

言語横断LDAモデルを用いた
統計的機械翻訳システム
西尾 拓 福富 崇博 貞光 九月 山本 幹雄

従来の統計的機械翻訳システム

提案手法

翻訳モデル

+

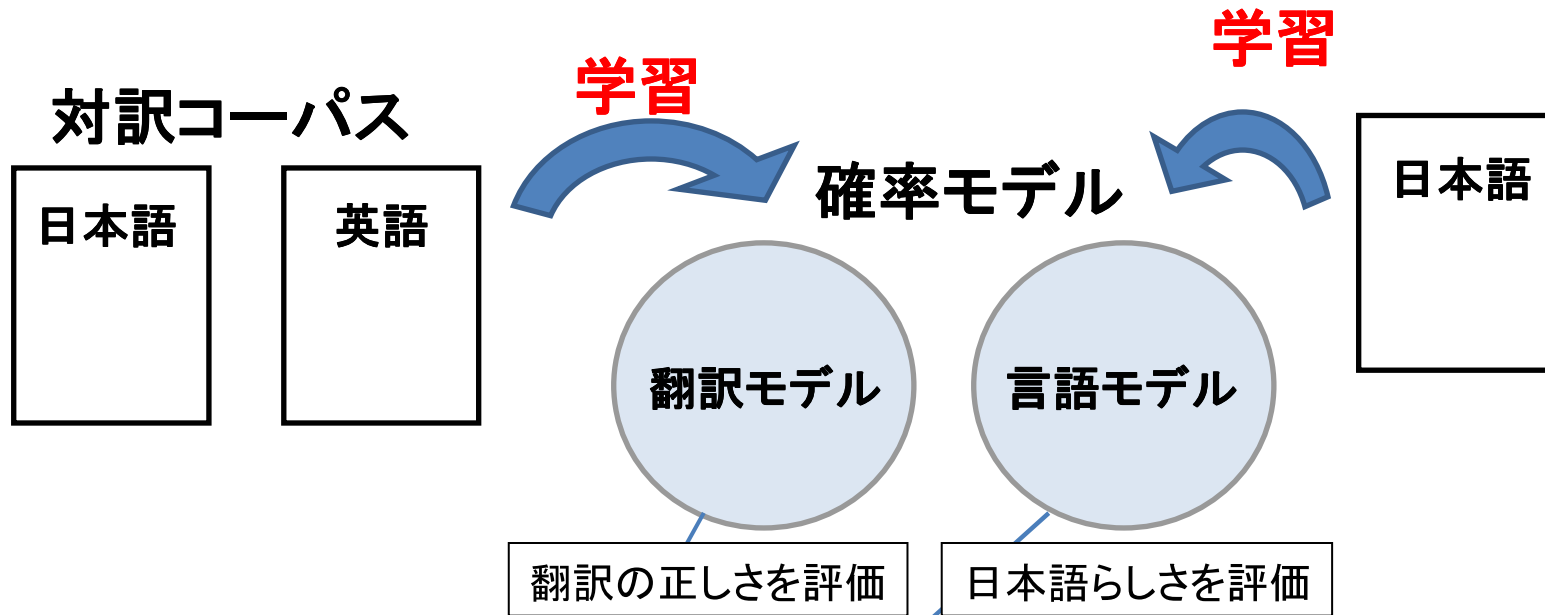
言語モデル

+

言語横断LDAモデル

言語横断LDAモデルの捉えるトピック情報を
統計的機械翻訳システムに持ち込むことで、
翻訳システム全体の翻訳精度向上を実現する。

統計的機械翻訳システム



翻訳したい文

e

入力

デコーダ
翻訳候補を探索

出力

翻訳結果

\hat{j}

最も高い評価

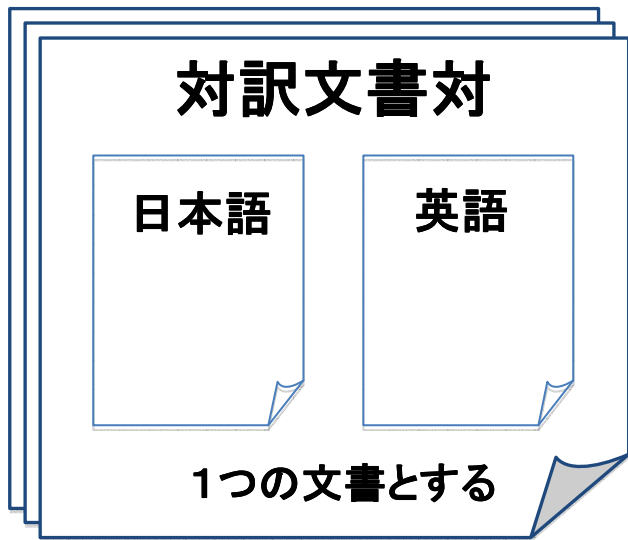
I watch every match
at the Olympics.

オリンピックの全ての縁談を観察。
オリンピックの全てのマッチ箱を警戒する。
オリンピックの全ての試合を観戦する。0.2
.....

0.17
