言語から画像を生成する深層学習モデルの挙動に関する考察

A Study on Behavior of Deep Neural Text-to-Image Generative Model

藤山 千紘^{*1} 小林 一郎^{*2} Chihiro Fujiyama Ichiro Kobayashi

*1お茶の水女子大学人間文化創成科学研究科理学専攻情報科学コース Advanced Sciences, Graduate School of Humanities and Sciences, Ochanomizu University

*2お茶の水女子大学 基幹研究院 自然科学系

Advanced Sciences, Graduate School of Humanities and Sciences, Ochanomizu University

In this study, we analyze the behavior of the computational mechanism and the structure of the feature representation space in a deep neural text-to-image generative model. This is a fundamental approach with a goal to construct artificial general intelligence reflecting the mechanism of human intelligence. First, we explore whether the model is capable of encoding captions and of generating valid images under the circumstance given input captions without word boundaries. Qualitative and quantitative evaluations demonstrate that it can generate compelling images, but the computational mechanism does not acquire the units of meaning. Secondly, we analyze the semantic compositionality in the embedding space. Our experimental result suggests that the semantic compositionality appears between words indicating positions.

1. はじめに

近年,汎用人工知能の構築を目指して,ヒトの知能に関する 知見を取り入れた人工神経回路網の研究が盛んに行われている [Cox 14, Lotter 17].また,機械学習モジュールが自然言語の 意味を理解することを,自然言語から画像へのグラウンディン グができることとして捉え,自然言語の意味理解,ひいては自 然言語からの概念の獲得を目指して,自然言語で記述された キャプションを入力とする画像生成モデルの提案が数多くなさ れている [Mansimov 16, Xu 18, Zhang 18].一方で,深層学 習分野での研究成果は経験的な成果として示されることが多 く,モデルの計算機構の挙動や特徴表現空間の構造についての 分析はあまり行われていない.

本研究では、ヒトの知能のメカニズムを模倣して構築され た、キャプションを入力とする画像生成モデルを対象に、入力 の粒度を変更した際の計算機構の挙動や、特徴表現空間の構造 を、ヒトの知能との親和性の観点から分析する.

2. alignDRAW

Mansimovら [Mansimov 16] は、ヒトが絵を描く際の,「特定 の言語表現に着目してそれに対応する部分を描く」というプロセ スの反復を、深層学習の枠組みで実現することを意図して構築し た画像生成モデル alignDRAW を提案している.このモデルは、 Variational AutoEncoder[Kingma 14a] を拡張し、生成画像 の段階的な高精細化および空間的注意を導入した画像生成モデ ル Deep Recurrent Attentive Writer(DRAW)[Gregor 15] を、 キャプションで条件付けて画像生成できるようにした Encoder-Decoder モデルである.alignDRAW は、単語単位のキャプショ ンを入力にとり、双方向 LSTM[Hochreiter 97] で構成される 言語エンコーダで処理を行った後、DRAW を画像デコーダと して、attention mechanism[Bahdanau 15] と合わせて用いる ことにより、反復的に画像生成を行う.alignDRAW の概要図 を図 1 に示す.

連絡先: 藤山千紘, お茶の水女子大学, 〒 112-8610 東京都文 京区大塚 2-1-1, fujiyama.chihiro@is.ocha.ac.jp



図 1: alignDRAW 概要図

3. 提案手法

3.1 alignDRAW を用いた単語分割タスクを含む画 像生成

先行研究 [Mansimov 16] では、単語分割済みのキャプショ ンから画像生成を行っているが、本手法では入力を単語分割 されていないキャプションに変更し、単語の境界情報が失われ た、すなわち意味の単位の情報が欠落した場合の alignDRAW の言語エンコード能力および画像生成能力を評価する.具体的 には、単語分割されていないキャプションに対して、妥当な画 像を生成し得るか、またその際の attention mechanism の挙 動が言語の意味の単位を表現しているかを考察する.

