

スタイル変換技術による対訳コーパスから 同時通訳コーパスへの拡張

2019/08/30

奈良先端科学技術大学院大学
二又航介、須藤克仁、中村哲

背景：同時通訳 (Simultaneous interpretation)

原言語の入力文の終了を待たずに目的言語への訳出を開始

- 通常の翻訳システムでは原言語の入力終了を待つため遅延が発生
- 講義や講演など主に訳出の遅延が許されない場面で使用
- 遅延を最小限にしつつ正確に部分訳出を行う

同時通訳による訳出例

A brand-new computer on the desk / which my father gave me on my birthday / doesn't work now.

原言語文

机の上にある新しいコンピュータですね、

これは父から誕生日にもらったものです、

ですが、今故障しています。

目的言語文

背景：通訳方法の違いによる遅延

英日翻訳のように語順が大きく異なる言語間の翻訳では 訳出開始までに遅延が発生

- 主要部先行型言語(head-initial)と主要部後続型言語(head-final)の違いによる遅延
- 長い修飾部を持つ英文の訳出開始までに大きな遅延が発生
- 原文の節や句の順序を守りながら原言語の語順に近い形で訳出(順送り)することで遅延が少なくなる^[1]
- 順送り方式による訳出では語順の洗練性は無いが助詞等による致命的な間違いが発生しない

[1] 水野的: 同時通訳の理論 — 認知的制約と訳出方略, 朝日出版社 (2015).

背景：通訳方法の違いによる遅延

訳出方法の違いによる遅延

A brand-new computer on the desk which **my father gave me** on my birthday doesn't work now .

父から誕生日に貰った、机上有る新しいコンピュータは今故障しています。

(待ち時間...) 父から誕生日に貰った、机上有る新しいコンピュータは...

訳出開始までに大きな遅延が発生する例

A brand-new computer **on the desk** which my father gave me on my birthday doesn't work now .

机上有る新しいコンピュータですね、これは父から誕生日に貰ったものです、ですが今故障しています。

(待ち時間...) 机上有る新しいコンピュータですね、これは父から誕生日に貰ったものです...

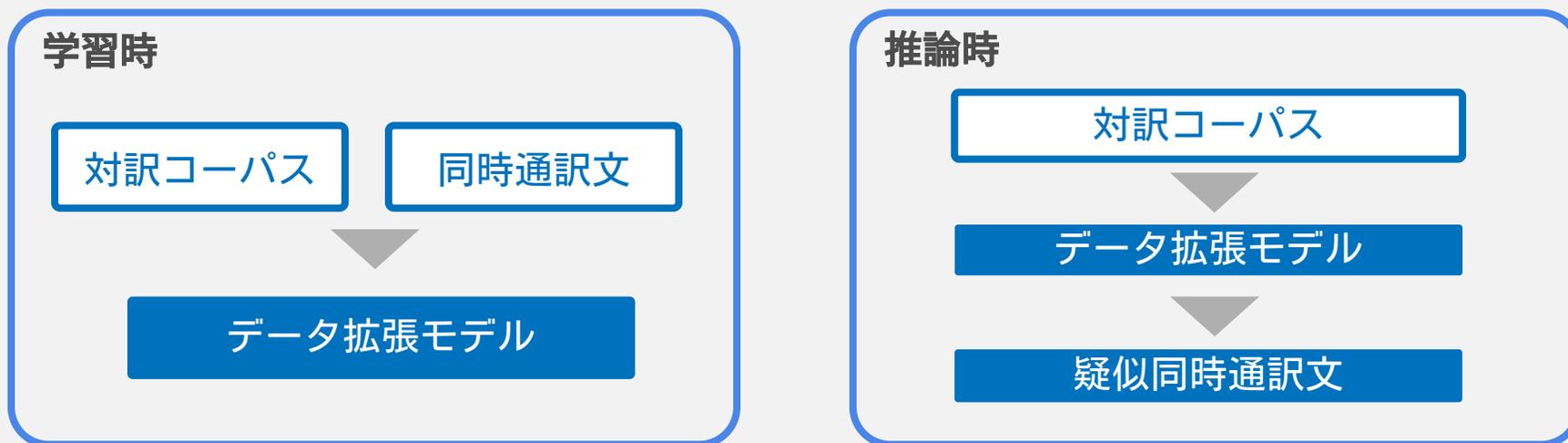
訳出開始までの遅延が少ない例(順送り方式)

