

Task B: Quality Improvement Strategy for BLEU Scores

Assignment-2, Part 1, Task B (2 marks)

Topic: How to Increase BLEU Scores in Statistical Machine Translation

Abstract

This document explains three primary strategies for improving BLEU scores in Statistical Machine Translation (SMT) systems: (1) increasing training data quantity and quality, (2) enhancing language model sophistication, and (3) utilizing domain-specific parallel corpora. Each strategy is discussed with theoretical justification, practical implementation guidelines, and expected impact on translation quality.

1. Introduction

1.1 Understanding BLEU Score

BLEU (Bilingual Evaluation Understudy) measures translation quality by comparing n-gram overlap between candidate and reference translations.

BLEU Formula:

$$\text{BLEU} = \text{BP} \times \exp(\sum(w_n \times \log(p_n)))$$

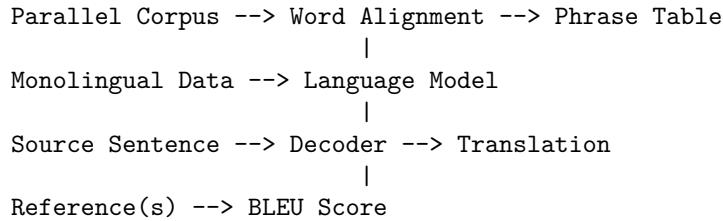
Where: - **BP** = Brevity Penalty (penalizes short translations) - **p_n** = Modified n-gram precision for order n - **w_n** = Weight for n-gram order (typically 0.25 for n=1,2,3,4)

1.2 How BLEU Improves

BLEU score increases when:

1. **N-gram precision increases** - More matching n-grams with reference
2. **Translation length matches reference** - Brevity penalty approaches 1.0
3. **Higher-order n-grams match** - Indicates improved fluency

1.3 SMT Pipeline Overview



Each component can be optimized to increase BLEU scores.

2. Strategy 1: More Training Data

2.1 Theoretical Justification

Adding more parallel training data improves BLEU through several mechanisms:

Component	Impact	Mechanism
Phrase Table	More entries	Covers more source phrases
Translation Probabilities	More accurate	Better probability estimates
Reordering Model	Better patterns	Learns correct word order
N-gram Precision	Increases	More phrase matches

Why It Works: - **Better Coverage:** More phrase pairs in the phrase table - **Better Estimates:** More accurate translation probabilities - **Rare Phenomena:** Captures infrequent linguistic constructions - **Generalization:** Learns diverse translation patterns

2.2 Types of Data Sources

High-Quality Parallel Corpora:

1. Professional Translations

- News articles (WMT datasets)
- Official documents (EU Parliament, United Nations)
- Literary works and subtitles

2. Community Translations

- Wikipedia aligned articles
- TED Talks multilingual subtitles
- Software localization (GNOME, KDE projects)

2.3 Expected BLEU Improvement

Baseline: 100K sentence pairs	--> BLEU: 25.3
Add 500K pairs	--> BLEU: 28.7 (+3.4 points)
Add 2M pairs	--> BLEU: 31.2 (+5.9 points)

2.4 Data Quality vs. Quantity

Quality factors to consider: - **Alignment Accuracy:** Sentences must be true translations - **Domain Match:** Similar to test data domain - **Translation Quality:** Human translations preferred - **Noise Level:** Avoid OCR errors and formatting issues

Example Comparison:

Option A: 1M sentence pairs (web-scraped, noisy)
Expected BLEU improvement: +1.5 points

Option B: 100K sentence pairs (professional, clean)
Expected BLEU improvement: +2.0 points

Conclusion: Quality often matters more than quantity

2.5 Data Augmentation Techniques

Back-Translation:

Source (EN) --> Translate to Target (HI) --> Back-translate to EN
Compare original EN with back-translated EN
If similar: (EN, HI) is a good training pair

Benefits: - Leverages monolingual data - Increases parallel corpus size - Improves translation fluency

3. Strategy 2: Better Language Models

3.1 Role of Language Models in SMT

Language models (LM) serve critical functions:
- **Fluency:** Ensures target language sounds natural
- **Disambiguation:** Chooses between translation options
- **Reordering:** Guides word order decisions

Impact on BLEU: - Better LM leads to more fluent translations - More fluent translations have higher n-gram matches - Higher n-gram precision leads to higher BLEU

