

時系列データの動向説明文生成における 参照時刻の不整合解消に向けた取り組み

濱園 侑美^{○,§} 上原 由衣[§] 能地 宏[§] 宮尾 祐介^{¶,§} 高村 大也^{§,†} 小林 一郎^{○,§}

[○]お茶の水女子大学 [¶]東京大学 [†]東京工業大学 [§]産業技術総合研究所

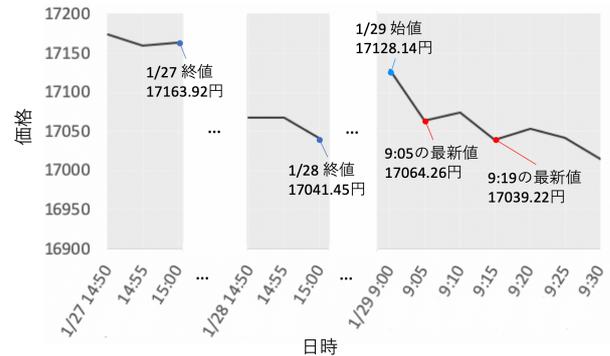
{hamazono.yumi,koba}@is.ocha.ac.jp, {yui.uehara,hiroshi.noji}@aist.go.jp,
takamura@pi.titech.ac.jp, yusuke@is.s.u-tokyo.ac.jp

1 はじめに

金融や医療、情報通信などの多くの分野において、様々な形式の非言語データを取り扱う機会が増えてきている。大規模で複雑なデータを専門知識のない人が解釈できるように、データを説明するテキストを自動的に生成する技術の必要性が高まっており、近年は様々なデータを題材に、Encoder-Decoder モデルを用いて end-to-end の学習を行うことで、高い生成性能を發揮している (漆原・小林, 2018; 村上他, 2017; Puzikov and Gurevych, 2018; Liu et al., 2018; Iso et al., 2019)。

日経平均株価の概況テキスト生成 (Murakami et al., 2017; Aoki et al., 2018) において, Murakami et al. (2017) は日経平均株価, Aoki et al. (2018) は日経平均株価に加えて, 日経平均株価の変化要因となりうる複数の経済市場データ計 L 種類を用意した。それぞれの経済市場データ data_i ($i = 1, 2, \dots, L$) に対して短期的な変動を捉えるために 5 分足で収集した直近の N 個の市況データ $\mathbf{x}_{\text{short}}^{\text{data}_i} = (x_{\text{short},1}^{\text{data}_i}, x_{\text{short},2}^{\text{data}_i}, \dots, x_{\text{short},N}^{\text{data}_i})$ と, 長期的な変動を捉えるための直近 M 取引日分の終値データ $\mathbf{x}_{\text{long}}^{\text{data}_i} = (x_{\text{long},1}^{\text{data}_i}, x_{\text{long},2}^{\text{data}_i}, \dots, x_{\text{long},M}^{\text{data}_i})$ を入力とした。

しかし, Murakami et al. (2017) および Aoki et al. (2018) において, $\mathbf{x}_{\text{short}}^{\text{data}_i}$ および $\mathbf{x}_{\text{long}}^{\text{data}_i}$ は, 概況テキストが発表された時刻を基準に取得しており, 発表時刻と参照時刻の不整合は無視されてきた。例えば, 図 1 に示した (I) から (III) の文章は, 全て「1/29 9:00」を参照時刻とする概況テキストであり, 概況テキストでの動向内容と参照時刻での値動きは一致している。(I) は, 発表時刻と参照時刻が同じため, 概況テキストの動向内容と記事の発表時刻を基準とした値動きは一致している。一方, (II) および (III) では, 発表時刻と参照時刻が異なるため, 概況テキストの動向内容と発表時刻での値動きは異なっている。また, サッカーの試合データを用いた速報テキスト生成 (Taniguchi et al., 2019) では, 試合データのイベントと速報テキストのア