目的：対訳コーパスから同時通訳コーパスの拡張

スタイル変換技術によって対訳コーパスから順送り方式の同時通訳コーパスを作成

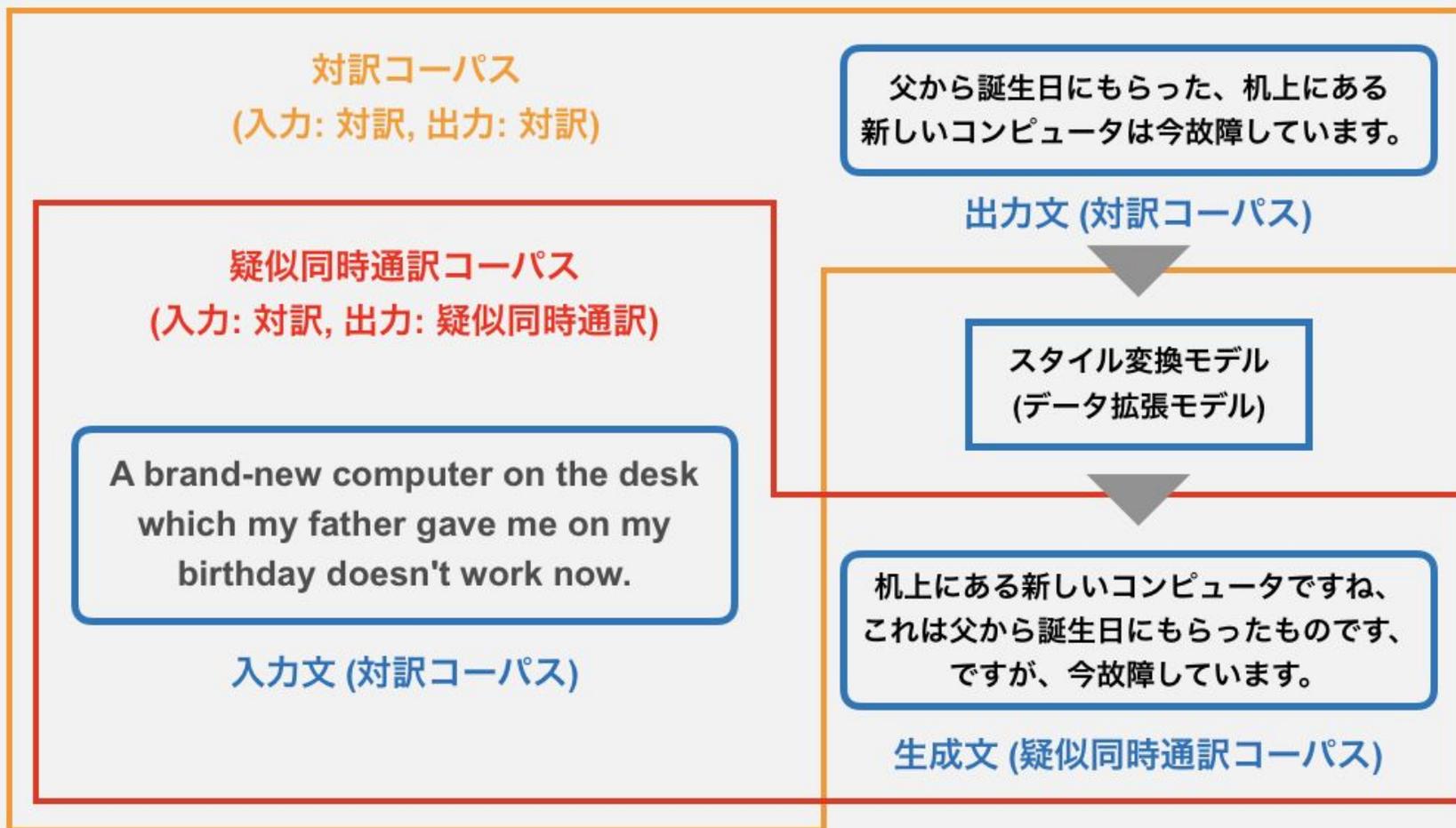
- 順送り方式で翻訳されたコーパスは少数
- 大量の対訳コーパスは利用可能
- 対訳コーパスの目的言語を順送りの同時通訳文にスタイル変換

