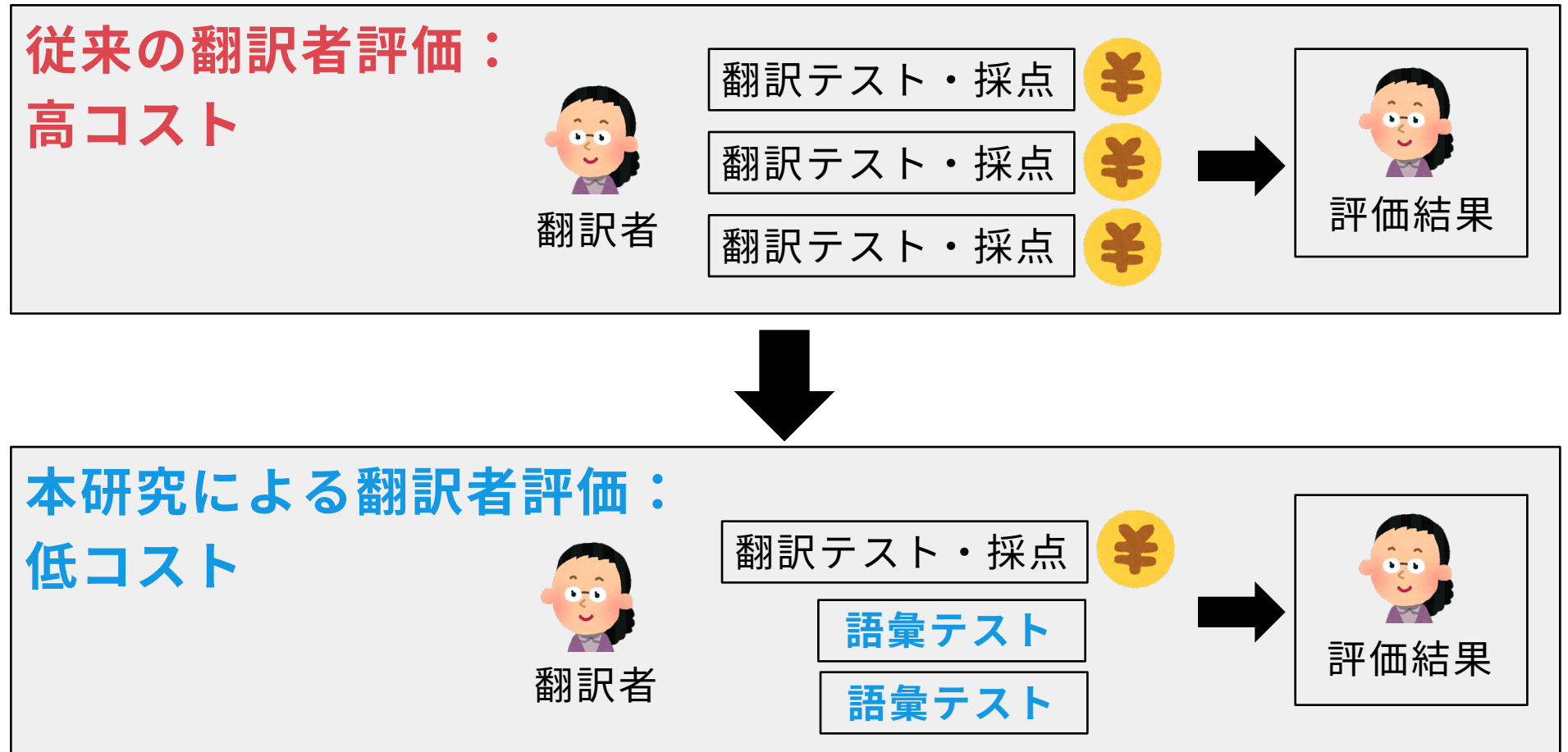


Assessing Translation Ability through Vocabulary Ability Assessment (appeared in IJCAI'16)

江原 遥 (産総研), 馬場 雪乃 (京大)
内山 将夫 (NICT), 隅田 英一郎 (NICT)