	発表時刻	発表時刻での値動き	概況テキスト (動向内容)
(I)	09:00	反発, 86 円高	日経平均、 <u>反発</u> で始まる 86 円高
(II)	09:05	反発, 22 円高	日経平均、 <u>反発</u> で始まる 86 円高、原油高を好感 トヨタ高い
(III)	09:19	続落, 2 円安	東証寄り付き、 <u>反発</u> 原油高を好感、ファンックは大幅安

図 1: 日経平均の値動きと「1/29 9:00」を参照時刻とする概況テキスト (表)。

ライメントが正しく取得できないことによる誤りが発生した。

このように時系列データを用いた研究では, 取得したデータに参照時刻の不整合が起こり得る。本研究では, このような時系列データを用いた data-to-text で発生し得る, データとテキストの不整合を解消する手法を, 日経平均株価の概況テキスト生成を例に提案する。

2 日経平均株価の概況テキスト生成

Murakami et al. (2017) は, Encoder-Decoder モデル (Sutskever et al., 2014) に基づいて, 日経平均株価データとその動きに言及しているニュースヘッドラインを例に, 時系列数値データから概況テキストを生成するモデルを開発した。また, Aoki et al. (2018) は, Murakami et al. (2017) のモデルを拡張し, 複数の経済市況データ

を用いて、変化要因に関する記述を含む概況テキストを生成する課題に取り組んだ。以降、Aoki et al. (2018) について詳しく述べる。

1 節で述べた通り、 $\mathbf{X}_{\text{short}} = (\mathbf{x}_{\text{short}}^{\text{data}_1}, \mathbf{x}_{\text{short}}^{\text{data}_2}, \dots, \mathbf{x}_{\text{short}}^{\text{data}_L})$ と、 $\mathbf{X}_{\text{long}} = (\mathbf{x}_{\text{long}}^{\text{data}_1}, \mathbf{x}_{\text{long}}^{\text{data}_2}, \dots, \mathbf{x}_{\text{long}}^{\text{data}_L})$ を入力として用いる。これらに Murakami et al. (2017) の研究で有効とされた 2 つの前処理を適用する：

$$x_{\text{std}_i} = \frac{x_i - \mu}{\sigma}, \quad (1)$$

$$x_{\text{move}_i} = x_i - r_i, \quad (2)$$

$$x_{\text{norm}_i} = \frac{2 \times x_{\text{move}_i} - (\bar{x}_{\text{max}} + \bar{x}_{\text{min}})}{\bar{x}_{\text{max}} - \bar{x}_{\text{min}}}. \quad (3)$$

1 つ目は式 (1) で表される平均値と標準偏差を用いて標準化する手法であり、この処理を $\mathbf{X}_{\text{short}}$ および \mathbf{X}_{long} に適用し、得られたベクトルをそれぞれ $\mathbf{X}_{\text{short}}^{\text{std}}$ および $\mathbf{X}_{\text{long}}^{\text{std}}$ と表す。2 つ目は、前日の終値からの価格の変動を捉えるために、式 (2) によってそれぞれの入力データに対して前取引日の終値 r_i からの差分を計算し、差分の最大値 \bar{x}_{max} 、最小値 \bar{x}_{min} を用いて、式 (3) によって $[-1, 1]$ へ正規化を行う手法であり、この処理で得られたベクトルをそれぞれ $\mathbf{X}_{\text{short}}^{\text{move}}$ および $\mathbf{X}_{\text{long}}^{\text{move}}$ と表す。これらのベクトルにそれぞれ MLP を用いて、 $\mathbf{h}_{\text{short}}^{\text{std}}$, $\mathbf{h}_{\text{long}}^{\text{std}}$, $\mathbf{h}_{\text{short}}^{\text{move}}$ および $\mathbf{h}_{\text{long}}^{\text{move}}$ を得る。これらを結合し、市況データをエンコードした隠れ状態 \mathbf{m} を得る：

$$\mathbf{m} = \mathbf{W}_m \left([\mathbf{h}_{\text{short}}^{\text{std}} ; \mathbf{h}_{\text{long}}^{\text{std}} ; \mathbf{h}_{\text{short}}^{\text{move}} ; \mathbf{h}_{\text{long}}^{\text{move}}] \right) + \mathbf{b}_m. \quad (4)$$