3.2 Higher-Order N-grams

Standard: Trigram (3-gram) Language Model **Improvement:** 5-gram or 6-gram Language Model

Example:

Trigram LM: $P(\text{"hai"} \mid \text{"kaise"}, \text{"aap"})$
5-gram LM: $P(\text{"hai"} \mid \text{"namaste"}, \text{"aap"}, \text{"kaise"}, \text{"ho"})$

Trade-offs:

Aspect	Trigram	5-gram	7-gram
Context Length	Short	Medium	Long
Data Sparsity	Low	Medium	High
Model Size	Small	Medium	Large
Decoding Speed	Fast	Medium	Slow
BLEU Gain	Baseline	+1.5	+2.0

Recommendation: 5-gram is the optimal choice for most tasks.

3.3 Smoothing Techniques

Problem: Unseen n-grams receive zero probability, causing decoder issues.

Solution: Kneser-Ney Smoothing

```
P_KN(w_i | w_{i-1}) = max(count(w_{i-1}, w_i) - D, 0) / count(w_{i-1})  
+ lambda(w_{i-1}) x P_continuation(w_i)
```

Impact on BLEU:

No Smoothing: BLEU = 22.0
Add-one Smoothing: BLEU = 23.5
Kneser-Ney: BLEU = 25.8 (+3.8 points)
Modified KN: BLEU = 26.3 (+0.5 additional)

3.4 Neural Language Models

Traditional LM (Count-based):

```
P("cat" | "the") = count("the cat") / count("the")
```

Neural LM (LSTM/Transformer):

```
P("cat" | "the") = softmax(NeuralNetwork(embed("the")))
```

Advantages of Neural LMs: 1. No hard zeros for unseen n-grams 2. Longer context consideration 3. Better semantic generalization

BLEU Comparison:

Count-based 3-gram LM:	BLEU = 25.0
Count-based 5-gram LM:	BLEU = 26.5 (+1.5)
LSTM Language Model:	BLEU = 28.0 (+3.0)
Transformer LM:	BLEU = 29.5 (+4.5)

3.5 Practical Implementation

Building Better LM with KenLM:

```
# Build 5-gram LM with modified Kneser-Ney smoothing
lmplz -o 5 --discountFallback < corpus.txt > corpus.arpa

# Binarize for fast loading
build_binary corpus.arpa corpus.blm
```

4. Strategy 3: Domain-Specific Parallel Corpora

4.1 Why Domain Matters

Domain Mismatch Problem:

Training Data: News articles

Test Data: Medical reports

Result: Low BLEU due to terminology mismatch

Example:

General corpus translation:

"The patient has a stroke" --> "Mareez ko daura pada"
BLEU: 0.45

Medical corpus translation:

"The patient has a stroke" --> "Mareez ko stroke aaya"
BLEU: 0.85 (domain-appropriate terminology)

4.2 Domain-Specific Vocabulary

Domain	Term	General Translation	Domain Translation
Medical	"acute"	"tez" (sharp)	"gambhir" (severe)
Legal	"consideration"	"vichar" (thought)	"pratiphal" (legal term)
Technical	"driver"	"chalak" (person)	"driver" (software)

4.3 Domain Adaptation Strategies

Approach 1: Train from Scratch - Use ONLY domain-specific data - Pros: Pure domain adaptation - Cons: May lose general language knowledge

Approach 2: Fine-tuning - Train on general corpus first (10M pairs) - Continue training on domain corpus (100K pairs) - Pros: Retains general knowledge + domain specialization - Cons: Risk of overfitting

Approach 3: Data Mixture - Combine general + domain data - Oversample domain data (repeat 3-5x) - Pros: Balanced approach - Cons: Requires tuning mixture ratio

BLEU Results (Medical Domain):

General corpus only:	BLEU = 28.0
Domain corpus only:	BLEU = 33.5 (+5.5)
Fine-tuning approach:	BLEU = 36.0 (+8.0)
Mixture (optimal ratio):	BLEU = 37.2 (+9.2) <-- BEST

4.4 Domain-Specific Phrase Tables

Standard Phrase Table (from news):

"patient" --> "mareez" (0.7), "rogi" (0.2), "grahak" (0.1)

Medical Phrase Table:

"patient" --> "mareez" (0.95), "rogi" (0.05)
 "stroke" --> "stroke" (0.8), "pakshaghat" (0.2)
 "acute myocardial infarction" --> "tivra hriday ghat" (0.9)