対訳コーパスから同時通訳文の拡張過程



目的： 擬似英日同時通訳コーパスの作成

英日対訳コーパスと英日疑似同時通訳コーパスの関係

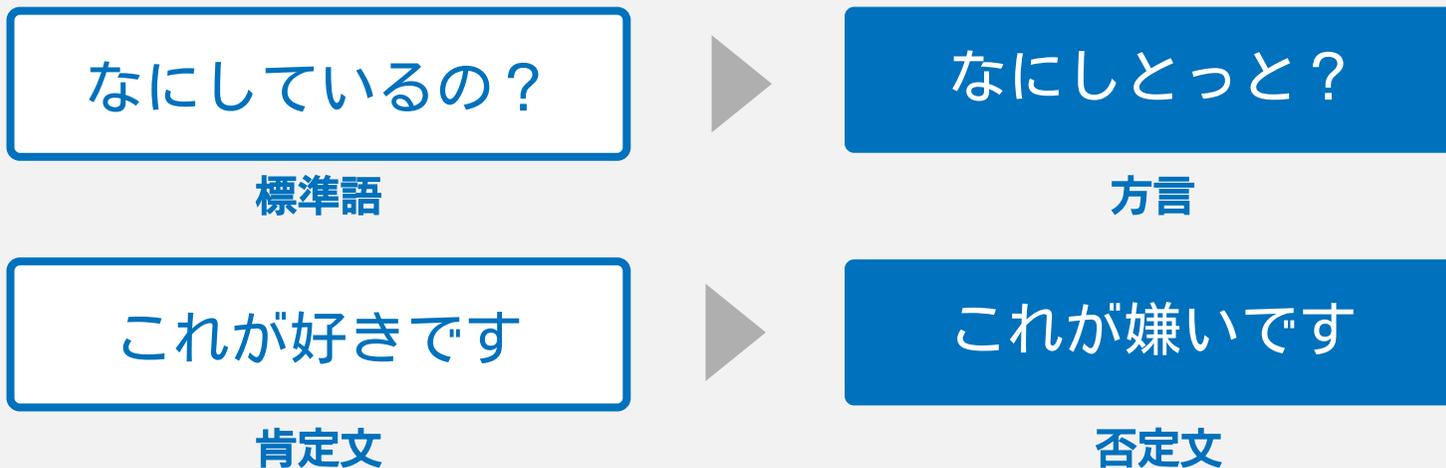


関連研究：スタイル変換

入力文における文意や意図を変更することなく文体(スタイル)を自動的に制御するタスク^[2]

- 入力文を意味的に等価な文へ書き換える言い換え生成の一種
- スタイル変換前後の対訳ペアを必要としない
- 変更するスタイル情報が明確な場合に特に有効

スタイル変換例



関連研究：スタイル変換の手法

1. 文意とスタイルを分離する手法^[3]

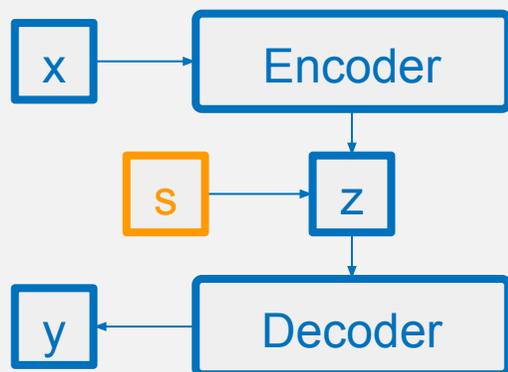
- スタイル変換の一般的な手法
- スタイル情報と独立した文意の潜在表現を学習

2. 文意とスタイルを分離せず直接スタイル変換する手法^[4]