デコーダは LSTM を用い、デコーダの初期値 \mathbf{s}_0 として、前述した隠れ状態 \mathbf{m} を用いる。時刻 t におけるデコーダの隠れ層の状態 \mathbf{s}_t は、ニュースヘッドラインの配信時刻 (9 時や 13 時などの 1 時間刻みの数値) を埋め込んだ時間帯情報埋め込みベクトル \mathbf{T} と、直前の単語埋め込み \mathbf{w}_{t-1} 、直前の隠れ層の状態 \mathbf{s}_{t-1} を用いて次のように計算される：

$$\mathbf{s}_t = \text{LSTM}([\mathbf{T} ; \mathbf{w}_{t-1}], \mathbf{s}_{t-1}). \quad (5)$$

ここで、時間帯情報埋め込みベクトル \mathbf{T} を各時点の隠れ状態に追加入力することで、生成文の制御を行う。

3 提案手法

Aoki et al. (2018) のモデルをベースとし、参照時刻の不整合解消のために、概況テキストが発表された時刻から取得できる直近のデータだけでなく、それ以前の時刻で取得できるデータを合わせて入力とする。図 2 にその概要を示す。具体的には、短期的な変動を捉えるために用意した 5 分足で収集した市況データを、直近 N 個の数値だけでなく、 k ステップ前 ($5 \times k$ 分前) で取得できる N 個の市況データ $\mathbf{x}_{\text{short-}k\text{step}}$ =

$(x_{\text{short-}k\text{step}, 1}, x_{\text{short-}k\text{step}, 2}, \dots, x_{\text{short-}k\text{step}, N})$ を用いる。これにより、市況データの直近の動向とそれ以前のステップとの動向差異や、概況テキストが発表される時間帯などを基に、参照すべき時系列データが判断できると考えられる。入力を 5 分足で収集した市況データの直近 N 個の数値 $\mathbf{X}_{\text{short-latest}}$ 、1 ステップ前から n ステップ前に取得できる各 N 個の数値 $\mathbf{X}_{\text{short-1step}}, \dots, \mathbf{X}_{\text{short-nstep}}$ 、および長期的な変動を捉えるための直近 M 取引日分の終値データ \mathbf{X}_{long} とする。Aoki et al. (2018) と同様に、各入力に前処理をおこない、 $\mathbf{X}_{\text{short-latest}}^{\text{std}}$, $\mathbf{X}_{\text{short-1step}}^{\text{std}}, \dots$, $\mathbf{X}_{\text{short-nstep}}^{\text{std}}$ および $\mathbf{X}_{\text{long}}^{\text{std}}$ と、 $\mathbf{X}_{\text{short-latest}}^{\text{move}}$, $\mathbf{X}_{\text{short-1step}}^{\text{move}}, \dots$, $\mathbf{X}_{\text{short-nstep}}^{\text{move}}$ および $\mathbf{X}_{\text{long}}^{\text{move}}$ を得る。これらのベクトルにそれぞれ MLP を用いて、 $\mathbf{h}_{\text{short-latest}}^{\text{std}}$, $\mathbf{h}_{\text{short-1step}}^{\text{std}}, \dots$, $\mathbf{h}_{\text{short-nstep}}^{\text{std}}$ および $\mathbf{h}_{\text{long}}^{\text{std}}$ と、 $\mathbf{h}_{\text{short-latest}}^{\text{move}}$, $\mathbf{h}_{\text{short-1step}}^{\text{move}}, \dots$, $\mathbf{h}_{\text{short-nstep}}^{\text{move}}$ および $\mathbf{h}_{\text{long}}^{\text{move}}$ を得、これらを結合し、市況データをエンコードした隠れ状態 \mathbf{m} を得る：