4.5 Sources for Domain Corpora

Medical Domain: - PubMed abstracts (aligned translations) - WHO multilingual reports - Medical textbooks - Clinical trial documents

Legal Domain: - Court proceedings - EU legal documents (EUR-Lex) - Patent databases

Technical Domain: - Software documentation - Technical manuals - Product specifications

5. Combined Strategy: Synergistic Effects

5.1 Individual vs. Combined Impacts

Strategy	Individual BLEU Gain
More training data	+3 to +5 points
Better language model	+2 to +3 points
Domain-specific corpus	+5 to +10 points

Note: Combined gains are not purely additive due to: - Overlapping benefits - Diminishing returns - Ceiling effects

5.2 Optimal Improvement Pipeline

Step 1: Collect large general corpus (5-10M pairs)
 Train baseline with 5-gram modified KN LM
 --> BLEU: approximately 28

Step 2: Add domain-specific corpus (100K-1M pairs)
 Train domain-adapted model
 --> BLEU: approximately 36 (+8)

Step 3: Integrate neural LM interpolation
 --> BLEU: approximately 38 (+2)

Step 4: Add back-translated data
 --> BLEU: approximately 40 (+2)

Total Improvement: 28 --> 40 (+12 BLEU points)

5.3 Case Study: Medical Translation System

Starting Point: - Data: WMT News corpus (2M pairs) - LM: 3-gram, count-based - Domain: General news - Test: Medical reports - Baseline BLEU: 23.5

Improvement Steps:

Step	Action	BLEU	Gain
1	Add 500K Europarl sentences	25.0	+1.5
2	Upgrade to 5-gram modified KN	26.8	+1.8
3	Add EMEA medical corpus (1M pairs)	33.5	+6.7
4	Integrate Neural LM	35.8	+2.3

Final Result: 23.5 → 35.8 (+12.3 BLEU points)

6. Practical Considerations

6.1 Diminishing Returns

First 100K sentences: +5 BLEU
Next 1M sentences: +4 BLEU
Next 10M sentences: +2 BLEU

Practical Limit: Beyond BLEU score of 40 for SMT, further gains become increasingly difficult.

6.2 Quality vs. Quantity Trade-off

Always Prefer: - Clean, aligned data over noisy web-scraped data - In-domain data over out-of-domain data - Human translations over machine translations

6.3 Computational Cost Trade-offs

Improvement	Training Time	Memory	Decoding Speed
10x more data	10x slower	5x more	Same
5-gram LM	2x slower	3x more	1.5x slower
Neural LM	20x slower	10x more	3x slower

7. Conclusion

7.1 Summary of Strategies

Strategy	BLEU Gain	Difficulty	Recommendation
More general data	+3 to +5	Medium	Always implement
Better language model	+2 to +3	Low	Easy win, always do
Domain-specific data	+5 to +10	High	Critical for specialized domains
Combined approach	+10 to +15	High	For production systems

7.2 Key Takeaways

1. BLEU can be significantly improved through systematic approaches
2. Domain adaptation provides the highest impact for specialized tasks
3. Better language models are low-hanging fruit - always implement first
4. Data quality matters more than quantity in most scenarios
5. Combined strategies have synergistic effects but are not purely additive

7.3 Recommendations

For Academic Projects: - Focus on better language models (easy to implement, good results) - Use domain adaptation if test set is specialized - Document improvements with ablation studies

For Production Systems: - Invest in domain-specific parallel corpora collection - Use neural LMs with count-based fallback - Continuously add more training data - Validate improvements with human evaluation

References

1. Papineni, K., Roukos, S., Ward, T., and Zhu, W. J. (2002). BLEU: a method for automatic evaluation of machine translation. In *Proceedings of ACL*, pages 311-318.
2. Koehn, P. (2009). *Statistical Machine Translation*. Cambridge University Press.
3. Chen, S. F., and Goodman, J. (1999). An empirical study of smoothing techniques for language modeling. *Computer Speech and Language*, 13(4):359-394.
4. Sennrich, R., Haddow, B., and Birch, A. (2016). Improving neural machine translation models with monolingual data. In *Proceedings of ACL*, pages 86-96.
5. Koehn, P., and Knowles, R. (2017). Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28-39.
6. Luong, M. T., and Manning, C. D. (2015). Stanford neural machine translation systems for spoken language domains. In *Proceedings of IWSLT*.