2016年8月9日 ERATO感謝祭 Season III

語彙テストを用いた、簡便な翻訳者評価法を提案



クラウドソーシング翻訳では品質が課題

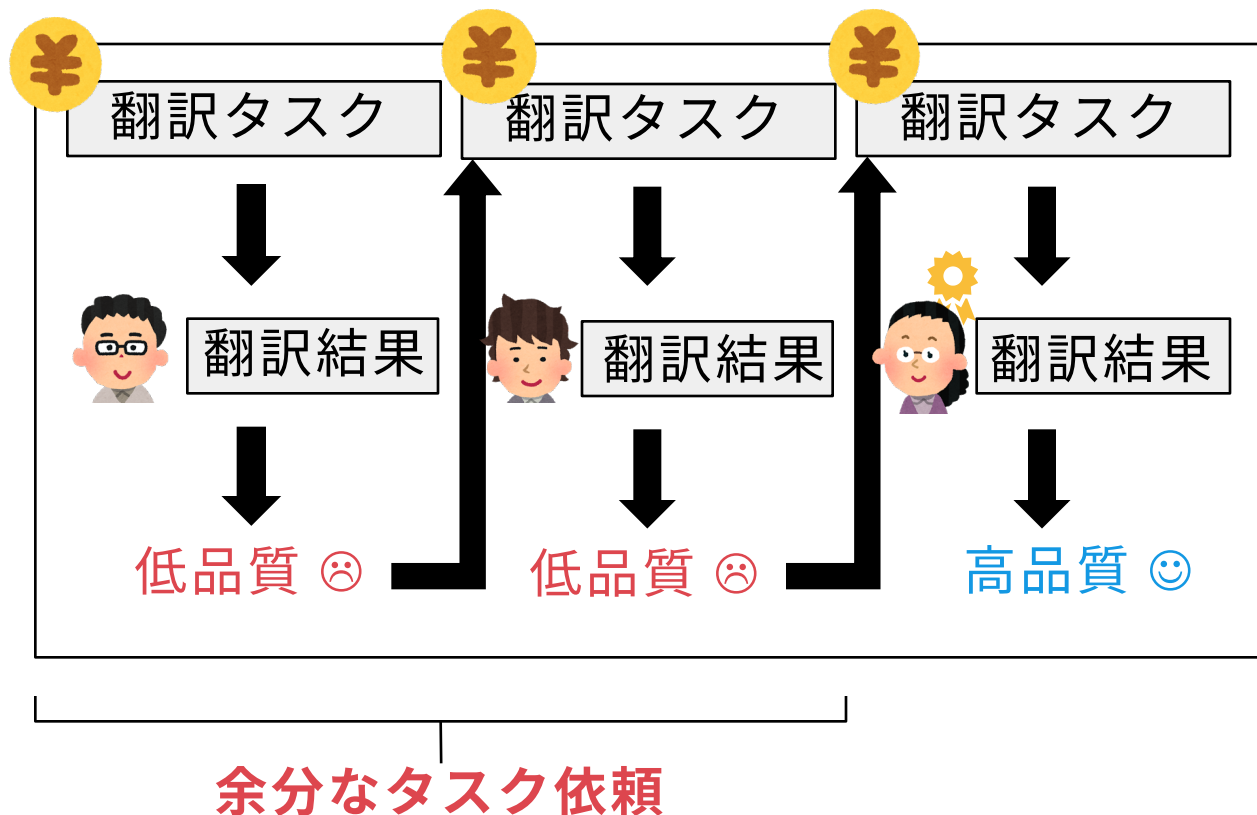
- クラウドソーシング翻訳では、非専門家も翻訳作業に参加
- 非専門家が、低品質の翻訳を返すことがしばしばある

クラウドソーシング翻訳の結果例

原文	The four brothers died of smallpox in 737 and there was a rumor that it was a punishment for driving the prince to suicide.
良い翻訳結果	4人の兄弟は737年に天然痘で死亡し、それは皇子を自殺に追い込んだことへの罰だといううわさが広まった。
悪い翻訳結果	4人の兄弟は737年スモールポックスで死んだ。そして自殺するためにプリンスを運転して消えたといううわさがたった。