$$\mathbf{m} = \mathbf{W}_m \left([\mathbf{h}_{\text{short-latest}}^{\text{std}} ; \mathbf{h}_{\text{short-1step}}^{\text{std}} ; \dots ; \mathbf{h}_{\text{short-nstep}}^{\text{std}} ; \mathbf{h}_{\text{long}}^{\text{std}} ; \mathbf{h}_{\text{short-latest}}^{\text{move}} ; \mathbf{h}_{\text{short-1step}}^{\text{move}} ; \dots ; \mathbf{h}_{\text{short-nstep}}^{\text{move}} ; \mathbf{h}_{\text{long}}^{\text{move}}] \right) + \mathbf{b}_m. \quad (6)$$

これをデコーダの初期値 \mathbf{s}_0 として用いる。

4 実験

4.1 データ

実験には、時系列データとして、Thomson Reuters DataScope Select*¹ から、2010 年 12 月から 2016 年 9 月までの期間の日経平均株価指数やダウ平均株価を含む 10 種類*² の時系列データを収集し、実験に使用した。概況テキストとして、日経 QUICK ニュース社が提供している日経平均株価に言及しているニュース記事のヘッドラインを使用した。2010 年 12 月から 2015 年 9 月の期間のデータを学習データ、2015 年 10 月から 2016 年 3 月の期間のデータを開発データ、2016 年 4 月から 2016 年 9 月の期間のデータを評価データとして使用した。実験に使用したデータについて、概況テキストに日経平均株価の数値の動向を説明する表現*³ が含まれているか、含まれている場合には、その表現と取得したデータの数値の動向が一致しているかを調査した。その統計値を表 1 に示す。

4.2 実験設定

5 分足で収集した記事発表時刻直近の 62 個の市況データ $\mathbf{X}_{\text{short-latest}}$ 、直近 7 取引日分の終値データ \mathbf{X}_{long} 、

*¹ <https://hosted.datascope.reuters.com/DataScope/>

*² 日経平均株価、東証株価指数、ダウ平均株価、S&P500、上海総合指数、香港ハンセン株価指数、FTSE100 種総合株価指数、USD/JPY、EUR/JPY、日経平均先物

*³ 上方向の動向：「続伸、反発、上げ」、下方向の動向：「続落、反落、下げ」

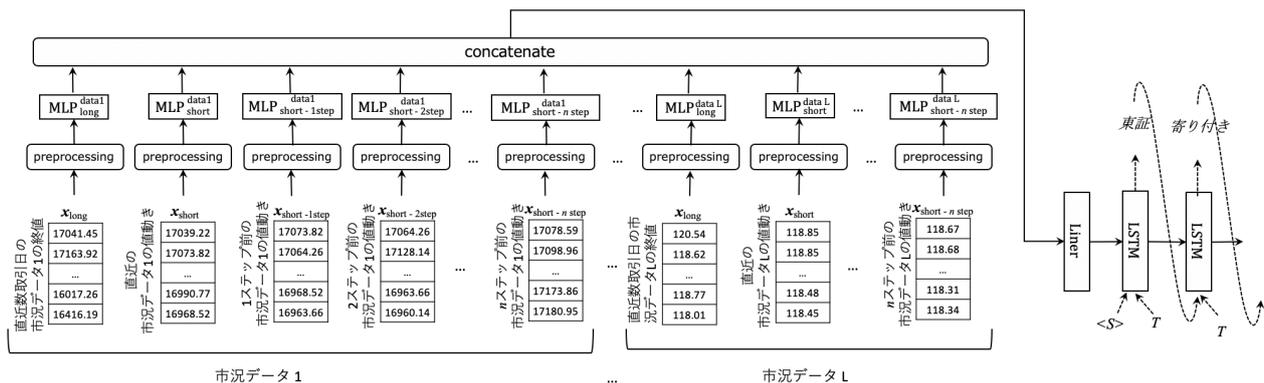


図 2: 提案モデルの全体概要図.

