

注意モデルの変遷と展開

浅川伸一(東京女子大学) asakawa@ieee.org

14/Jun/2020

1 エピグラフ

蕉門に千歳不易の句、一時流行の句と云ふ有り。
是を二つに分て教へ給へる、其元は一つ也。
不易を知らざれば基たちがたく、流行を知らざれば風新たならず

—去来抄

2 謝辞

本日、このような機会を与えていただきましたエクサウイザーズ、藤井亮宏様、遠藤太一郎様に感謝申し上げます。



credit: {‘Beatles’:https://youtu.be/dsxtImDVMig, ABBA:https://youtu.be/92cwKCU8Z5c, ‘トランസフォーマー’:https://youtu.be/cwpXeH90qfE), ‘ELMO’:https://giftsandwich.com/christmas-gifts-for-kids/playskool-friends-sesame-street-tickle-me-elmo/, ‘BERT’: https://sesamestreet.tumblr.com/post/5772176064/bert-found-his-socks-what-a-great-monday}

3 自己紹介



図1 師匠エルマンとUSCDにて

浅川伸一:博士(文学) 東京女子大学情報処理センター勤務。早稲田大学在学時はピアジェの発生論的認識論に心酔する。卒業後エルマンネットの考案者ジェフ・エルマンに師事、薰陶を受ける。以来人間の高次認知機能をシミュレートすることを通して知的であるとはどういうことかを考えていると思っていた。著書に「AI白書2019, 2018」(2019年, アスキー出版, 共著), 「深層学習教科書ディープラーニングG検定(ジェネラリスト)公式テキスト」(2018年, 翔泳社, 共著), 「Pythonで体験する深層学習」(コロナ社, 2016), 「ディープラーニング, ビッグデータ, 機械学習あるいはその心理学」(新曜社, 2015), 「ニューラルネットワークの数理的基礎」「脳損傷とニューラルネットワークモデル, 神経心理学への適用例」いずれも守一雄他編「コネクショニストモデルと心理学」(2001)北大路書房など

4 アウトライン

1. どこにでも現れる注意
2. BERT概説
3. 流行の句あり
4. 不易の句あり
5. まとめ

5 第1部

どこにでも現れる注意

6 多頭=自己注意 Multi-Head Self-Attention:MHSA

- 自然言語処理 NLP **Transformer** Vaswani et al. (2017); **BERT** Devlin, Chang, Lee, & Toutanova (2018); **RoBERTa** Liu et al. (2019); **distilBERT** Sanh, Debut, Chaumond, & Wolf (2020); and more ...
- 画像処理 Ramachandran et al. (2019); **A2-Net** Chen, Kalantidis, Li, Yan, & Feng (2018); **U-GAT-IT** Kim, Kim, Kang, & Lee (2019)
- 強化学習, メタ学習 **SNAIL** Mishra, Rohaninejad, Chen, & Abbeel (2018)
- 敵対生成ネットワーク **SAGAN** Zhang, Goodfellow, Metaxas, & Odena (2019)

7 多頭=自己注意 Multi-Head Self-Attention

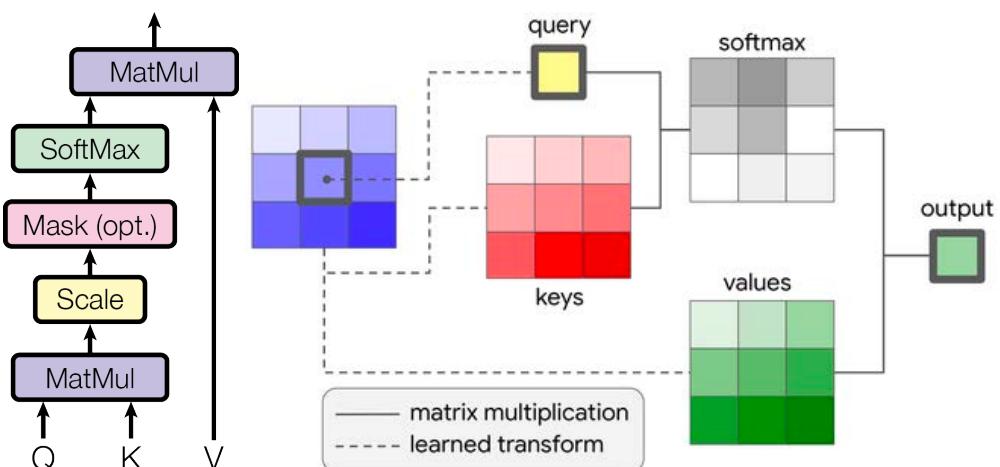


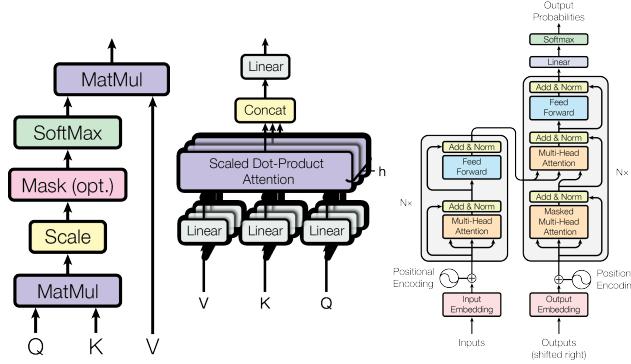
図2 Left: Vaswani et al. (2017), Right: Ramachandran et al. (2019)

$$\text{自己注意}(\mathbf{X}_{t,:}) = \text{ソフトマックス}(\mathbf{A}_{t,:}) \mathbf{X} \mathbf{W}_{\text{パリュ}_v}, \quad (1)$$

$$\mathbf{A} = \mathbf{X} \mathbf{W}_{\text{クエリ}_q} \mathbf{W}_{\text{キ}-k}^T \mathbf{X}^T \quad (2)$$

$$\mathbf{A} := (\mathbf{X} + \mathbf{P}) \mathbf{W}_{\text{クエリ}_q} \mathbf{W}_{\text{キ}-k}^T (\mathbf{X} + \mathbf{P})^T, \quad \mathbf{P} \text{ は 位置符号化器 PE} \quad (3)$$