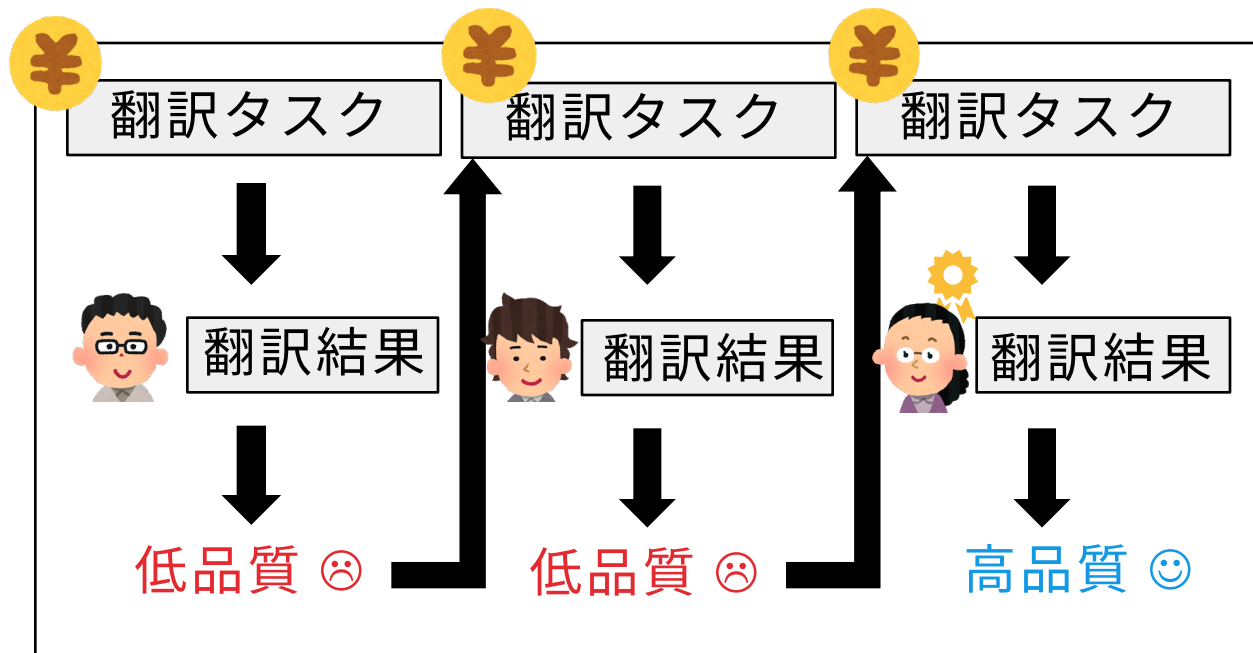
低品質な翻訳は依頼回数（＝費用）を増加させる

- 十分な品質の翻訳が得られるまでタスクを投げ続ける場合、低品質な翻訳は余分なタスク依頼を生じさせる



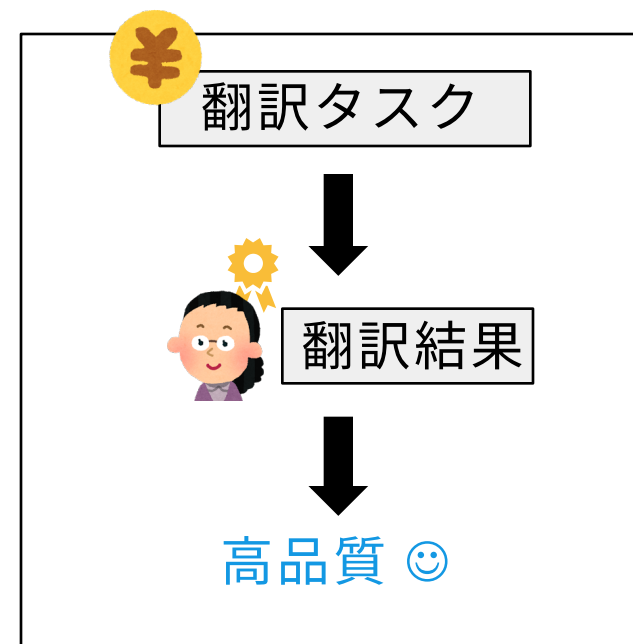
良い翻訳者が事前にわかれば依頼を効率化できる

「良い翻訳者」が
事前にわからない場合



3回のタスク依頼が必要

「良い翻訳者」が
事前にわかる場合



1回のタスク依頼で済む

翻訳テストによる翻訳者評価は正確だが高コスト

- 翻訳者評価の従来法：翻訳テスト
- 翻訳テストは正確だが採点コストが掛かりスケールしない
- 簡便な翻訳者評価方法が求められている

翻訳テストの例

Translate the following sentences:
Q1: "The accident triggered the ..."
その事件がきっかけとなり...

Q2: "His style was handed down ..."
彼の手法は江戸時代の終わりから続く...



専門家が採点



自己申告語彙テストの利用で採点コストを減らす

- 自己申告語彙テスト
 - 各単語の意味を知っているかどうかを答える
 - 自己申告なので、専門家による採点は不要
- 翻訳能力と語彙テストの回答を関連づける
確率モデルを構築

自己申告語彙テストの例

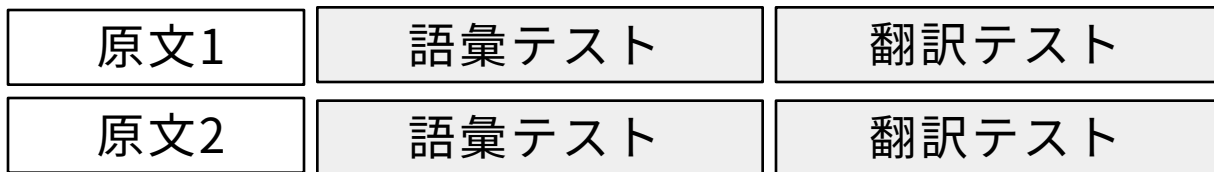
意味を知っている単語をチェックしてください

Q1. “The accident triggered the establishment of the Law for the ...”