0.18
0.2
.....

オリンピックの全ての
試合を観戦する。

言語横断LDAモデル



共通のトピックにおける
両言語の単語出現確率を学習

LDA学習



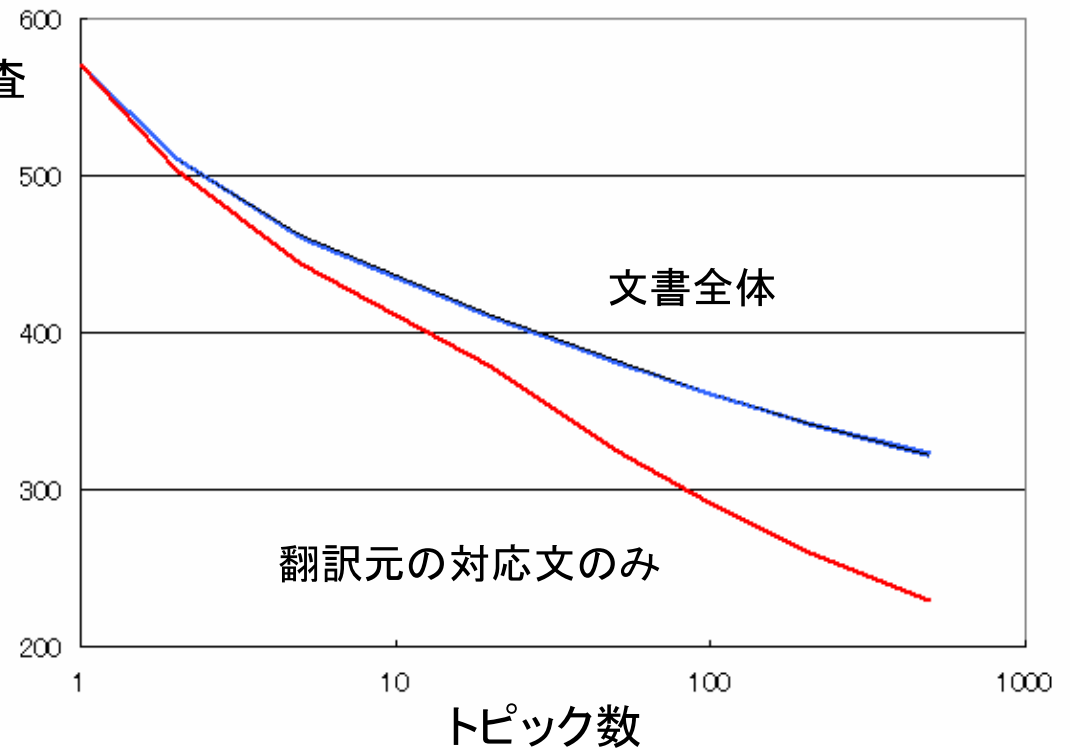
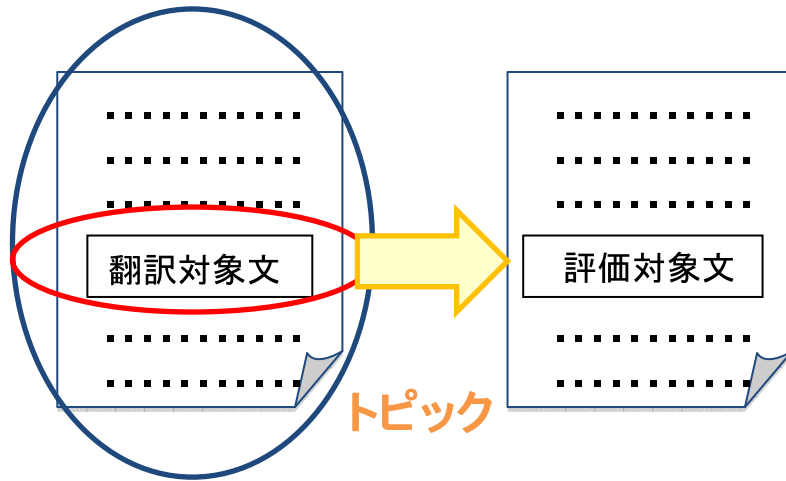
トピック	出現確率の高い単語
WEB	キーワード, 文章, 手続き, サーバ, クライアント menu, job, message, password, script, icon, client
画像	シェーディング, dpi, ポリゴン, 画素, 輝度, RGB gray, gamma, picture, luminance, dct, rgb
機械工学	内圧, 弁, 噴射, 点火, 燃焼, NO _x , アクセル, 燃料 manifold, spark, ignition, purge, exhaust, nox, fuel