表 1: 使用したデータの統計値

期間	動向表現			合計
	なし	あり		
		一致	不一致	
学習データ	3,990	11,891	396	16,277
開発データ	370	1,450	46	1,866

$$\begin{aligned}
 move_{text}(w) &= \begin{cases} \text{上方向} & (w \in \{\text{続伸, 反発, 上げ}\}) \\ \text{下方向} & (w \in \{\text{続落, 反落, 下げ}\}) \\ \text{なし} & (\text{otherwise}) \end{cases} \\
 move(x) &= \begin{cases} \text{上方向} & (x \text{ の直近の数値動向} > 0) \\ \text{下方向} & (x \text{ の直近の数値動向} < 0) \\ \text{なし} & (\text{otherwise}) \end{cases}
 \end{aligned}$$

および概況テキストのペアを使用する。また記事発表時刻以前のデータとして、1 ステップ前 $X_{\text{short-1step}}$ から 6 ステップ前 $X_{\text{short-6step}}$ までを段階的に用いる。なお、時系列株価データからテキスト中で言及された価格を適切に出力するために、Murakami et al. (2017), Aoki et al. (2018) と同様に、概況テキスト中の数値表現を汎化タグに置換した。株価等の時系列数値データをベクトルへ変換する各 MLP の隠れ状態の次元は 256、これらの隠れ状態を結合した後、線形変換により次元を 256 とした。デコーダの隠れ状態の次元は 256、単語埋め込みベクトルの次元は 128、時間帯情報埋め込みベクトルの次元は 64 とした。各パラメータの最適化手法には Adam (Kingma and Ba, 2015) を使用し、学習率は 5×10^{-4} とした。ミニバッチサイズは 100、エポック数は 60 とし、学習時に開発データで計算された BLEU が連続して下がった場合には学習を終了し、各エポックのモデルの中で、開発データに対する BLEU が最大となったモデルを使用して、評価を行う。

4.3 評価方法

評価指標には、実際の株価の概況テキストと生成されたテキストの一致度合いを測る目的として BLEU を使用する。また、正解概況テキストおよび生成した概況テキストの、日経平均株価の短期的・長期的な動向を説明する表現と、発表時刻直近の数値データの動向の一致・不一致を取得して評価する。開発データの i 番目の事例とモデルの生成文を合わせて $(x_i, w_i^{\text{gold}}, w_i^{\text{pred}})$ とし、次のように定義する:

$$\begin{aligned}
 S_{\text{gold}} &= \{i | move_{text}(w_i^{\text{gold}}) = move(x_i)\} \\
 S_{\text{pred}} &= \{i | move_{text}(w_i^{\text{pred}}) = move(x_i)\} \\
 D_{\text{gold}} &= \{i | move_{text}(w_i^{\text{gold}}) \neq move(x_i)\} \\
 D_{\text{pred}} &= \{i | move_{text}(w_i^{\text{pred}}) \neq move(x_i)\}
 \end{aligned}$$

これらを用いて、次の式により評価する:

$$\text{一致}_{\text{precision}} = \frac{|S_{\text{gold}} \cap S_{\text{pred}}|}{|S_{\text{pred}}|} \quad (7)$$

$$\text{一致}_{\text{recall}} = \frac{|S_{\text{gold}} \cap S_{\text{pred}}|}{|S_{\text{gold}}|} \quad (8)$$

$$\text{不一致}_{\text{precision}} = \frac{|D_{\text{gold}} \cap D_{\text{pred}}|}{|D_{\text{pred}}|} \quad (9)$$

$$\text{不一致}_{\text{recall}} = \frac{|D_{\text{gold}} \cap D_{\text{pred}}|}{|D_{\text{gold}}|} \quad (10)$$