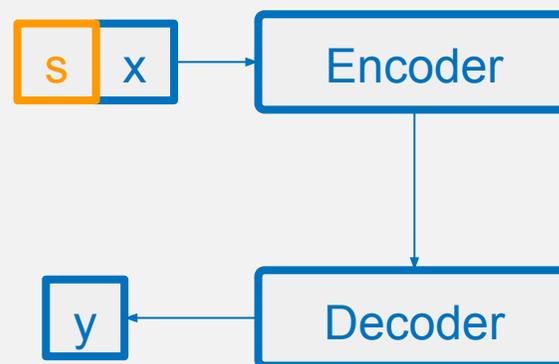
- 潜在表現には文意とスタイル情報が含まれる
- あるスタイルから異なるスタイルへ直接変換

x: 文章(スタイル1)
y: 文章(スタイル2)
s: スタイル識別子
z: 潜在表現

2種類のスタイル変換手法



文意とスタイルを分離する手法



直接スタイル変換する手法

[3] Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, Rui Yan, Style Transfer in Text: Exploration and Evaluation, 2017

[4] Guillaume Lample, Sandeep Subramanian, Eric Smith, Ludovic Denoyer, Marc'Aurelio Ranzato, Y-Lan Boureau, Multiple-Attribute Text Rewriting, 2019

関連研究: Style Transformerによるスタイル変換

Style Transformer^[5]

- Transformerをベース
- 直接スタイルを変換
- 意味的な単語を変換する

x : 入力文(スタイル1)

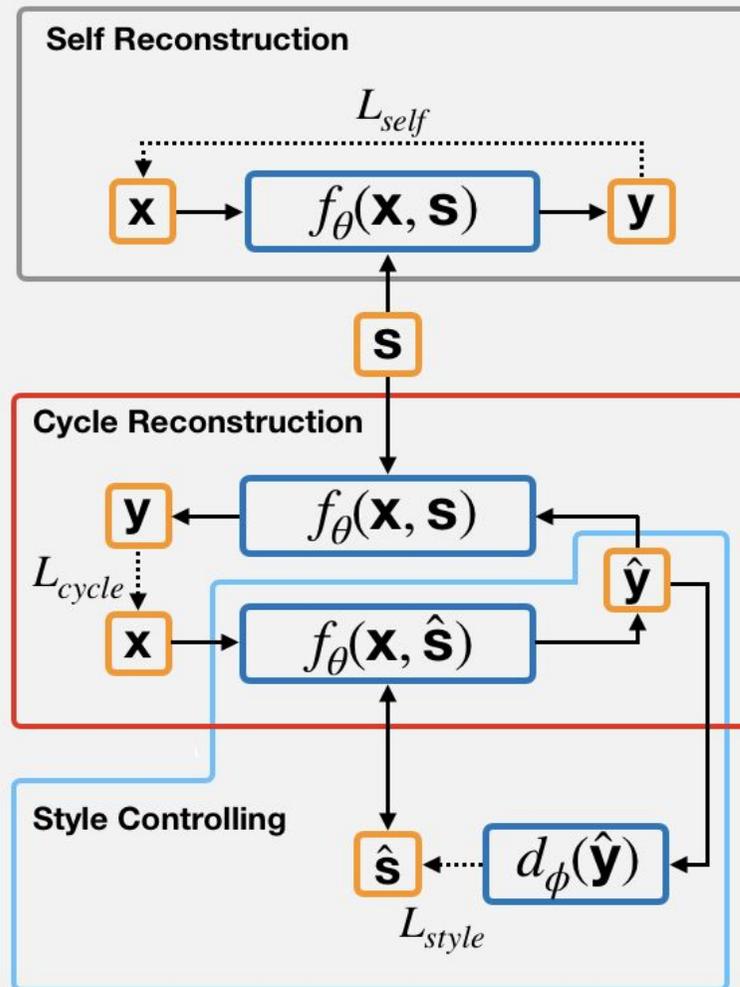
y : 変換文(スタイル1)

\hat{y} : 変換文(スタイル2)

s : 変換前のスタイル識別子

\hat{s} : 変換後のスタイル識別子

Style Transformerの構成



[5] Ning Dai, Jianze Liang, Xipeng Qiu, Xuanjing Huang, Style Transformer: Unpaired Text Style Transfer without Disentangled Latent Representation, 2019

関連研究: Style Transformerによるスタイル変換

Self Reconstruction

入力文(\mathbf{x})とスタイル(\mathbf{s})は同じ

$$\mathcal{L}_{self}(\theta) = -p_{\theta}(\mathbf{y} = \mathbf{x} | \mathbf{x}, \mathbf{s})$$