8 Multi-head self-attention: MHSA(2)



9 Multi-head self-attention: MHSA (3) SAGAN (Self-Attention GAN)

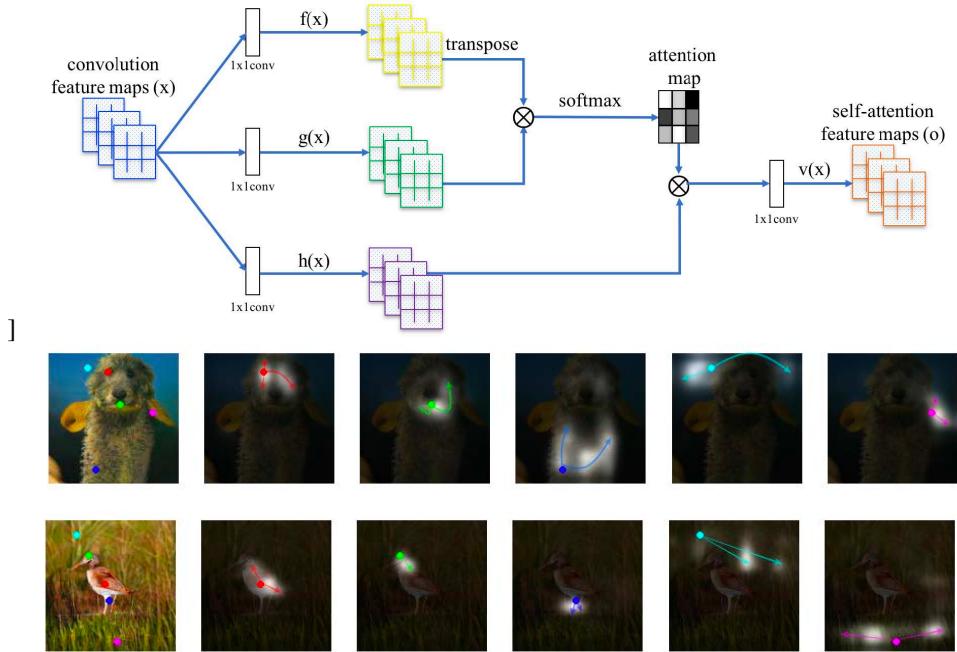


図3 From Zhang et al. (2019) Fig. 1, and 3.

画像生成において、近傍画素から情報だけでなく、関連する遠距離の特徴を利用して生成することにより一貫性のある対象やシナリオを生成可能。各行の左の元画像上のカラーポイントは5つの代表的なクエリの場所を示す。右側の5画像は各クエリ位置における注意地図。最も注目されている領域が、色分けされた矢印で示されている。

10 Multi-head self-attention: MHSA (4) Non-Local Net

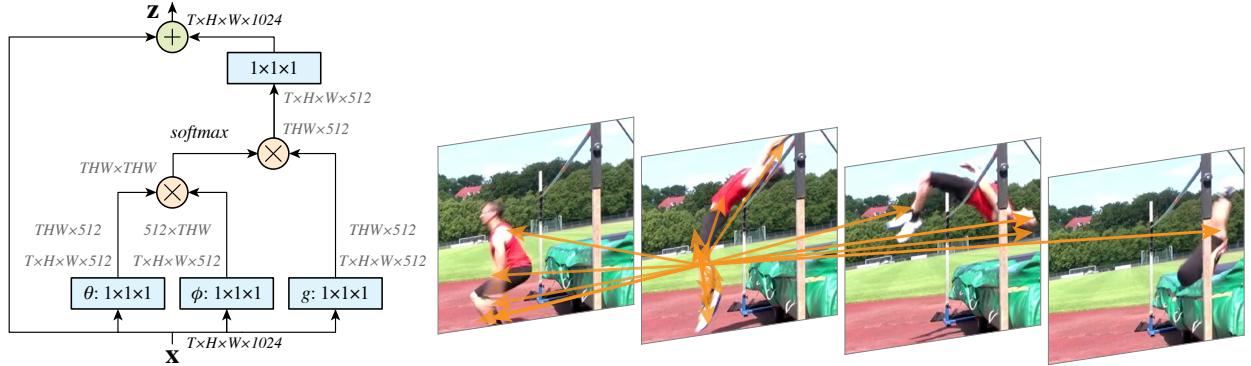


図4 時空の非局所ネットワークの概念図。特徴地図はテンソルとして示されている。例えば 1024 チャンネルの場合は $T \times H \times W \times 1024$ である。 \otimes は行列積を、 \oplus は要素和を示す。ソフトマックス演算は各行に対して実行される。青いボックスは $1 \times 1 \times 1 \times 1$ の畳み込みを表す。512 チャンネルのボトルネックを持つ埋め込みガウシアン版が示されている。バニラガウス版は θ と ϕ を除去することで ドット積版は $1/N$ のスケーリングでソフトマックスを置き換えることで行うことができる。From Wang et al. (2018)

11 Multi-head self-attention: MHSA (4) SNAIL

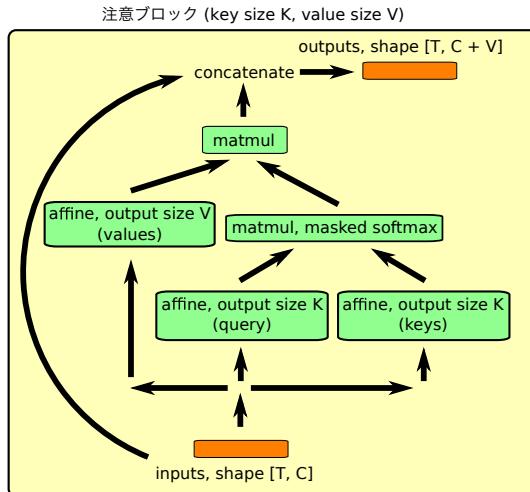


図5 From Mishra et al. (2018) Fig. 2

トランスフォーマーはリカレント構造や畳み込み構造を持たず埋め込みベクトルに位置符号化器を加えることで系列情報を処理する。しかし、逐次的な順序情報が貧弱であるとの批判がある。とりわけ強化学習のような位置依存性に敏感な課題では問題。トランスフォーマーモデルにおける位置問題を解決するため、自己注意機構と時間的な畳み込み temporal convolution を組み合わせたモデルが Simple Neural Attention Meta-Learner (SNAIL)Mishra et al. (2018)。SNAIL は、メタ学習、強化学習の両方の課題に優れていることが実証された。

12 注意用語集 Taxosonomy of attention

- 文脈ベース 注意 context-base attention: $\text{score}(s_t, h_i) = \cos(s_t, h_i)$ Graves, Wayne, & Danihelka (2014)
- 加算的 (連結的) 注意 Additive : $\text{score}(s_t, h_i) = v_a^\top \tanh(W_a [s_t; h_i])$ Bahdanau, Cho, & Bengio (2015)

- Luong, Pham, & Manning (2015) では 連結 concatenated, Vaswani et al. (2017) では 加算 additive と表記されている
- **場所ベース 注意 Location-Base:** $a_{t,i} = \text{softmax}(\mathbf{w}_a \mathbf{s}_t)$ Luong et al. (2015)
- Note: This simplifies the softmax alignment to only depend on the target position.
- **一般的 注意 general:** $\text{score}(s_t, h_i) = s_t^\top \mathbf{W}_a h_i$ Luong et al. (2015)
- \mathbf{W}_a は学習可能な結合係数行列
- **ドット積 注意 dot-product:** $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = s_t^\top h_i$ Luong et al. (2015)
- **スケール化ドット積 注意 scaled dot-product(?)**: $\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \frac{\mathbf{s}_t^\top \mathbf{h}_i}{\sqrt{n}}$ Vaswani et al. (2017)
 - スケール化規格化因子 $1/\sqrt{n}$ を用いる