実験条件

統合手法	リランキング法	log-linear model
使用コーパス	特許文180万文 (NTCIR7)	
	英語語彙サイズ	121,815
	日本語語彙サイズ	139,491
トピック数	500	2,5,10,20,50,100,200
デコーダ	Moses	Pharaohコピー
Reordering table	あり	なし
言語モデル	SRI Language Modeling Toolkit により学習した 5gramモデル	
翻訳モデル	フレーズベースモデル	
N-best	100	—

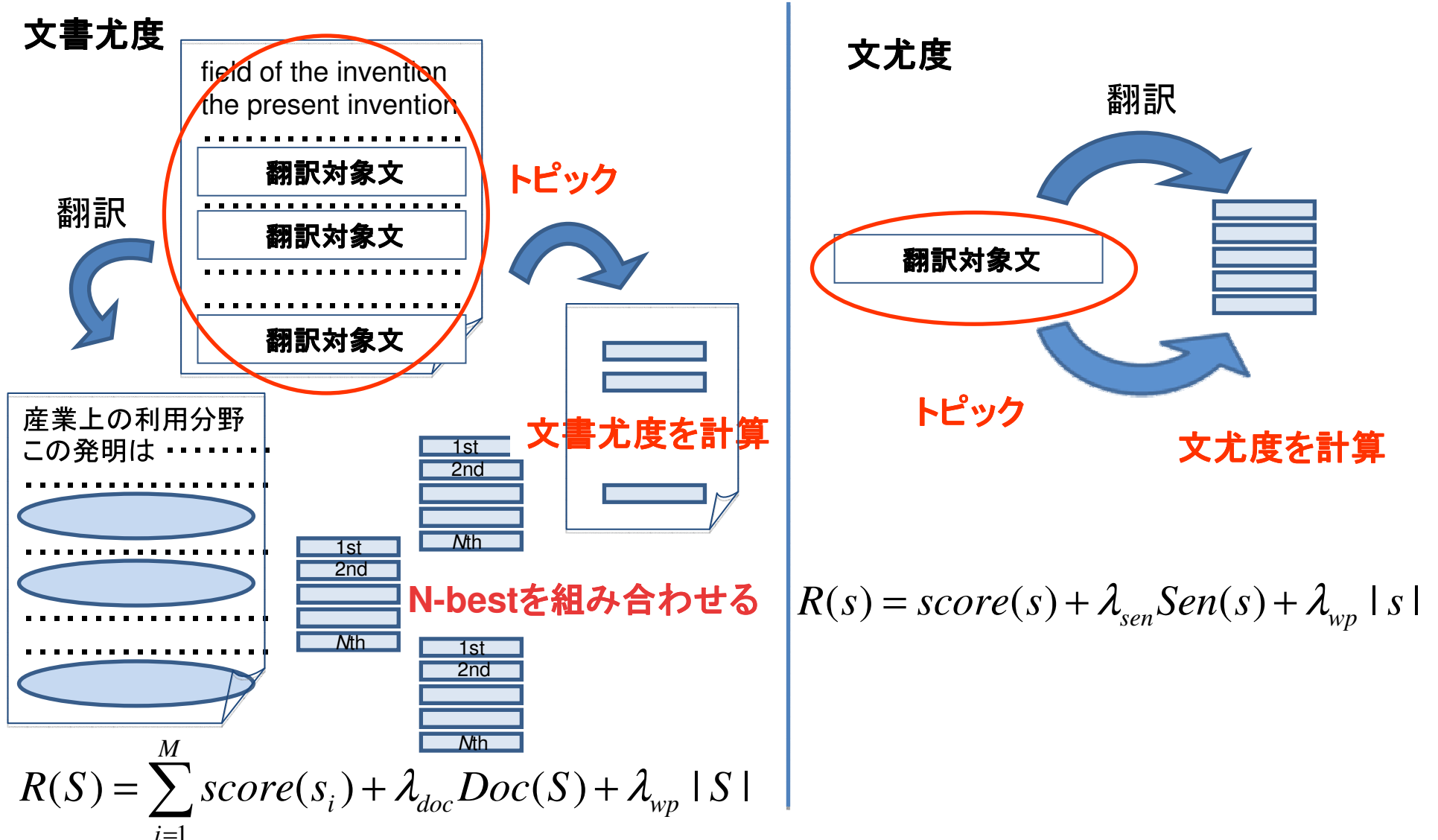
言語横断LDAモデル 適応範囲の違いによる比較

各適応範囲のパープレキシティの推移を調査



**適応の範囲を翻訳元の対応文のみに絞った場合に
言語モデルとしての性能が最も良くなる**

リランキング法 文書尤度/文尤度



log-linear modelによる統合

従来の翻訳システムにおける翻訳候補の評価式を変更する

翻訳らしさの評価式

従来手法で用いられる特徴関数

$$score(s) = \lambda_{LM} L(s) + \lambda_{TM} T(s) + \lambda_D D(s)$$

言語横断LDAモデルを追加

$$- \lambda_{WP} |s| - \lambda_{UNK} \cdot unk(s) + \lambda_{LDA} LDA(s)$$

原言語文

I watch every match
at the Olympics .

トピックを捉える

翻訳候補

オリンピックの全ての縁談を観察する。
オリンピックの全てのマッチ箱を警戒する。
オリンピックの全ての試合を観戦する。
.....