Cycle Reconstruction

入力文(\mathbf{x})とスタイル($\hat{\mathbf{s}}$)が異なる

入力文を異なるスタイル($\hat{\mathbf{s}}$)で変換後($\hat{\mathbf{y}}$)
スタイル(\mathbf{s})によって入力文を復元(\mathbf{y})

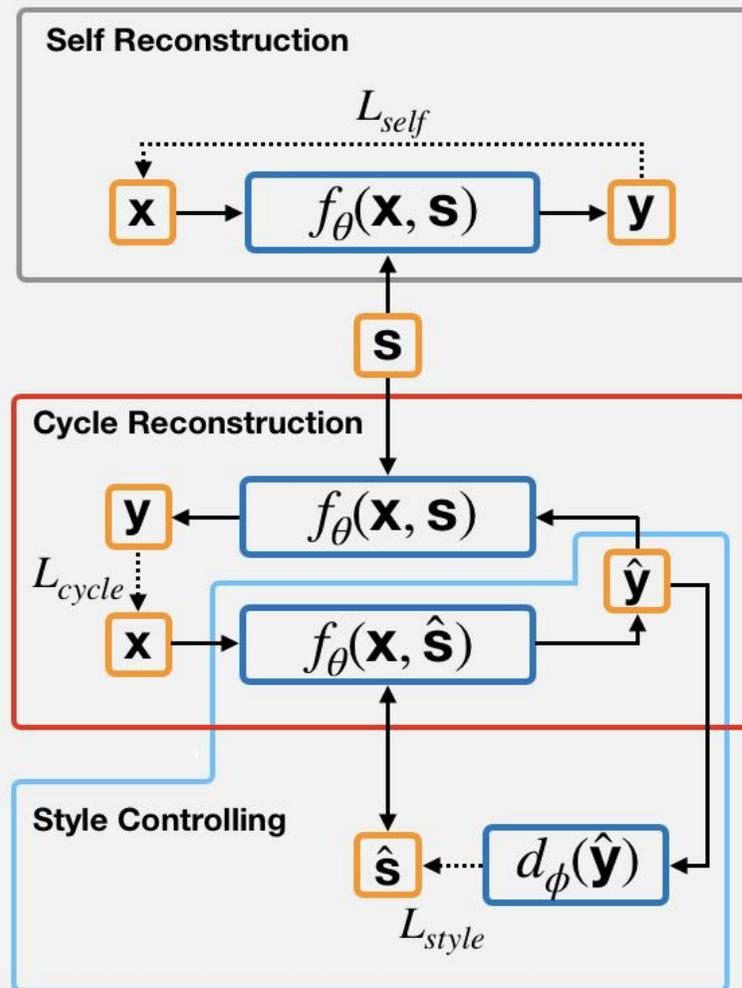
$$\mathcal{L}_{cycle}(\theta) = -p_{\theta}(\mathbf{y} = \mathbf{x} | f_{\theta}(\mathbf{x}, \hat{\mathbf{s}}), \mathbf{s})$$

Style Controlling

Cycle Reconstructionにより $\hat{\mathbf{y}}$ の
文意が大きく変わるのを防ぐ

$$\mathcal{L}_{style} = -p(\mathbf{c} = \hat{\mathbf{s}} | f_{\theta}(\mathbf{x}, \hat{\mathbf{s}}))$$