13 第1部 Multi-head self-attention: MHSA のまとめ

- 自然言語処理, 画像処理, 強化学習, メタ学習の4分野でほぼ同様の MHSA が取り入れられている。
- クエリ, キー, バリュー 各テンソルを学習することが行われている
- 従来手法である 置き換える動きがある。
- ただし, SAGAN と SNAIL (non-local net) では 入力情報を concatenate して上位層に伝える点が他と異なる

14 補足 注意が現れるに至った歴史

- BOW, TFIDFJones (1972), SMTManning & Schütze (1999), N-gram モデル, Dimensionality would increase w.r.t. V^N
- RNN Elman (1990), Mikolov, Karafiat, Burget, Černocký, & Khudanpur (2010), Mikolov, Kombrink, Burget, Černocký, & Khudanpur (2011)
- LSTM Hochreiter & Schmidhuber (1997), Gers, Schmidhuber, & Cummins (1999), Greff, Srivastava, Koutník, Steunebrink, & Schmidhuber (2015), Seq2seqSutskever, Vinyals, & Le (2014), 注意モデル Bahdanau et al. (2015), Transformer Vaswani et al. (2017)
- BERT Devlin et al. (2018)

それぞれ有名なので説明はしません

15 第2部

BERT 概説

16 Mnih and Graves (2014)

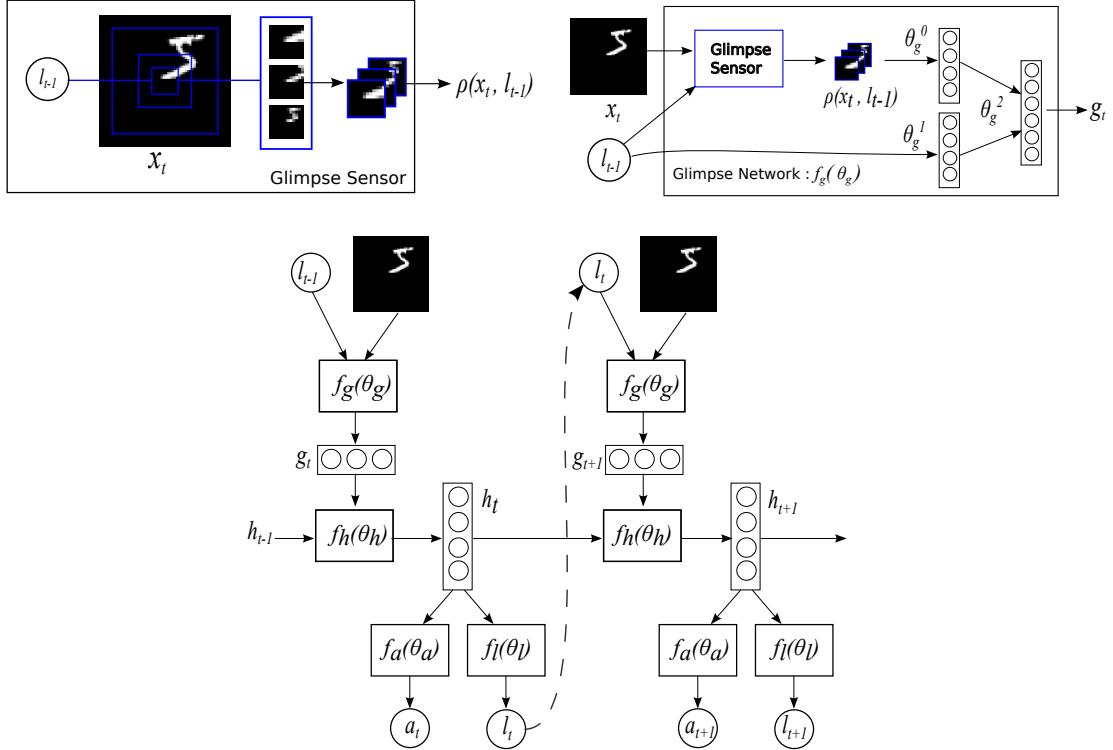


図 6 From Mnih et al. (2014)

17 Show and Tell (2014)

Figure 2. Attention over time. As the model generates each word, its attention changes to reflect the relevant parts of the image. “soft” (top row) vs “hard” (bottom row) attention. (Note that both models generated the same captions in this example.)

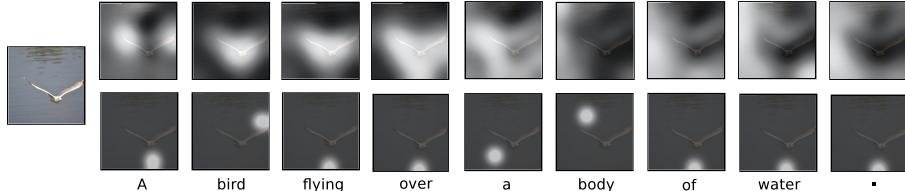


Figure 3. Examples of attending to the correct object (white indicates the attended regions, underlines indicated the corresponding word)

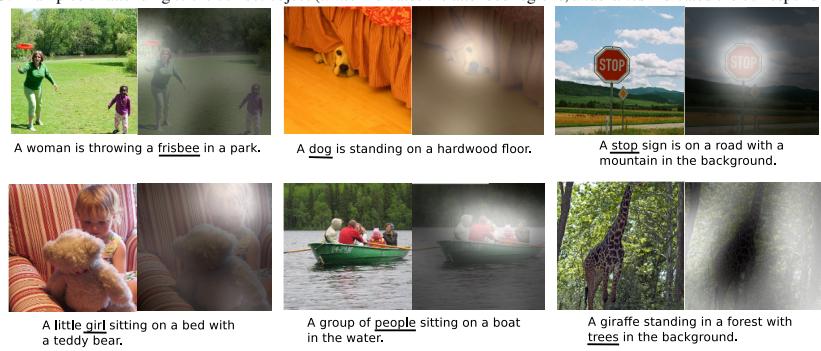


図 7 Attention for neural image captioning. From Xu et al. (2015)

18 Seq2seq model

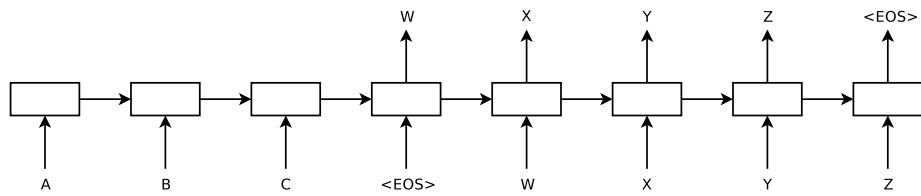


図 8 From Sutskever et al. (2014) Fig. 1, 翻訳モデル “seq2seq” の概念図

“eos”は文末を表す。中央の“eos”的前がソース言語であり、中央の“eos”的後はターゲット言語の言語モデルであるSRNの中間層への入力として用いる。

注意すべきは、ソース言語の文終了時の中間層状態のみをターゲット言語の最初の中間層の入力に用いることであり、それ以外の時刻ではソース言語とターゲット言語は関係がない。逆に言えば最終時刻の中間層状態がソース文の情報を全てを含んでいるとみなしうる。この点を改善することを目指すことが2014年以降盛んに行われてきた。顕著な例が後述する**双方向RNN**、**LSTM**採用したり、**注意**機構を導入することであった。