3.2 alignDRAW における言語の意味の構成的特性 の分析

本手法では、モデルの特徴表現空間において言語の意味の構成的特性が表現されるかを評価する.先行研究 [Mansimov 16] では、キャプションに含まれる各単語を one-hot ベクトルとして双方向 LSTM に入力しているが、本手法では各単語を分散表現に埋め込んだ後、双方向 LSTM に入力するよう、埋め込み層を追加する形でモデルの拡張を行う.モデルは単語分割済みのキャプションを用いて学習し、埋め込み層の分散表現について、単語の意味の構成的特性を分析する.

4. 実験

4.1 実験設定

データセットは、表1に示すテンプレートと手書き数字画 像のデータセット MNIST^{*1}を用いて作成した.キャプション はプレースホルダーを含むテンプレートを各実験8種類用意 し、MNIST から無作為に抽出した画像および正解ラベルの組 のうちラベル情報を埋め込む形で作成した.画像は、ラベル と対応する MNIST 画像をキャプション内容に適合する領域 に、無作為に4ピクセルのゆらぎを持たせて配置した60×60 ピクセルのグレースケールの画像とした.両実験ともに、学習 データ40,000 事例、開発データ4,000 事例、評価データ4,000 事例を用いた.モデルの実装には深層学習のフレームワーク TensorFlow^{*2}を用い、学習は各実験、表2の設定で行った.

表 1: キャプション作成時のテンプレート

単語分割タスクを含む画像生成	構成的特性の分析
すうじ_がすうじ_のひだりにある.	数字 _ が 画像 の 左 に ある.
すうじ_がすうじ_のみぎにある.	数字 _ が 画像 の 右 に ある .
すうじがすうじのうえにある.	数字 _ が 画像 の 上 に ある.
すうじがすうじのしたにある.	数字 _ が 画像 の 下 に ある .
すうじががぞうのひだりうえにある.	数字 _ が 画像 の 左上 に ある .
すうじ_ががぞうのひだりしたにある.	数字 が 画像 の 左下 に ある .
すうじががぞうのみぎしたにある.	数字 _ が 画像 の 右下 に ある .
すうじががぞうのみぎうえにある.	数字 が 画像 の 右上 に ある .

表 2: alignDRAW 学習時のハイパーパラメータ

	単語分割タスクを含む画像生成	構成的特性の分析	
入力語彙サイズ	33	28	
言語	one-hot ベクトル \rightarrow	32 次元分散表現 →	
エンコーダ	128 ユニット双方向 LSTM	128 ユニット双方向 LSTM	
attention mechanism	512 ユニット	256 ユニット	
	Bahdanau Attention	Bahdanau Attention	
デコーダ	300 ユニット DRAW LSTM		
描画反復回数	32 ステップ		
潜在変数 z	150 次元		
最適化アルゴリズム	RMSProp [Tie	rop [Tieleman 12]	
	初期学習率: 0.001	初期学習率: 0.001	
学習率	110 エポック以降	75 エポック以降	
	0.0001 に減衰	15 エポック毎に 0.5 倍	
パラメータ初期値	平均: 0, 分散: 0.1 の正規分布乱数		
学習エポック数	200	150	

生成画像の定量評価としては、生成された画像を入力キャ プション毎に分類する分類器を用い、弁別性の自動評価を行っ た.分類器の学習は alignDRAW の学習に用いたデータセッ トと同一のデータセット上で行い、ハイパーパラメータは表3 に示す設定とした.なお、キャプションは単語分割タスクを含 む画像生成で440 種類、構成的特性の分析で80 種類であるの で、各々440 クラス分類、80 クラス分類となる.

