

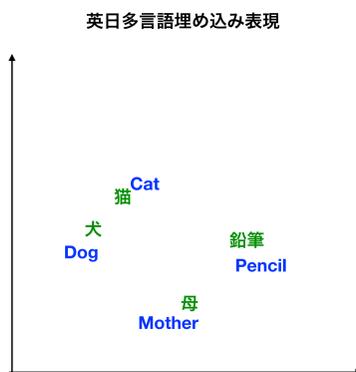
## 背景

### 自然言語処理での資源における言語間格差

- ・ 共通意味空間上で単語を扱う多言語埋め込み表現の登場により資源の豊富な言語からのモデル転移が可能

### 多言語埋め込み表現学習の問題点

- ・ 線形写像での学習には限界がある
- ・ 学習に必要なコンパラブルコーパスの量は限られている

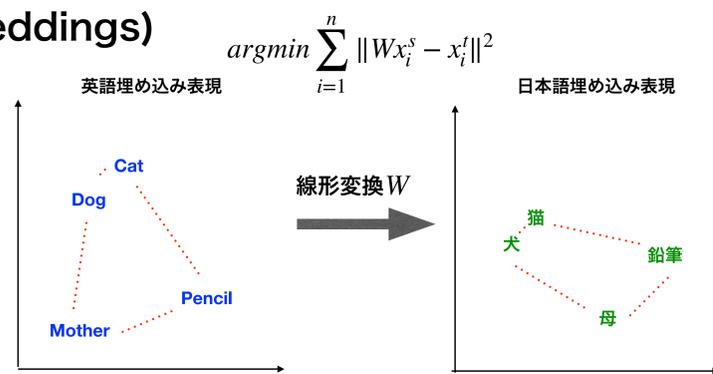


## 関連研究

### 多言語埋め込み表現 (Bilingual Word Embeddings)

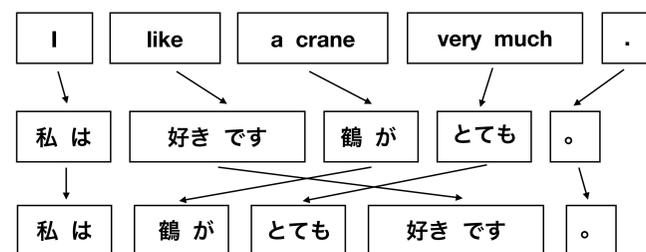
- ・ 別々に学習した異なる言語の埋め込み表現は単語埋め込み空間の構造が似る

→ 線形写像を用いて同一空間上に写像 [Artetxe et al., 2019]



### 教師なし機械翻訳モデル

- ・ 対訳コーパスなしで機械翻訳を実現 [Lample et al., 2018]
- ・ 多言語埋め込み表現による単語対応情報と言語モデルによる文法情報による実現



### 教師なし機械翻訳モデルの出力結果を利用した関連研究

- ・ 擬似対訳コーパスにbilingual skip-gramを適用する手法 (joint-learning) [Marie et al., 2019]

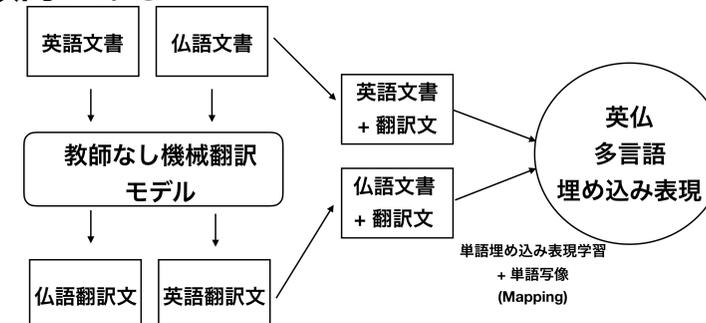
## 提案手法

### 機械翻訳の翻訳文は翻訳元の性質を反映する傾向にある

→ 翻訳文を学習データとして利用

↓  
単語埋め込み表現学習時の単語共起情報が似通い単語埋め込み空間の構造が似通うことを期待

### 提案モデル



## 実験

### コーパス

- ・ Wikipedia Comparable Corpora (英語・フランス語・ドイツ語・日本語) より1000万文を利用

### 翻訳モデル

- ・ 教師なしフレーズベースモデルを利用
- ・ Bilingual Lexicon Induction にて単語写像性能評価
- ・ Eigen Similarityで単語埋め込み空間の類似度評価

### 評価指標

### 実験結果

既存手法とのBLIスコアの比較

データ学習方法	en→fr		fr→en		en→de		de→en		en→ja		ja→en	
	BLI	Eigen Sim	BLI	Eigen Sim	BLI	Eigen Sim	BLI	Eigen Sim	BLI	Eigen Sim		
訓練データのみ Mapping	0.664	0.636	0.561	0.567	0.451	0.357						
擬似対訳文 Joint-learning	0.620	0.594	0.527	0.520	0.263	0.273						
訓練+翻訳文 Mapping (提案手法)	<b>0.696</b>	<b>0.669</b>	<b>0.637</b>	<b>0.612</b>	<b>0.488</b>	<b>0.418</b>						

翻訳文ごとのBLIスコアとEigen Similarity

	Fr		De		Ja	
	BLI	Eigen Sim	BLI	Eigen Sim	BLI	Eigen Sim
訓練データのみ	0.655	22	0.548	31	0.451	237
訓練+擬似文 (Fr)	<b>0.704</b>	<b>10</b>	0.561	33	0.448	188
訓練+擬似文 (De)	0.647	23	<b>0.610</b>	<b>15</b>	0.445	242
訓練+擬似文 (Ja)	0.598	25	0.496	152	<b>0.463</b>	<b>130</b>

- ・ 翻訳元の性質を反映した翻訳文を学習に用いることで写像精度が向上