19 Seq2seq (2)

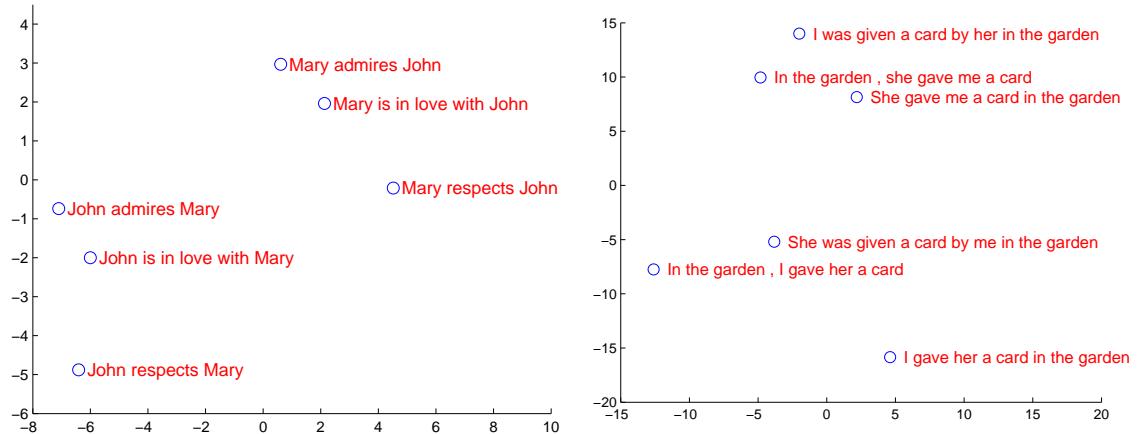


図 9 From Sutskever et al. (2014) Fig. 2

20 自然言語系の注意

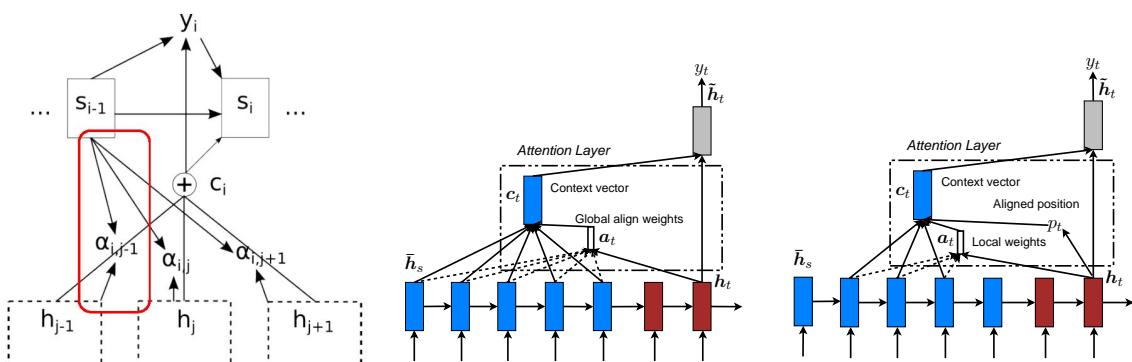


図 10 左:Bahdanau et al. (2015), 中:Luong et al. (2015) Fig. 2, 右:Luong et al. (2015) Fig. 3

21 BERT の特徴

BERT の特徴を 3 つにまとめると以下の通り

1. トランスフォーマー Transformer に基づく MHSA を用いた多層ニューラルネットワークモデル
2. 2 つの事前訓練: **マスク化言語モデル** と 次文予測課題
3. Fine tuning によるマルチタスクで性能向上 GLUE スコアボード, SuperGLUE を参照のこと

22 BERT の入力表現

埋め込みトークンの総和, 位置符号器, 分離埋め込みの 3 者 FromDevlin et al. (2018) Fig. 2

23 BERT の事前訓練: マスク化言語モデル

全入力系列のうち 15% をランダムに [MASK] トークンで置き換える

- 入力はオリジナル系列を [MASK] トークンで置き換えた系列
- ラベル: オリジナル系列の [MASK] 部分にの正しいラベルを予測
- 80%: オリジナル入力系列を [MASK] で置換
- 10%: [MASK] の位置の単語をランダムな無関連語で置き換える
- 10%: オリジナル系列

24 BERT の事前訓練: 次文予測課題

言語モデルの欠点を補完する目的, 次の文を予測

[SEP] トークンで区切られた 2 文入力

- 入力: the man went to the store [SEP] he bought a gallon of milk.
- ラベル: IsNext
- 入力: the man went to the store [SEP] penguins are flightless birds.
- ラベル: NotNext

25 BERT: ファインチューニング

(a), (b) は文レベル課題, (c),(d) はトークンレベル課題, E: 入力埋め込み表現, T_i : トークン i の文脈表象。

From Devlin et al. (2018) Fig.3

26 GLUE: General Language Understanding Evaluation

- **CoLA:** 入力文が英語として正しいか否かを判定
- **SST-2:** スタンフォード大による映画レビューの極性判断
- **MRPC:** マイクロソフトの言い換えコーパス。2 文が等しいか否かを判定
- **STS-B:** ニュースの見出し文の類似度を 5 段階で評定
- **QQP:** 2 つの質問文の意味が等価かを判定
- **MNLI:** 2 入力文が意味的に含意, 矛盾, 中立を判定
- **QNLI:** Q and A
- **RTE:** MNLI に似た 2 つの入力文の含意を判定

- **WNI:** ウィノグラッド会話チャレンジ

その他

- **SQuAD:** スタンフォード大による Q and A ウィキペディアから抽出した文
- **RACE:** 中学入試、高校入試に相当するテスト多肢選択回答 #BERT モデルの詳細
- データ: Wikipedia (2.5B words) + BookCorpus (800M words)
- バッチサイズ: 131,072 words (1024 sequences * 128 length or 256 sequences * 512 length)
- 訓練時間: 1M steps (~40 epochs)
- 最適化アルゴリズム: AdamW, 1e-4 learning rate, linear decay
- BERT-Base: 12 層, 各層 768 ニューロン, 12 多頭注意
- BERT-Large: 24 層, 各層 1024 ニューロン, 16 多頭注意
- 4x4 / 8x8 TPU で 4 日間