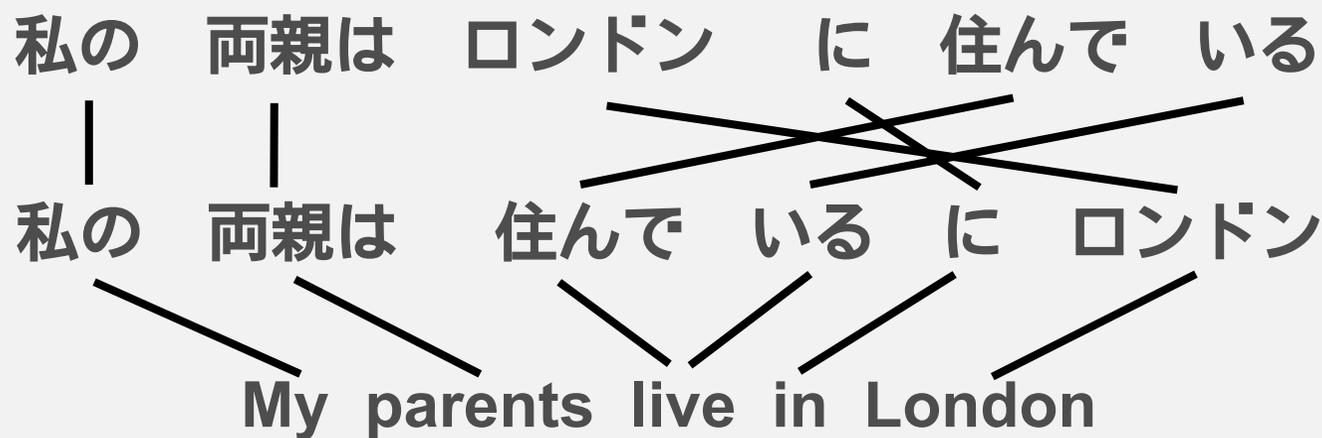
Style Transformerの構成



翻訳機に入力する前に原言語における文の語順を目的言語の語順に近づくように並び替える手法

- 主に統計的機械翻訳(SMT)で使用
- 順送り方式の同時通訳文は原言語の語順と類似している
- 事前並べ替えによって原言語との単語間の交差が少なくなる

日本語文に対する事前並べ替えの適用例



[6] Tetsuji Nakagawa, 2015, Efficient Top-Down BTG Parsing for Machine Translation Preordering, in Proceedings of ACL, pages 208-218.

提案手法：スタイル変換と事前並べ替えによるデータ拡張

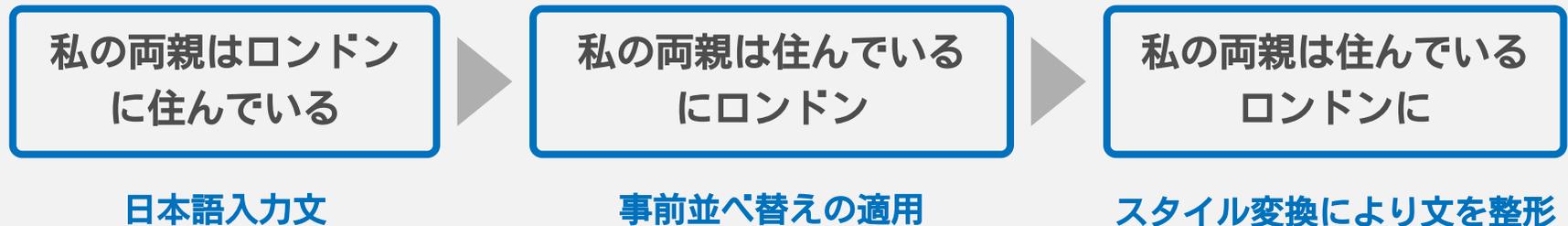
スタイル変換前の対訳コーパスに事前並び替えを適用

- スタイル変換は単語の置換や単語の削除など表現力が限定的であり語順の並べ替えには不向き
- スタイル変換前に語順を並べ替えることでスタイル変換を容易に

スタイル変換を用いたデータ拡張の過程

1. 日本語入力文を英文の語順のように事前並べ替え
2. スタイル変換により文を正しく整形

提案手法による同時通訳コーパスの拡張方法



実験設定：使用データ

Style Transformerと事前並び替えによって英日対訳コーパスから英日同時通訳コーパスへのスタイルを変換

- 事前並べ替えの適用あり/なしの2種類で実験
- 対訳コーパス：日本語話し言葉コーパス(CSJ)
- 同時通訳コーパス：TEDコーパス(独自に収集した同時通訳文)
- 事前並べ替え：ASPEC

CSJ,TED,ASPECのコーパスサイズ

Corpus	Number of sentences		
	Train	Val	Test
CSJ	3509	500	500
TED	28606	500	500
ASPEC	1000000	1790	1812

スタイル変換を行うコーパスペア

Corpus pair	Detail
CSJ-TED	CSJ から TED
CSJ(preordered)-TED	事前並べ替えした CSJ から TED
CSJ-CSJ(preordered)	CSJ から事前並べ替えした CSJ