表 3: 定量評価のための分類器学習時のハイパーパラメータ (両実験 とも共通)

アーキテクチャ	(CNN+最大値プーリング) × 4 層 + 全結合 × 1 層
	フィルタサイズ: 5 × 5 (全層)
畳み込みに	チャネル数: 入力層から順に 32, 64, 128, 256
関する設定	0 パディング
	ストライド: 1
全結合層	1024 ユニット
活性化関数	ReLU 関数
損失関数	交差エントロピー損失
勾配クリッピング	[-10.0, 10.0]
最適化アルゴリズム	Adam [Kingma 14b]
初期学習率	0.0001
パラメータ	平均: 0, 標準偏差: 0.1 の
初期値	±2σ の切断正規分布乱数
学習エポック数	100

^{*1} http://yann.lecun.com/exdb/mnist/

*2 https://www.tensorflow.org/

4.2 単語分割タスクを含む画像生成

単語分割されていないキャプションからの生成画像例を図2 に示す.単語分割済みのキャプションからの画像生成と比較し て、単語の境界情報が失われている、つまり意味の単位の情報 が欠落している点で、画像生成タスクとしてはより困難になっ ていると考えられるが、生成結果からキャプション内容に適合 する画像を生成できていることが確認できる.

キャプション	生成画像	参照画像
すうじぜろががぞうのみぎしたにある.	0	0
すうじなながすうじよんのしたにある.	4 7	4 7
すうじろくががぞうのひだりしたにある.	6	6

図 2: 単語分割されていないキャプションからの生成画像例

分類における評価データ上での正解率と alignDRAW による 生成画像データ上での正解率を表4に示す.生成画像データ上 での正解率が評価データ上での正解率を上回っており, align-DRAW によって生成された画像がどのキャプションに由来す るかを十分に識別できる質を達成できていることが分かる.こ れは,生成モデル学習の結果,データセットに存在するノイズ や手書き文字の個人差が吸収され,正規化された数字を生成で きるようになったためと考えられる.

表 4: 評価データおよび alignDRAW による生成画像データ上での分 類正解率

評価データ	生成画像データ
0.719	0.737

またテンプレートに従わないキャプションからの生成画像例 を図3に示す.テンプレートに含まれるキャプションからの生 成結果(図2)と比較して、描かれる数字の質が劣化している 事例が見受けられるが、おおよそキャプション内容に適合する 画像を生成できており、「すうじ」という言語表現の省略や主 語の位置変更に対して、モデルがある程度頑健に動作している ことが認められる.

キャプション	生成画像
すうじぜろのひだりにすうじきゅうがある.	90
さんががぞうのひだりしたにある.	3

図 3: テンプレートに従わないキャプションを入力とする生成画像例

続いて、画像を生成する際の attention mechanism の挙動 を図4に示す。図4上段については、「いち」という言語表現付 近に attention が当たっているときに画像空間上で1に相当す る部分の生成が進んでおり、続いて数字2を表す「に」という 言語表現付近に attention が移って画像空間での2に相当する 部分の鮮明化が進んでいる様子が見てとれる.2箇所のプレー スホルダーを持つテンプレートに由来するキャプションでは, 一方の数字のアイデンティティを示す言語表現およびその周囲 に attention が当たっているときにそれに対応する部分が画像 空間で鮮明化し, attention が他方の数字に移ると生成部位も 移るという傾向が共通して見られ,数字のアイデンティティを 特定する言語表現と画像空間での表現に大まかな対応がとれ ていることが確認された.一方で,この場合には空間を表す言 語表現についてはほとんど attention が当たっていなかった. また図 4 下段については,生成画像には3を描くことができ ているが,数字のアイデンティティを表現する「さん」にはほ とんど attention が当たっていない。この傾向はプレースホル ダーが1箇所のみのテンプレートに由来するキャプションに 共通して見られ,数字のアイデンティティを表す言語表現と画 像空間での表現に対応関係が認められなかった.



図 4: attention mechanism の挙動: 左に入力キャプション,上に画 像の生成過程 (左から t=0, 7, 15, 23, 31) を示す.