27 BERT: ファインチューニング手続きによる性能比較

Masking Rates			Dev Set Results		
MASK	SAME	RND	MNLI		NER
			Fine-tune	Fine-tune	Feature-based
80%	10%	10%	84.2	95.4	94.9
100%	0%	0%	84.3	94.9	94.0
80%	0%	20%	84.1	95.2	94.6
80%	20%	0%	84.4	95.2	94.7
0%	20%	80%	83.7	94.8	94.6
0%	0%	100%	83.6	94.9	94.6

図 11 マスク化言語モデルのマスク化割合の違いによる性能比較

マスク化言語モデルのマスク化割合はマスクトークン:ランダム置換:オリジナル=80:10:10 だけでなく、他の割合で訓練した場合の 2 種類下流課題、MNLI と NER で変化するかを下図 11 に示した。80:10:10 の性能が最も高いが大きな違いがあるわけではないようである。

28 BERT: モデルサイズ比較

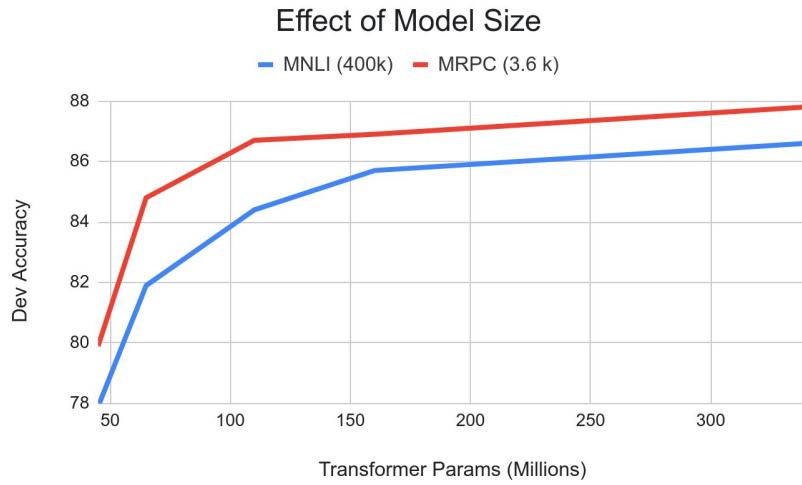


図 12 モデルのパラメータ数による性能比較

パラメータ数を増加させて大きなモデルにすれば精度向上が期待できる。下図では、横軸にパラメータ数で MNLI は青と MRPC は赤で描かれている。パラメータ数增加に伴い精度向上が認められる。図に描かれた範囲では精度が天井に達している訳ではない。パラメータ数が増加すれば精度は向上していると認められる。

29 BERT: モデル単方向、双方向モデル比較

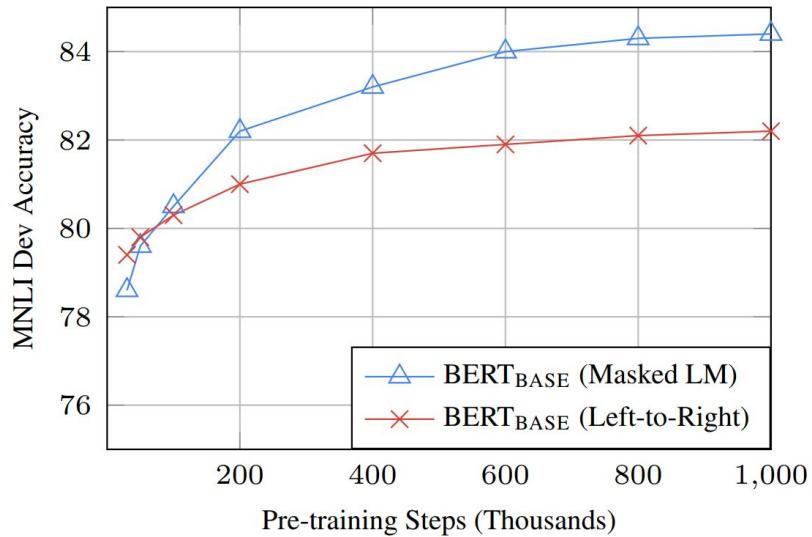


図 13 言語モデルの相違による性能比較

言語モデルをマスク化言語モデルか次単語予測の従来型の言語モデルによるかの相違による性能比較を下図 13 に示した。横軸には訓練ステップである。訓練が進むことでマスク化言語モデルとの差は 2 パーセントではあるが認められるようである。

30 BERT: 事前訓練比較

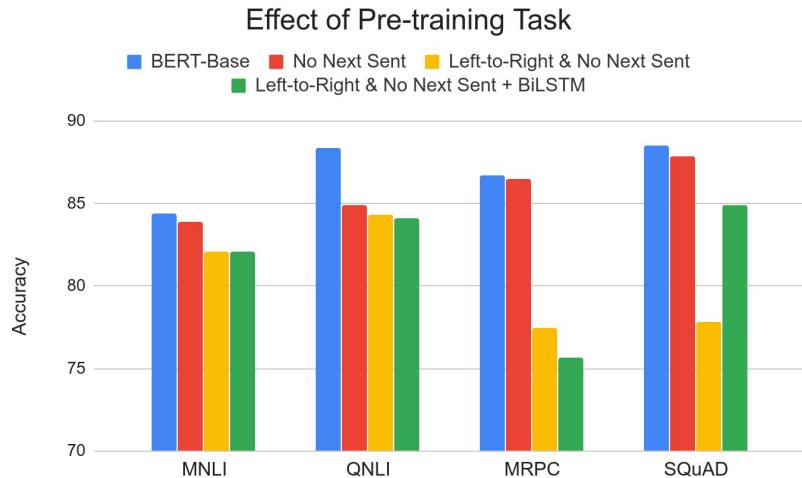


図 14 事前訓練の効果比較

図には事前訓練の比較を示している。全ての事前訓練を用いた場合が青、次文訓練を除いた場合が赤、従来型言語モデルで次文予測課題をした場合を黄、従来型言語モデルで次文予測課題なしを緑で描かれている。4種類の下流課題は MNLI, QNLI, MRPC, SQuAD である。下流のファインチューニング課題ごとに精度が分かれるようである。

31 各モデルの特徴

- RoBERTa: BERT の訓練コーパスを巨大(173GB)にし、ミニバッチサイズを大きした
- XLNet: 順列言語モデル。2ストリーム注意
- MT-DNN: BERT ベースの転移学習に重きをおいたモデル
- GPT-2: BERTに基づく。人間超えて 2019 年 2 月時点で炎上騒ぎ
- BERT: Transformerに基づく言語モデル。**マスク化言語モデル**と**次文予測**に基づく事前訓練、各下流課題をファインチューニング。事前訓練されたモデルは一般公開済。
- DistillBERT: BERT の蒸留版
- ELMo: 双方向 RNN による文埋め込み表現
- Transformer: 自己注意に基づく言語モデル。多頭注意、位置符号器。