翻訳・語彙テストを同時実施、一部の翻訳だけ採点

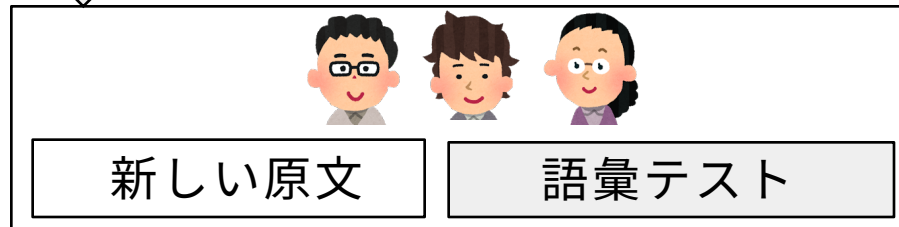
準備フェーズ



少数の翻訳テスト
結果のみ
専門家が採点

実行フェーズ

新しい翻訳タスク



モデル



翻訳・語彙テストの結果から品質予測器を学習

- 入力

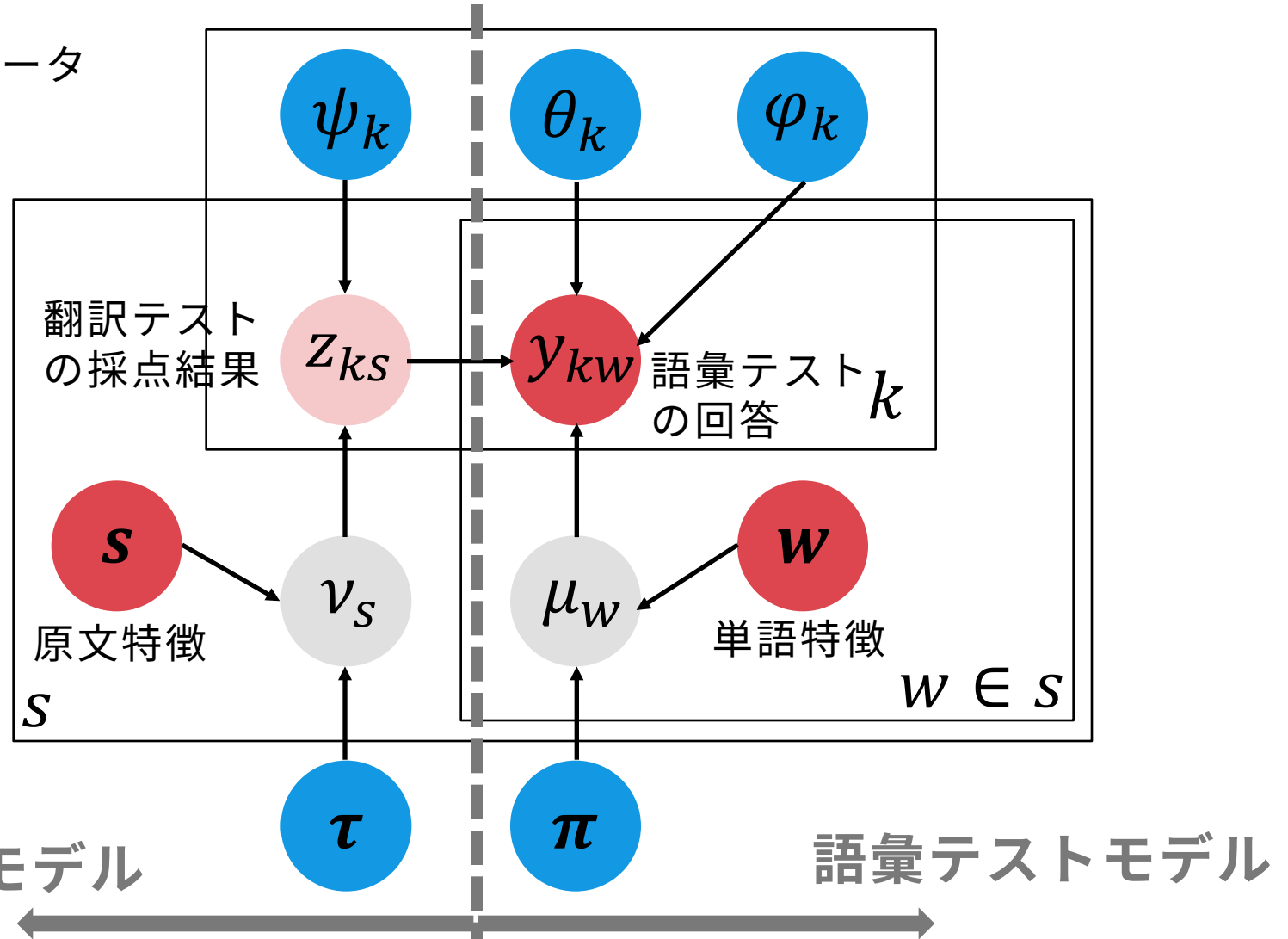
$\{k\}_k$	翻訳者集合
$\{s\}_s$	翻訳・語彙テストの原文集合, 原文は単語の集合で表現
$\{z_{ks}\}_{k,s}$	翻訳テストの採点結果の集合, $z_{ks} \in \{0, 1, "?\}$
$\{y_{kw}\}_{k,w \in s}$	語彙テストの回答結果の集合, $y_{kw} \in \{0, 1\}$
$\{s\}_s$	原文特徴の集合
$\{w\}_{w \in s}$	単語特徴の集合

- 出力：品質予測器 $\Pr [z_{ks^*} = 1 \mid \{y_{kw}\}_{w \in s^*}]$
(新しい原文 s^* に対して翻訳者 k が、
高品質の翻訳を返す確率を予測)