単語分割されていないキャプションから妥当な画像を生成で きたことから、attention mechanism の挙動が意味の単位を表 現していることが期待されたが、図4のように、本設定下で はそのような傾向は見られなかった.これは、言語エンコーダ として用いた双方向LSTMの表現能力が高いことが一因であ ると考えられるため、言語エンコーダを簡素化した追加実験を 行って分析を加える.本稿では系列処理を行わず、one-hot ベ クトルを直接 attention mechanism の入力として画像生成を 行った場合について述べる.なおこの場合、系列処理を行わな いため入力の各文字は集合として扱われ、並びの順序は無視さ れる.したがって、キャプション内の主体と対象の関係を画像 空間上に適切に反映させ得ない場合があることが想定される. キャプションの系列処理を行わない場合の生成画像例を図5 に示す.図5上段のように妥当な画像を生成できる場合もある 一方で,下段のように位置関係が反転する事例があることが確 認された.キャプションを系列処理しなかった場合の attention mechanism の挙動を図6に示す.系列処理を行わない,つま りキャプションを文字の集合として扱った場合には,数字のア イデンティティと描画位置を一意に特定するために必要な言語 表現に attention が集中する傾向があったが,これは必ずしも 単語単位での表現とは一致しないため,attention mechanism の挙動が意味の単位を表現している様子は見られなかった.一 方で,双方向 LSTM を用いた場合と比較すると,数字のアイ デンティティや空間を表す言語表現の一部に確実に attention が当たっていることが確認できた.



図 5: キャプションを系列処理しない場合の生成画像例



図 6: キャプションを系列処理しない場合の attention mechanism の 挙動: 左に入力キャプション,上に画像の生成過程 (左から t=0, 7, 15, 23, 31) を示す.

alignDRAW では画像生成過程に数字を描く順序の制約がな く自由度が高いため、このことに起因して attention の学習が 困難である.また現在のモデルには、言語表現の連続したセ グメントが意味の単位となり得るという情報が一切与えられ ていないため、attention mechanism の挙動で言語の意味の 単位を表現することが困難であったと考えられる.attention mechanism で意味の単位を捉えつつ、ヒトが絵を描く過程を 計算機構に反映して画像生成を行うためには、損失関数の検討 を含めてモデルの拡張が必要であると考えられ、これは今後の 課題の一つである.

4.3 言語の意味の構成的特性の分析

学習したモデルを用いた生成画像の例を図7に示す.本実験においても、キャプション内容に適合する画像が生成されていることが確認できる.

分類における評価データ上での正解率と alignDRAW によ る生成画像データ上での正解率を表5に示す.本実験では,生 成画像データ上での正解率が評価データ上での正解率を下回る 結果となり,学習した生成モデルがデータセットの分布を同質 といえる程には近似できていなかったことが分かる.図8に



図 7: 構成的特性の分析を目的とした実験での生成画像例

正しく分類できなかった事例を示す.間違って分類された生成 画像は、入力キャプションに含まれる数字のアイデンティティ を正しく画像内に表現できていない事例が多かった.本分析で は空間を意味する言語表現に着目するため、数字が描かれる 位置のみを分類する、すなわち数字のアイデンティティの区別 を分類の対象としない 8 クラス分類を表 3 の設定で学習した 結果、評価データ上、生成画像データ上ともに分類正解率が 1.00 となり、数字を描く位置を間違えた事例はなかったと推 測される.

表 5: 評価データおよび alignDRAW による生成画像データ上での分 類正解率

評価データ	生成画像データ
0.985	0.882



図 8: 正しく分類できなかった生成画像の例: 左から生成モデルへの 入力キャプション,参照画像,生成画像,分類器によって予測された キャプションを示す.