スタイル変換で主に使用される3つの評価指標

Style accuracy

- 二値分類器(対訳or同時通訳)によってスタイル分類精度を計測
- スタイル変換された文(疑似同時通訳文)が同時通訳のスタイルとして識別されることを期待
- CSJとTED(同時通訳文)により学習

Bleu score

- Bleu scoreによってスタイル変換後の文意の保持具合を計測
- スタイル変換された文(疑似同時通訳文)が入力文と同一の文意を保持していることを期待

Perplexity

- Perplexityによってスタイル変換後の文章の流暢さを計測
- スタイル変換された文(疑似同時通訳文)が入力文と同様に流暢であることを期待

自動評価による実験結果

CSJ-TED: 語尾や単語を変更する傾向

- 文意と流暢さが保持されたため高BLEU, 低Perplexity

CSJ(preordered)-TED: 語順変化&単語追加する傾向

- 語順変化&単語追加されたため低BLEU, 高Perplexity

CSJ-CSJ(preordered): 長い文を短い単位に区切る傾向

- 語順変化&流暢さが損なわれたため低BLEU, 高Perplexity

自動評価指標による実験結果

Corpus pair	ACC	BLEU	PPL
CSJ-TED	89.0	48.8	73.5
CSJ(pre-ordered)-TED	96.0	21.8	361.5
CSJ-CSJ(preordered)	30.2	28.3	242.7

実験結果: CSJ-TEDの生成例

CSJからTEDへのスタイル変換

- 語尾や単語を同時通訳コーパスに現れるものに変換する傾向
- Style accuracy: 89.0, Bleu score: 48.8, Perplexity: 73.5

入力文(CSJ)	世界 戦争 がヨーロッパ から 始まりました 。
変換文(TED)	世界 戦争 がヨーロッパ から 始まり ます 。

入力文(CSJ)	私 が 言わなくちゃいけない 内容 に入る 前に, 少し 私 の 自己 紹介 です 。
変換文(TED)	私 が 言わなくちゃいけない 内容 に入る 前に, 少し 僕 の 自己 紹介 です 。

入力文(CSJ)	そして, この 町 の 中で 最悪 の 裁判 所 です 。
変換文(TED)	そして, この アフリカ の 中で 最悪 の 裁判 所 です 。

実験結果: CSJ(preordered)-TEDの生成例

CSJ(preordered)からTEDへのスタイル変換

- 語順変更や単語追加が行われるが文意が大きく変わる傾向
- Style accuracy: 96.0, Bleu score: 21.8, Perplexity: 73.5

原文(CSJ)	学習 データ は、こちらと同じものです。
入力文(CSJ(preordered))	学習 データ は、です 同じものはこちら。
変換文(TED)	学習 データ は ない です、同じものはこちら。

原文(CSJ)	一方 アジア 人は、日本人 学部 生 は、僕 一人で少ない方でした。
入力文(CSJ(preordered))	た 一方 アジア 人 は、日本人 学部 生 は、でし 僕 一人で方 少ない。
変換文(TED)	私 一方 アジア 人は、日本人 学部 生 、一人 で少ない。

実験結果: CSJ-CSJ(preordered)の生成例

CSJからCSJ(preordered)へのスタイル変換

- 長い文を短い単位に区切り同時通訳らしい文を生成する傾向
- Style accuracy: 30.2, Bleu score: 28.3, Perplexity: 242.7

入力文(CSJ)	講演 音声 認識 の 識別 率 は 、 今 の ところ 七十 パーセント 程度 です 。
参照文(CSJ)	識別 率 の 講演 音声 認識 は 、 です 程度 と ころ の 今 七十 パーセント 。
変換文 (CSJ(preordered))	講演 音声 認識 の 識別 率 は です 今 の ところ 七十 パーセント 程度 。

入力文(CSJ)	次に本 研究 におけます システム の 概要 を 説明 いたします 。
参照文(CSJ)	ます いたし 次に 概要 の おけます に本 研究 システム を 説明 。
変換文 (CSJ(preordered))	ついに 本 研究 におけます システム の 概要 を きちんと 説明 て 。

実験結果: CSJ-CSJ(preordered)の生成例

入力文(CSJ)	使用する関係としては、動詞 目的語 名詞 というこの三種類を使用します。
参照文(CSJ)	まずしは、を使用この三種類 いうと動詞 目的語 名詞 しと関係する使用で。
変換文 (CSJ(preordered))	使用する関係とし まず 動詞 目的語 名詞 というこの三種類を使用して。

