Proceedings of the 7th International Conference on Spoken Language Processing (ICSLP)*, pages 901–904, Denver.
- Masao Utiyama and Hitoshi Isahara. 2007. A comparison of pivot methods for phrase-based statistical machine translation. In *Proceedings of Human Language Technologies (HLT)*, pages 484–491, New York, USA.
- Hua Wu and Haifeng Wang. 2007. Pivot Language Approach for Phrase-Based Statistical Machine Translation. In *Proceedings of 45th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 856–863, Prague, Czech Republic.