BLEU による評価は開発データおよび評価データに対して行い、概況テキストの動向表現と数値データとの比較による評価は開発データに対して行った。なお、学習時の seed を 0, 5, 10, 50, 100, 500 の 6 種類とし、各モデルによる評価結果の平均を 4.4 節の実験結果として用いた。

4.4 実験結果

表 2 の実験結果によると、発表時刻直前に取得できる数値データのみを用いる既存手法に比べて、それ以前の 1 から 6 ステップ前で取得できる数値データも用いる提案手法で、BLEU の向上がわずかに見られる。また、提案手法により、概況テキストの動向表現と数値データの動向が一致する場合の precision / recall が、既存手法に比べてわずかに下がる場合 (表 2 下線) があるものの、概況テキストの動向表現と数値データの動

表 2: BLEU (%) と生成文中の動向表現と数値データの動向一致・不一致評価 (%)

手法	開発データ					評価データ BLEU
	BLEU	一致		不一致		
		precision	recall	precision	recall	
既存手法	25.2	98.1	96.0	14.2	5.2	24.0
提案手法:						
+ $X_{\text{short-1step}}$	25.7	98.2	96.0	27.6	8.3	25.3
+ $X_{\text{short-2step}}$	26.7	<u>98.0</u>	96.6	30.5	10.1	26.7
+ $X_{\text{short-3step}}$	26.6	98.1	<u>95.9</u>	23.8	10.4	26.7
+ $X_{\text{short-4step}}$	26.8	98.2	96.0	22.1	9.4	26.8
+ $X_{\text{short-5step}}$	27.1	98.1	96.3	23.4	9.0	27.2
+ $X_{\text{short-6step}}$	25.8	<u>98.0</u>	96.2	18.0	7.6	26.2

表 3: 生成された文 1 (記事発表時の数値動向=23.15)

手法	生成文	BLEU	動向との一致
正解文	日経平均、下げに転じる	—	×
既存手法	東証後場寄り、小動き アジア株高が支え、先物に買い	1.6	(記述なし)
提案手法:			
+ $X_{\text{short-1step}}$	東証後場寄り、上げ幅拡大 アジア株高が支え	2.0	○
+ $X_{\text{short-2step}}$	日経平均、下げに転じる	100.0	×
+ $X_{\text{short-3step}}$	日経平均、下げに転じる	100.0	×
+ $X_{\text{short-4step}}$	日経平均、下げに転じる	100.0	×
+ $X_{\text{short-5step}}$	日経平均、下げに転じる	100.0	×
+ $X_{\text{short-6step}}$	日経平均、下げに転じる	100.0	×

向が不一致の場合には、全ての場合で precision / recall が共に向上した。特に、 $X_{\text{short-2step}}$ までの数値を用いた場合に precision が、 $X_{\text{short-3step}}$ までの数値を用いた場合に recall が、最も高くなった (表 2 太字)。

各モデルでの生成結果を表 3 および表 4 に示す。表 3 の生成例では、 $X_{\text{short-2step}}$ よりも以前の数値データを合わせて用いることで、参照時刻の不整合を解決できていることがわかる。一方、表 4 の生成例では、 $X_{\text{short-2step}}$ まで、もしくは $X_{\text{short-5step}}$ までの数値を用いることで、参照時刻の不整合を解決できている。

これらの実験結果より、参照時刻の不整合解消には、直近の数値だけでなく、以前の時刻で取得できるデータを合わせて使用することが効果的と言えるが、 $X_{\text{short-6step}}$ までを入力とした場合に BLEU, precision / recall とも最大値となっていないことより、より大きい時間幅を持たせて入力とすれば、より良く不整合を解決できるとは限らないと言える。

5 結論

本研究では、時系列データを用いた動向説明文生成において、直近の数値データだけでなく、それ以前の複数のタイムステップで取得できる数値データをあわせて考慮することで、参照時刻の不整合を解消して概況を生成するモデルを提案した。6 ステップ前までに取