入力文(CSJ)	これは、係り受けの係り先が間違っておりますので誤りとなっております。
参照文(CSJ)	まず おりなつと誤りてでのまず これ おり先 係りの は、係り受けがっ間違て。
変換文 (CSJ(preordered))	これは、係り受け です 係り先が間違っておりますので誤りとなり。

入力文(CSJ)	次に実際に決定木を構築しての選択し検索を行なう実験というのを行ないました。
参照文(CSJ)	たし選択のし次に実際決定木を構築にてまし行ないのいうと実験行なう検索をを。
変換文 (CSJ(preordered))	次に実際に決定木を構築し まず 選択し検索を行ない まず 実験というのを行ないました。

実験設定：人手評価

既存の自動評価指標で疑似同時通訳文を評価するのは困難

- 長い文を短く区切る傾向にあり最も同時通訳文らしいCSJ-CSJ(preordered)のペアを対象に評価実験
- 7人の被験者が50サンプルを評価
- スタイル変換前の文(入力文)とスタイル変換後の文(生成文)を提示
- 評価指標に当てはまるか1から5の5段階で評価

評価指標

- 入力文と比較し生成文が短い単位に区切られているか(Segmentation)
- 生成文は日本語として自然で流暢であるか(Fluency)
- 入力文と生成文が意味的に同一であるかどうか(Identity)

人手評価による実験結果

SegmentationとFluencyの間, SegmentationとIdentityの間にトレードオフの関係性

- Segmentation ≥ 3.0 , Fluency ≥ 3.0 , Identity ≥ 3.0 はそれぞれの平均値が3.0以上であったサンプル
- Segmentationの値が上がるとFluencyとIdentityの値が下がる
- Segmentationの値が下がるとFluencyとIdentityの値が上がる

7人の被験者による評価平均値

	Num samples	Segmentation	Fluency	Identity
All samples	50	3.03	2.42	3.11
Segmentation ≥ 3.0	30	3.72	2.02	2.95
Fluency ≥ 3.0	10	2.37	4.00	3.74
Identity ≥ 3.0	28	2.96	2.76	3.73

人手評価による実験結果

Segmentationの平均値が特に高くFluencyとIdentityの平均値が特に低い変換例

- 入力文を短い単位に区切る傾向にあるが日本語が不自然

入力文(CSJ)	第四声は、発音が高いところから最も低いところへ急激に移行していきます。
参照文(CSJ)	まずいきし第四声に急激移行へところは最も低いからところ、発音が高くて。
変換文 (CSJ(preordered))	四声は、発音が高い まず 、最も低い 事故 へ移行して。

入力文(CSJ)	大体このくらいで終わりにしたいと思います。
参照文(CSJ)	まず大体に終わりしたいと思いでこのくらい。
変換文 (CSJ(preordered))	まず 、 思い くらいで、終わりにしたいと思います よく 。

人手評価による実験結果

Segmentationの平均値が特に低くFluencyとIdentityの平均値が特に高い変換例

- 入力文が短い単位に区切られないが日本語として自然

入力文(CSJ)	特にですね自転車を使った旅が好きです。
参照文(CSJ)	です特に旅を使っねです自転車たが好き。
変換文 (CSJ(preordered))	特にですね自転車を使った旅が好き。

入力文(CSJ)	でこの町の中心である公園について御説明したいと思います。
参照文(CSJ)	でます思いとたいしについて御説明に公園ある中心のこの町。
変換文 (CSJ(preordered))	でこの町の中心である公園について御説明したいと思います。

スタイル変換手法によって対訳コーパスから同時通訳コーパスへと拡張する手法を提案

- スタイル変換と事前並べ替えにより疑似同時通訳文を生成
- 提案手法により長い入力文を短い単位に区切ることが可能
- 提案手法ではSegmentationとFluency, Identityの間にトレードオフ

同時通訳文の特性を考慮したモデル及び指標の導入

- 実際の同時通訳文には語順変化や単語の追加など様々な要因が関与
- 流暢性、意味的同一性を担保した疑似同時通訳文の生成が必要
- 疑似同時通訳文を評価するための新たな指標の導入