32 事前訓練とマルチ課題学習

From Liu, He, Chen, & Gao (2019) Fig. 1

33 Transformer: Attention is all you need

$$\text{attention}(Q, K, V) = \text{dropout} \left(\text{softmax} \left(\frac{QK^\top}{\sqrt{d}} \right) \right) V \quad (4)$$

From Vaswani et al. (2017) Fig. 2

34 位置符号器 Position encoders

トランسفォーマーの入力には、上述の単語表現に加えて、位置符号器からの信号も重ね合わされる。位置 i の信号は次式で周波数領域へと変換される：

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad (5)$$

位置符号器による位置表現は、 i 番目の位置情報をワンホット表現するのではなく、周波数領域に変換することで周期情報を表現する試みと見なし得るだろう。

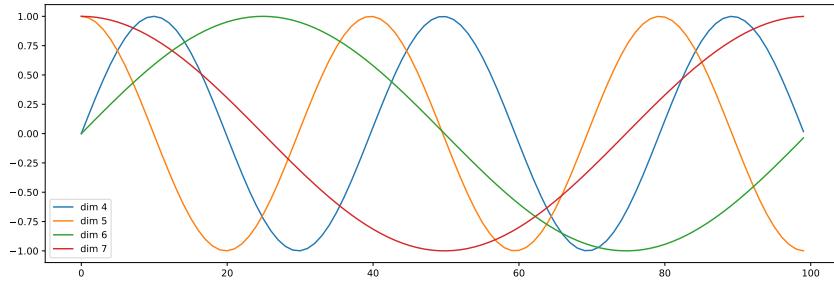


図 15 位置符号化に用いられる符号化

このようにしてできた値を入力側と出力側で下図のように連結させたものがトランسفォーマーである。

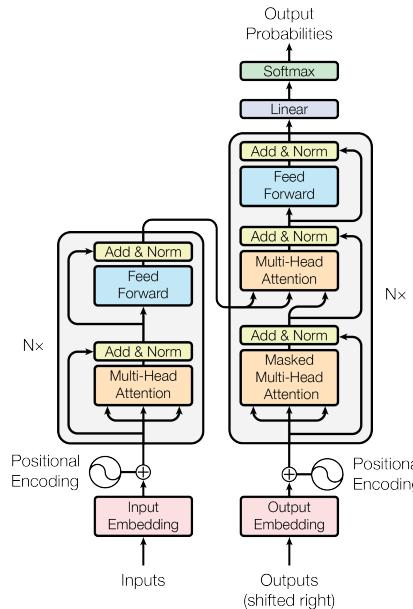


図 16 From Vaswani et al. (2017) Fig. 1

これまで見てきたように、トランسفォーマーでは入力信号に基づいて情報の変換が行なわれる。この意味ではトランسفォーマーにおける多頭自己注意 MHSA とはボトムアップ注意の変形であるとみなしうる。逆言すれば、RNN のように過去の履歴をすべて保持しているわけではないので、系列情報については、position encoders に頼っている側面が指摘できる。

35 BERT, GPT, ELMo 事前訓練の違い

- BERT: トランسفォーマー、マスク化言語モデル、次文予測課題

- GPT: 順方向トランスフォーマー
- ELMo: 双方向 RNN による中間層の連結

36 多言語対応

From Lample & Conneau (2019) Fig. 1

37 BERT の発展

From <https://towardsdatascience.com/a-review-of-bert-based-models-4fffdc0f15d58>

38 BERT: 埋め込みモデルによる構文解析

BERT の構文解析能力を下図示した。各単語の共通空間に射影し、単語間の距離を計算することにより構文解析木と同等の表現を得ることができることが報告されている Hewitt & Manning (2019)。

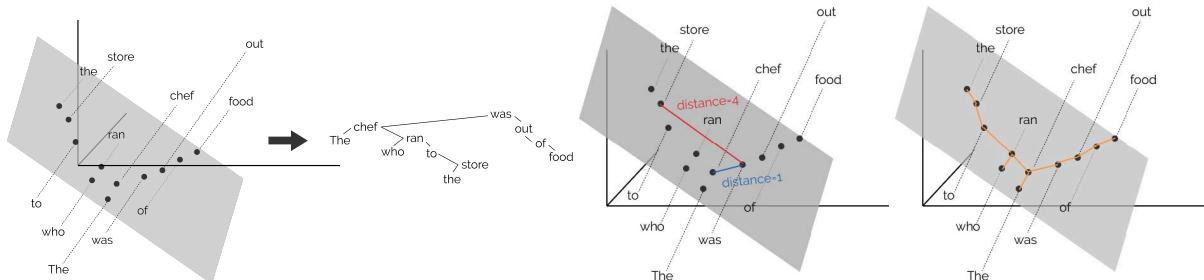


図 17 BERT による構文解析木を再現する射影空間 From <https://github.com/john-hewitt/structural-probes>

word2vec において単語間の距離は内積で定義されていた。このことから、文章を構成する単語で張られる線形内積空間内の距離が構文解析木を与えると見なすことは不自然ではない。そこで構文解析木を再現するような射影変換を見つけることができれば BERT を用いて構文解析が可能となる。例えば上図における chef と store と was の距離を解析木を反映するような空間を見つけ出すことに相当する。2 つの単語 w_i, w_j とし単語間の距離を $d(w_i, w_j)$ とする。適当な変換を施した後の座標を h_i, h_j とすれば、求める変換 B は次式のような変換を行なうことに相当する:

$$\min_B \sum_l \frac{1}{|s_\ell|^2} \sum_{i,j} \left(d(w_i, w_j) - \|B(h_i - h_j)\|^2 \right) \quad (6)$$

ここで ℓ は文 s の訓練文のインデックスであり、各文の長さで規格化することを意味している。

39 BERT 実装

BERT 実装のパラメータを以下に示した。現在配布されている BERT-base あるいは性能が良い BERT-large は各層のニューロン数と全体の層数である。

- データ: Wikipedia (2.5B words) + BookCorpus (800M words)
- バッチサイズ: 131,072 words (1024 sequences \times 128 length or 256 sequences \times 512 length)
- 訓練ステップ: 1M steps (40 epochs)
- 最適化アルゴリズム: AdamW, 1e-4 learning rate, linear decay
- BERT-Base: 12 層、各層 768 ニューロン、12 多頭注意
- BERT-Large: 24 層、各層 1024 ニューロン、16 多頭注意
- 訓練時間: 4x4 / 8x8 の TPU で 4 日間