提案モデル

翻訳・語彙テストの結果を関連づける確率モデル

- モデルパラメータ
- 観測変数



翻訳能力と原文難易度に従い翻訳テスト結果が決まる

- 翻訳者 k が原文 s に対して高品質の翻訳を返す確率をラッシュモデルで表現

$$\Pr [z_{ks} = 1] = \frac{1}{1 + \exp(-(\psi_k - \nu_s))}$$

翻訳能力

原文難易度

語彙能力と単語難易度に従い語彙テスト回答が決まる

- 翻訳者 k が単語 w を「知っている」と答える確率をラッシュモデルで表現

$$\Pr [y_{kw} = 1 \mid z_{ks} = 1] = \frac{1}{1 + \exp(-(\theta_k - \mu_w))}$$

$$\Pr [y_{kw} = 1 \mid z_{ks} = 0] = \frac{1}{1 + \exp(-(\phi_k - \mu_w))}$$

語彙能力

単語難易度

原文・単語の難易度は各特徴に従って決まる

- 原文・単語難易度のモデル：

原文難易度

$$\nu_s = \pi^\top s$$

原文特徴への重み 原文特徴

単語難易度

$$\mu_w = \tau^\top w$$

単語特徴への重み 単語特徴

- 最終的なモデルパラメータ：
各翻訳者 k の**翻訳能力** (ψ_k), **語彙能力** (θ_k, ϕ_k),
原文特徴への重み (π), **単語特徴への重み** (τ)

翻訳能力と語彙能力を近づける事前分布を導入

- モデルの尤度：
 $L(\{\psi_k\}, \{\theta_k\}, \{\phi_k\}, \boldsymbol{\pi}, \boldsymbol{\tau})$
 $= \prod_k \mathcal{N}(\psi_k | \theta_k, \lambda^{-1}) \mathcal{N}(\theta_k | 0, \lambda^{-1}) \mathcal{N}(\phi_k | 0, \lambda^{-1})$
 $\times \mathcal{N}(\boldsymbol{\pi} | \xi^{-1} \mathbf{I}) \times \mathcal{N}(\boldsymbol{\tau} | \xi^{-1} \mathbf{I})$
 $\times \prod_k \prod_{s \in \mathcal{S}_k} \prod_{w \in s} \Pr[y_{kw} | z_{ks}] \Pr[z_{ks}]$

S_k : 翻訳者 k が翻訳した原文集合
 λ and ξ : ハイパーパラメータ

翻訳能力 ψ_k が
語彙能力 θ_k に近づく
ような事前分布を導入
- EMアルゴリズムでモデルパラメータを推定
 - 翻訳テストの採点結果 z_{ks} が既知の場合、
 $\mathbb{E}[z_{ks}] = \Pr[z_{ks} = 1 | \{y_{kw}\}_{w \in s}]$ の代わりに z_{ks} を直接使う

データセット

クラウドソーシングで1,498件の英日翻訳結果を収集

- 「Wikipedia日英京都関連文書対訳コーパス」から英文を選択
- ランサーズで英日翻訳タスクを依頼（1件10円）
 - － 翻訳時に語彙テストも受験させる
- 各翻訳結果に品質ラベル（「高品質」 or 「低品質」）を付与