言語の意味の構成的特性の分析では、キャプションに含まれ る空間を意味する単語「左」「右」「上」「下」「左上」「左下」 「右下」「右上」の8単語を対象とした.我々は例えば「左上」 を,意味的に「左」と「上」を足したものとして解釈している と考えられる.この加法構成性を、埋め込み層の分散表現にお いて獲得できるかを評価するため、まず学習された分散表現の うち「左」「右」「上」「下」に対応するものの和をとり、各々 「左上」「左下」「右下」「右上」の推定分散表現を作成した.例 えば、「左上」の推定分散表現は、「左」と「上」それぞれに対 応する実分散表現の和として構成した. このようにして推定し た分散表現に対して、学習で得た実分散表現それぞれとの cos 類似度を算出した結果を表6に示す.また,各推定分散表現と cos 類似度が高かった分散表現に対応する単語上位5件を表7 に示す.本設定下では、「右上」以外の3単語については高い cos 類似度が得られ、表7においても推定分散表現と実分散表 現が比較的近くに確認されることから、埋め込み空間において 言語の意味の構成的特性が表現される可能性を示唆する結果を 得たと言える.

表 6:	推定分散表現と実分散表現の	\cos	類似度
------	---------------	--------	-----

	左上	左下	右下	右上
cos 類似度	0.89	0.80	0.92	0.29

表 7: 推定分散表現との cos 類似度が高い分散表現上位 5 件

	左上	左下	右下	右上
1	左上	左	右	右
2	上	左上	1	1
3	左	左下	右下	右下
4	0	0	の	の
5	2	2	7	4

5. おわりに

本研究では、キャプションからの画像生成モデルについて、 計算機構の挙動や特徴表現空間の構造の分析を行った.入力の キャプションをヒトが言語を獲得する過程を一段階遡る形で、 単語分割されたキャプションから単語の境界情報の欠落した キャプションに変更した場合には、定性的にも定量的にもキャ プション内容に適合する画像を生成するモデルを学習できた. しかし、その際の attention mechanism の挙動において意味 の単位を獲得している様子は確認されなかった.また、特徴表 現空間における言語の意味の構成的特性の分析については、埋 め込み空間の分散表現について空間を意味する単語間で、意味 の加法構成性を得られる可能性を示唆する結果を得た.

今後の課題としては、より精緻な分析を行うことや、ヒトの 知能のメカニズムを反映して動作する生成モデルの構築が挙げ られる.

参考文献

- [Bahdanau 15] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate.", In ICLR, 2015.
- [Cox 14] D. D., Cox, and T., Dean, "Neural networks and neuroscience-inspired computer vision.", Current Biology, 24(18), R921-R929, 2014.
- [Gregor 15] K. Gregor, I. Danihelka, A. Graves, D. J. Rezende, and D. Wierstra, "DRAW: A recurrent neural network for image generation.", In ICML, 2015.
- [Hochreiter 97] S. Hochreiter, and J. Schmidhuber, "Long shortterm memory.", Neural Computation, 9(8), pp. 1735-1780, 1997.
- [Kingma 14a] D. P., Kingma, and M., Welling, "Auto-encoding variational bayes.", In ICLR, 2014.
- [Kingma 14b] D. P., Kingma, and J., Ba, "Adam: A method for stochastic optimization.", arXiv preprint arXiv:1412.6980, 2014.
- [Lotter 17] W., Lotter, G., Kreiman, and D., Cox, "Deep predictive coding networks for video prediction and unsupervised learning.", In ICLR, 2017.
- [Mansimov 16] E. Mansimov, E. Parisotto, J. L. Ba., and R. Salakhutdinov, "Generating images from captions with attention.", In ICLR, 2016.
- [Tieleman 12] T., Tieleman, and G., Hinton, Lecture 6.5-RMSProp, COURSERA: Neural Networks for Machine Learning. University of Toronto, Tech. Rep, 2012.
- [Xu 18] T., Xu, P., Zhang, Q., Huang, H., Zhang, Z., Gan, X., Huang, and X., He, "Attngan: Fine-grained text to image generation with attentional generative adversarial networks.", In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 1316-1324), 2018.
- [Zhang 18] Z., Zhang, Y., Xie, and L., Yang, "Photographic textto-image synthesis with a hierarchically-nested adversarial network,", In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.