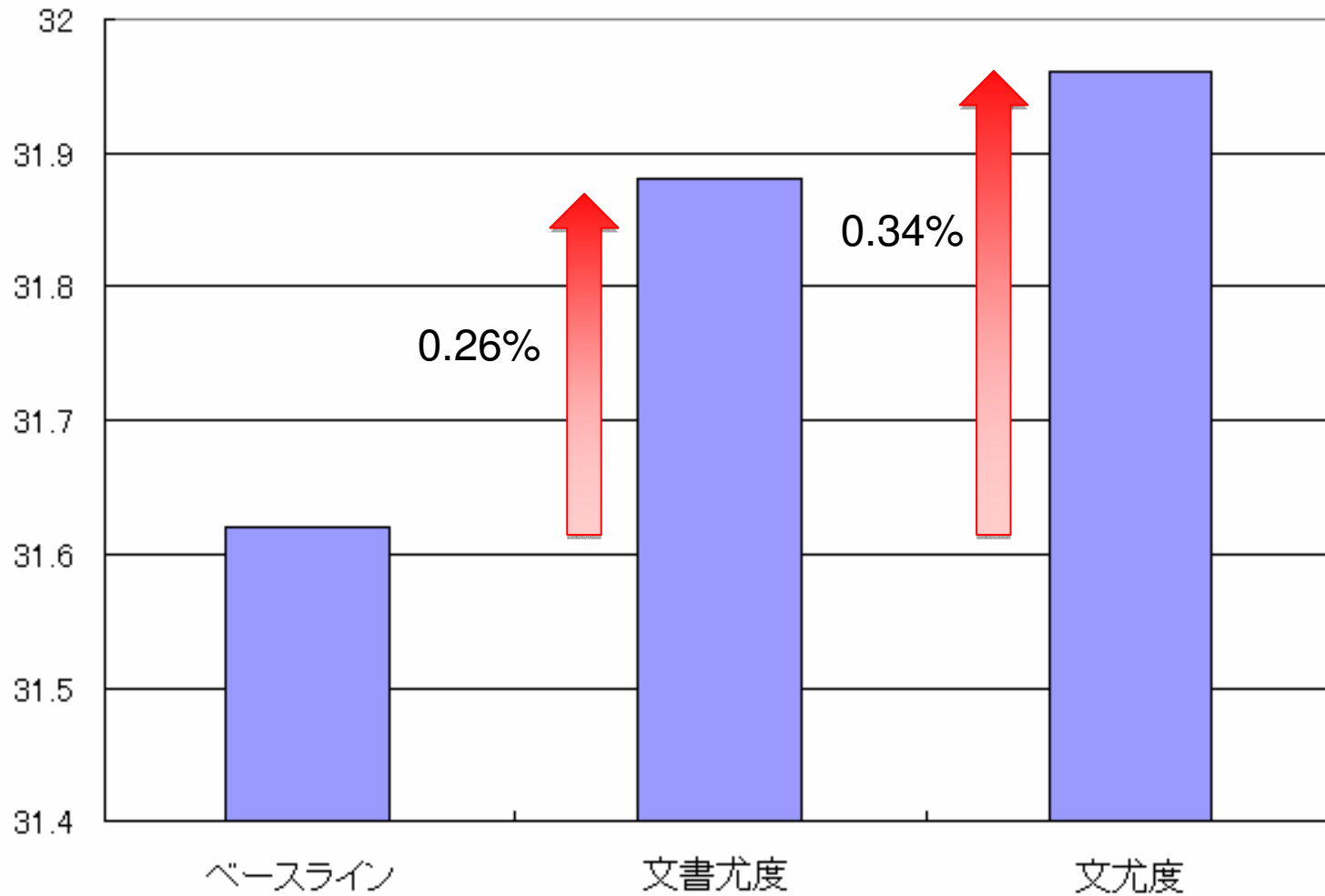
Score(s)

0.02
0.001
0.4
.....