表 4: 生成された文 2 (記事発表時の数値動向=-5.49)

手法	生成文	BLEU	動向との一致
正解文	日経平均、小幅続伸で始まる 始値は 14 円高の 18961 円	—	×
既存手法	日経平均、小動きで始まる 米株高や円安を好感	8.7	(記述なし)
提案手法:			
+ $X_{\text{short-1step}}$	日経平均、反落で始まる 米株安や円高で利益確定売	8.6	○
+ $X_{\text{short-2step}}$	日経平均、小幅続伸で始まる 米株高や円安で買い先行	41.4	×
+ $X_{\text{short-3step}}$	日経平均、反落で始まる 下げ幅 0 円超	8.5	○
+ $X_{\text{short-4step}}$	日経平均、反落で始まる 米株安円高で売り先行	8.7	○
+ $X_{\text{short-5step}}$	日経平均、小幅続伸で始まる 米株高や円安を好感	41.6	×
+ $X_{\text{short-6step}}$	日経平均、反落で始まる 米株安や円高で	8.6	○

得できる数値データを段階的に入力とし、生成モデルを作成、評価したところ、直近の数値データだけでなく、以前のタイムステップによる数値データを考慮することで、概況生成の精度が向上することが示された。

謝辞 この成果は、国立研究開発法人新エネルギー・産業技術総合開発機構 (NEDO) の委託業務の結果得られたものです。

参考文献

- Aoki, Tatsuya, Akira Miyazawa, Tatsuya Ishigaki, Keiichi Goshima, Kasumi Aoki, Ichiro Kobayashi, Hiroya Takamura, and Yusuke Miyao (2018) “Generating Market Comments Referring to External Resources,” in *Proc. of INLG 2018*, pp. 135–139.
- Iso, Hayate, Yui Uehara, Tatsuya Ishigaki, Hiroshi Noji, Eiji Aramaki, Ichiro Kobayashi, Yusuke Miyao, Naoaki Okazaki, and Hiroya Takamura (2019) “Learning to Select, Track, and Generate for Data-to-Text,” in *Proc. of ACL 2019*, pp. 2102–2113.
- Kingma, Diederik P. and Jimmy Ba (2015) “Adam: A Method for Stochastic Optimization,” in *Proc. of ICLR 2015*.
- Liu, Tianyu, Kexiang Wang, Lei Sha, Baobao Chang, and Zhifang Sui (2018) “Table-to-text generation by structure-aware seq2seq learning,” in *Proc. of AAAI*.
- Murakami, Soichiro, Akihiko Watanabe, Akira Miyazawa, Keiichi Goshima, Toshihiko Yanase, Hiroya Takamura, and Yusuke Miyao (2017) “Learning to Generate Market Comments from Stock Prices,” in *Proc. of ACL 2017*, pp. 1374–1384.
- Puzikov, Yevgeniy and Iryna Gurevych (2018) “E2E NLG Challenge: Neural Models vs. Templates,” in *Proc. of INLG 2018*, pp. 463–471.
- Sutskever, I, O Vinyals, and QV Le (2014) “Sequence to sequence learning with neural networks,” in *Proc. of NIPS*.
- Taniguchi, Yasufumi, Yukun Feng, Hiroya Takamura, and Manabu Okumura (2019) “Generating Live Soccer-Match Commentary from Play Data,” in *Proc. of AAAI*.
- 漆原理乃・小林一郎 (2018) 「人の動作および物体認識に基づく動画画像からの文生成」、『言語処理学会第 24 回年次大会 (2018)』, 160–163 頁。
- 村上総一郎・笹野遼平・高村大也・奥村学 (2017) 「数値予報マップからの天気予報コメントの自動生成」、『言語処理学会第 23 回年次大会 (2017)』, 1121–1124 頁。