データセット概要

原文数	104
翻訳結果数	1,498
翻訳者数	55

言語非依存の単純な特徴を利用

原文特徴

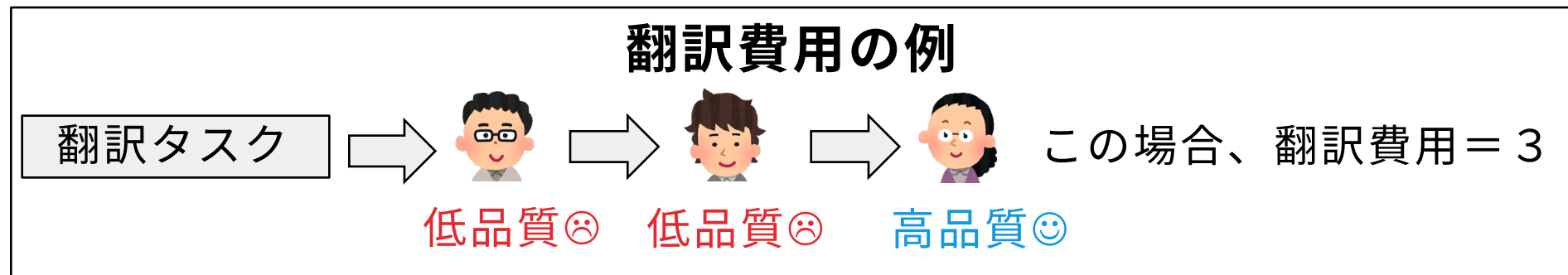
文章中の語数
文章中のカンマ数
文のパープレキシティ
文中の未知語数

単語特徴

Google n-gramコーパスにおける、1-gram確率の負の対数尤度
現代アメリカ英語コーパスCOCAにおける、1-gram確率の負の対数尤度
Brownコーパスそれぞれにおける、1-gram確率の負の対数尤度
12段階の単語難易度指標 (SVL 12000)

高品質の翻訳を得るまでの翻訳費用で評価

- 提案法により「高品質の翻訳を返す確率が高い」順に翻訳者を並べることができる
- 出力された順に翻訳を依頼した場合の翻訳費用で評価
 - － 翻訳費用：
「高品質」ラベル付きの翻訳に到達するまでに必要な翻訳者数



ロジスティック回帰とSVMをベースラインに使用

ベースラインと提案法のバリエーション

	各手法で用いる情報		
	原文特徴	単語特徴	語彙テスト結果
ランダム			
SVM-S	✓		
SVM-SW	✓	✓	✓
LR-S	✓		
LR-SW	✓	✓	✓
OURS-in	✓	✓	✓
OURS	✓	✓	✓

OURS-in:

提案法のバリエーション。翻訳能力と語彙能力を関連づけない

結果 (1)

提案法は、翻訳テストの採点率が小さい場合（採点コストを掛けない場合）に翻訳費用の削減効果大

翻訳費用の対ランダム比

	翻訳テスト結果の採点率			
	0%	1%	6%	12%
ランダム	1.000	1.000	1.000	1.000
SVM-S	-	0.825	0.703	0.672
SVM-SW	-	0.756	0.679	0.659
LR-S	-	0.820	0.706	0.669
LR-SW	-	0.731	0.678	0.660
OURS-in	0.908	0.736	0.699	0.661
OURS	0.673	0.660	0.675	0.664

結果 (2)

翻訳能力と語彙能力の関連付けは 予測精度向上に効果アリ

翻訳費用の対ランダム比

	翻訳テスト結果の採点率			
	0%	1%	6%	12%
ランダム	1.000	1.000	1.000	1.000
SVM-S	-	0.825	0.703	0.672
SVM-SW	-	0.756	0.679	0.659
LR-S	-	0.820	0.706	0.669
LR-SW	-	0.731	0.678	0.660
OURS-in	0.908	0.736	0.699	0.661
OURS	0.673	0.660	0.675	0.664

結果 (3)

語彙テスト結果は他手法の予測精度向上にも有効

翻訳費用の対ランダム比

	翻訳テスト結果の採点率			
	0%	1%	6%	12%
ランダム	1.000	1.000	1.000	1.000
SVM-S	-	0.825	0.703	0.672
SVM-SW	-	0.756	0.679	0.659
LR-S	-	0.820	0.706	0.669
LR-SW	-	0.731	0.678	0.660
OURS-in	0.908	0.736	0.699	0.661
OURS	0.673	0.660	0.675	0.664

語彙テストを用いた、簡便な翻訳者評価法を提案

- 自己申告語彙テストを用い、採点コストを掛けずに翻訳者能力ひいては翻訳品質を予測
- 翻訳・語彙テストの結果を関連づける確率モデルを提案
- 結果
 - － 提案法はクラウドソーシング翻訳費用の削減に貢献
 - － 翻訳能力と語彙能力の関連付けが有効