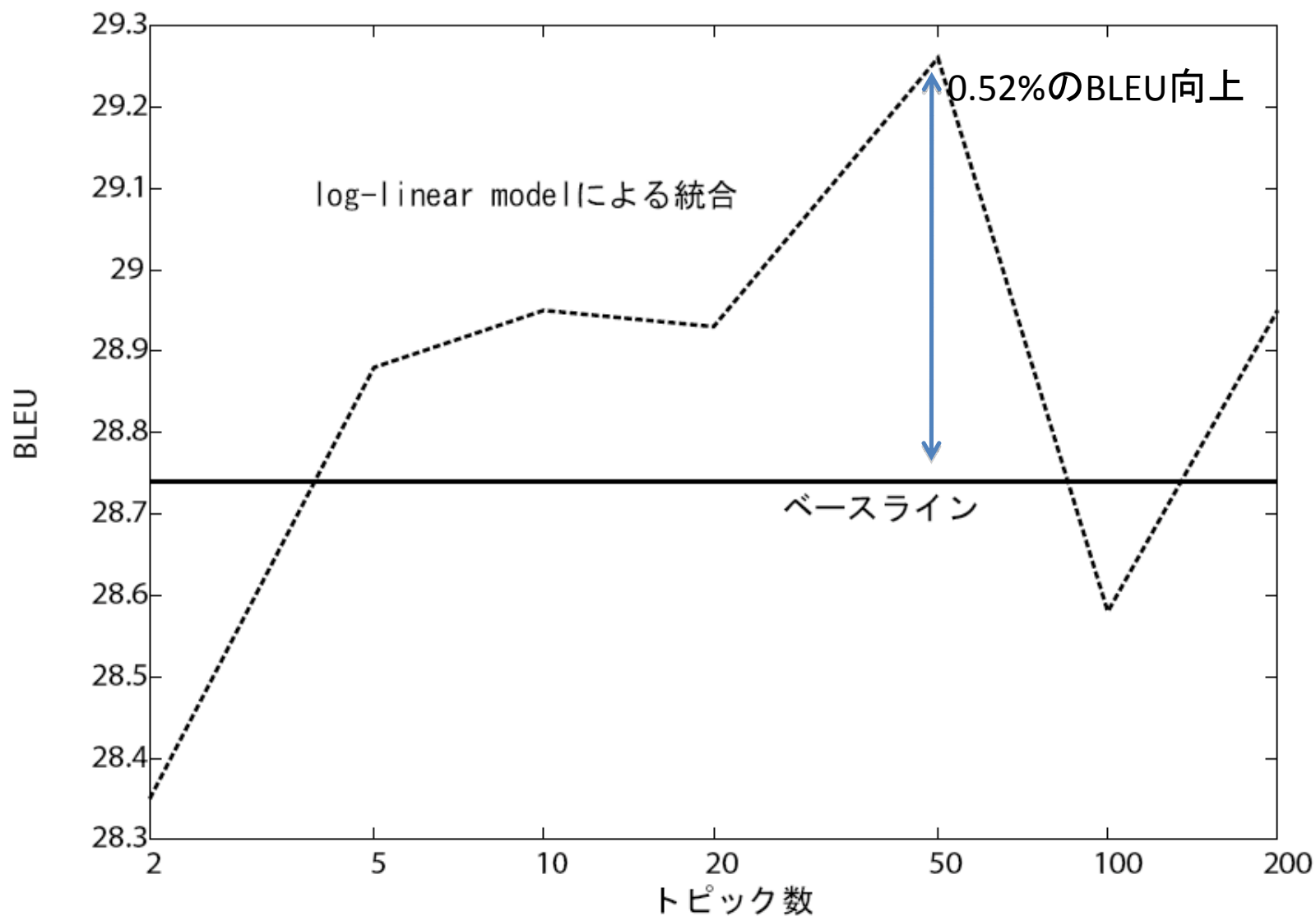
LDAモデルによる
確率評価を含む

翻訳結果として出力

翻訳実験(リランキング法) 実験結果



翻訳実験(log-linear modelによる統合) 実験結果



翻訳結果(log-linear modelによる統合) 改善した例

- 英文
 - On the other hand , when the spindle motor 3 is started at one of current values i_2 , i_3 , and i_4 , the control skips step s_4 .
- 正解文
 - なお、電流値 i_2, i_3, i_4 のいずれかでスピンドルモータ3が起動した場合には、ステップ s_4 はスキップされる。
- 従来手法
 - 一方、スピンドルモータ3の一方の電流値 i_2, i_3, i_4 , スキップ制御が開始されると(ステップ s_4)、
- 提案手法
 - 一方、スピンドルモータ3が起動されると、ステップ s_4 で、制御電流値の i_2, i_3, i_4 である。

翻訳結果(log-linear model)による統合)

悪化した例

- 英文
 - The reader unit 1 is further provided with an **operation** unit 115 for effecting various settings on the **composite image input / output apparatus** .
- 正解文
 - また、リーダ部 1 には、本 **複合画像入出力装置** に対して各種設定を行うための **操作部** 115 が設けられている。
- 従来手法
 - また、リーダ部 1 が設けられた **操作部** 115 の **複合画像入出力装置** の各種設定を行う。
- 提案手法
 - また、リーダ部 1 の **複合画像入出力装置** の各種設定を行うための **演算部** 115 が設けられている。

まとめと今後の課題

- まとめ

- 言語横断LDAモデルを用いた統計的機械翻訳システム

- リランキング法
- log-linear model による統合
- どちらもベースラインをわずかに上回る性能

- 文脈適応範囲

- 翻訳元の対応する文のみ適応で良い性能

- 今後の課題

- ドメインが狭いとトピックが効きにくい

- 多様なドメインを含むコーパスによる実験