40 LSTM

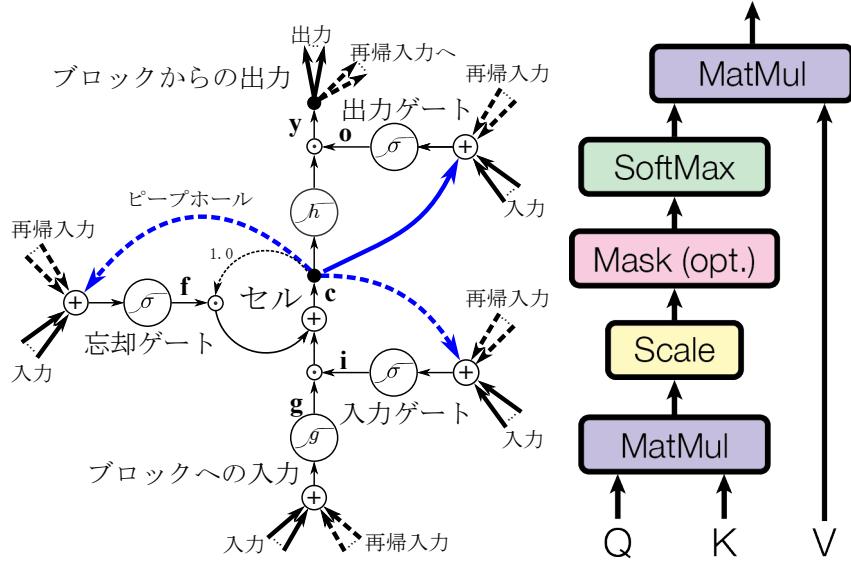


図 18 左: LSTM (浅川, 2015) より, 右: トランスフォーマー Vaswani et al. (2017)

入力ゲートと入力は Q, K と同一視, 出力ゲートと V とは同一視可能?

41 BERT embeddings

42 BERT inside

43 第3部

44 流行りの句

arXiv:2001.04451v2 [cs.LG] 18 Feb 2020

Published as a conference paper at ICLR 2020

REFORMER: THE EFFICIENT TRANSFORMER

Nikita Kitaev*
U.C. Berkeley & Google Research
kitaev@cs.berkeley.edu

Lukasz Kaiser*
Google Research
{lukasz.kaiser,levskaya}@google.com

ABSTRACT

Large Transformer models routinely achieve state-of-the-art results on a number of tasks but training these models can be prohibitively costly, especially on long sequences. We introduce two techniques to improve the efficiency of Transformers. For one, we replace dot-product attention by one that uses locality-sensitive hashing, changing its complexity from $O(L^2)$ to $O(L \log L)$, where L is the length of the sequence. Furthermore, we use reversible residual layers instead of the standard residuals, which allows storing activations only once in the training process instead of N times, where N is the number of layers. The resulting model, the Reformer, performs on par with Transformer models while being much more memory-efficient and much faster on long sequences.

Published as a conference paper at ICLR 2020

ON THE RELATIONSHIP BETWEEN SELF-ATTENTION AND CONVOLUTIONAL LAYERS

Jean-Baptiste Cordonnier, Andreas Loukas & Martin Jaggi
Ecole Polytechnique Fédérale de Lausanne (EPFL)
{first.last}@epfl.ch

ABSTRACT

Recent trends of incorporating attention mechanisms in vision have led researchers to reconsider the supremacy of convolutional layers as a primary building block. Beyond helping CNNs to handle long-range dependencies, Ramachandran et al. (2019) showed that attention can completely replace convolution and achieve state-of-the-art performance on vision tasks. This raises the question: do learned attention layers operate similarly to convolutional layers? This work provides evidence that attention layers can perform convolution and, indeed, they often learn to do so in practice. Specifically, we prove that a multi-head self-attention layer with sufficient number of heads is at least as expressive as any convolutional layer. Our numerical experiments then show that self-attention layers attend to pixel-grid patterns similarly to CNN layers, corroborating our analysis. Our code is publicly available¹.

Are Transformers universal approximators of sequence-to-sequence functions?

Chulhee Yun*
MIT
chulheey@mit.edu

Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J. Reddi, Sanjiv Kumar
Google Research New York
{bsrinadh,ankitsrawat,sashank,sanjivk}@google.com

Abstract

Despite the widespread adoption of Transformer models for NLP tasks, the expressive power of these models is not well-understood. In this paper, we establish that Transformer models are universal approximators of continuous *permutation equivariant* sequence-to-sequence functions with compact support, which is quite surprising given the amount of shared parameters in these models. Furthermore, using positional encodings, we circumvent the restriction of permutation equivariance, and show that Transformer models can universally approximate *arbitrary* continuous sequence-to-sequence functions on a compact domain. Interestingly, our proof techniques clearly highlight the different roles of the self-attention and the feed-forward layers in Transformers. In particular, we prove that fixed width self-attention layers can compute *contextual mappings* of the input sequences, playing a key role in the universal approximation property of Transformers. Based on this insight from our analysis, we consider other simpler alternatives to self-attention layers and empirically evaluate them.

45 Residual attention

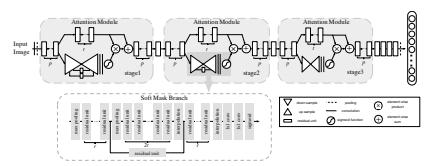
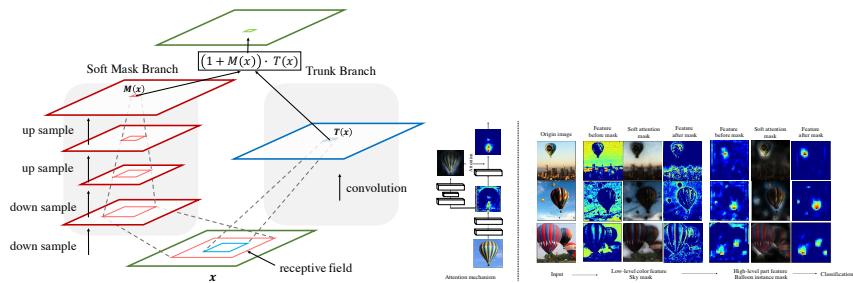


図 19 Wang et al. (2017) Fig. 1, 2, 3

46 A2 net

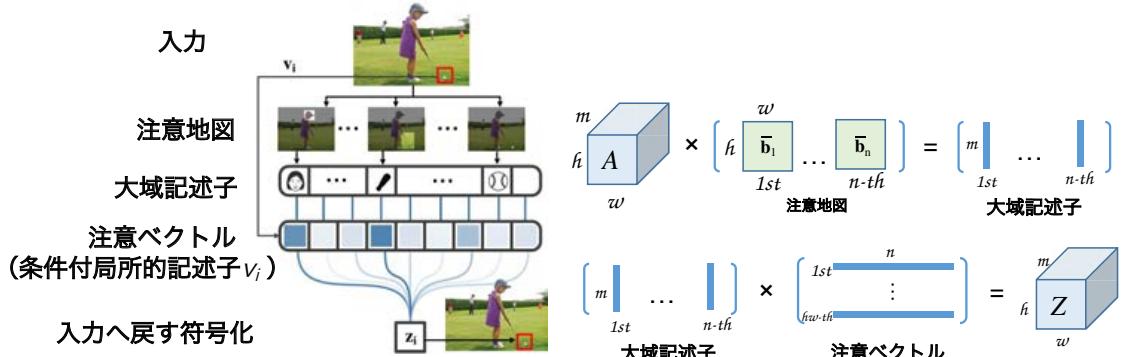
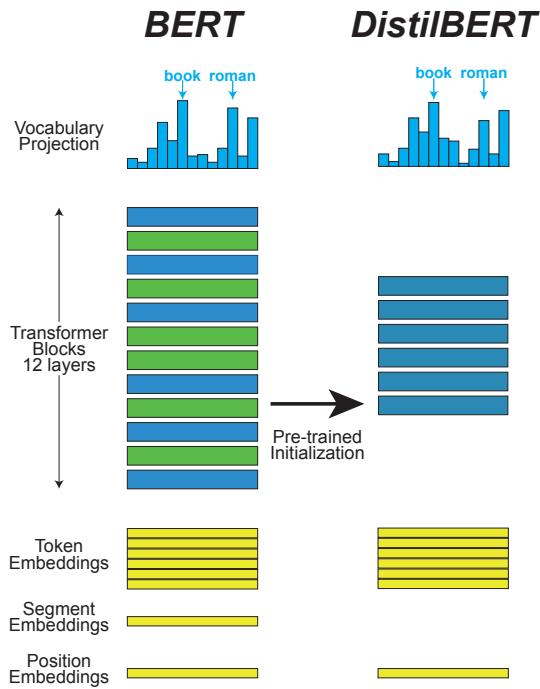


図 20 From Chen et al. (2018) Fig. 1

47 DistilBERT



3つの損失関数 Sanh et al. (2020):

1. 知識蒸留損失
2. マスク化言語モデル損失
3. コサイン損失

48 Relationship between self-attention and convolution

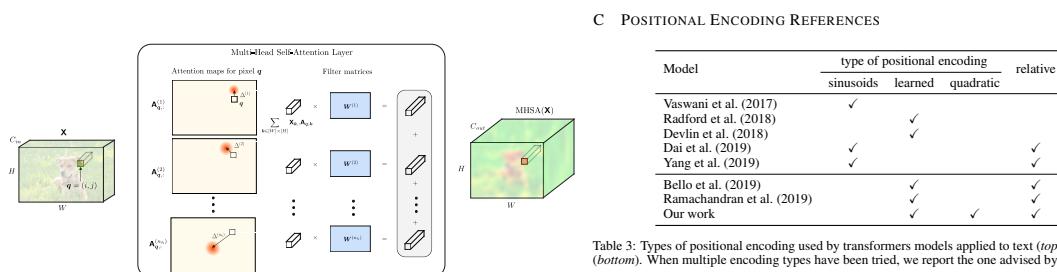


Table 3: Types of positional encoding used by transformers models applied to text (top) and images (bottom). When multiple encoding types have been tried, we report the one advised by the authors.

From

図 21 Cordonnier et al. (2020)

49 第3部まとめ

- MHSA は畠み込みと同等の能力がありそうである。
- Reformer に見られるように position encodings を工夫する余地は残されているように思われる。

50 第4部

不易の句

51 Dicotomy

- ボトムアップとトップダウン
- 何と何処(腹側 背側)
- 特徴、対象、場所へ向けられるの注意
- 外発的、内発的 注意

52 関連脳領域

- FEF 前頭眼野 Monosov & Thompson (2009)
- Lateral Intraparietal area (LIP) 側頭頭頂領域 Wardak, Olivier, & Duhamel (2004)
- Superior Colliculus(SC) 上丘 Krauzlis, Lovejoy, & Zénon (2013)
- PFC 前頭皮質 Miller & Cohen (2001)
- VPA Bichot, Heard, DeGennaro, & Desimone (2015)

53 認知心理学分野

- フィルタリング Broadbent (1958), 減衰説 Treisman (1969)
- 特徴統合理論 Treisman & Gelade (1980);Treisman (1988)
- Guided Search 2.0 Wolfe (1994)
- 目標／妨害刺激類似性: Duncan & Humphreys (1989, 1992)
- サーチライト(スポットライト) 仮説 Crick (1984), ズームレンズ Eriksen & St.James (1986)
- 勝者占有回路 Koch & Ullman (1985) = softmax

54 計算モデル (Implementation)

- Milanese, Wechsler, Gill, Bost, & Pun (1994)
- Itti, Koch, & Niebur (1998)
- Borji & Itti (2013) SOTA

55 総説論文

- Itti & Koch (2001)
- Knudsen (2007)
- Petersen & Posner (2012)
- Kimura, Yonetani, & Hirayama (2013)
- Itti & Borji (2015) Oxford Handbook of attention

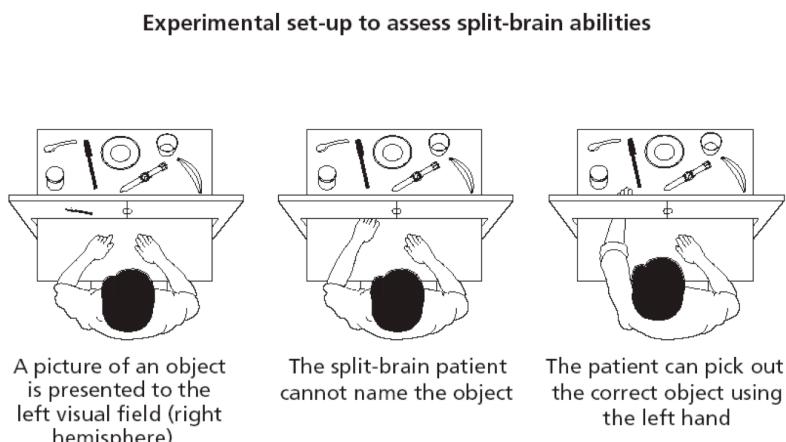
56 深層学習系

- 自動翻訳 Bahdanau et al. (2015); Luong et al. (2015)
- 画像脚注付け Vinyals, Toshev, Bengio, & Erhan (2015)
- 注意 Wang & Shen (2018)

57 温故知新

- 脳梁切断患者による分離脳 Sperry (1961)
- 半側空間無視 Heilman & Valenstein (1979)
- 頭頂葉損傷患者の注意のディスエンゲージメント Posner (1980)
- 両耳分離聴実験, カクテルパーティ効果 Broadbent (1958); Treisman (1964)
- 特徴統合理論 Treisman (1988); Treisman & Gelade (1980)
- 計算論的モデル サーチライト (スポットライト) 仮説 Crick (1984)
- モデルとデータセット公開, 競技会 Itti & Koch (2001); Itti & Borji (2014)
- DeepGazeII Kümmerer, Wallis, Gatys, & Bethge (2017)

58 分離脳 Split brain



From Sperry (1968)

図 22 Fig. 5

59 半側空間無視

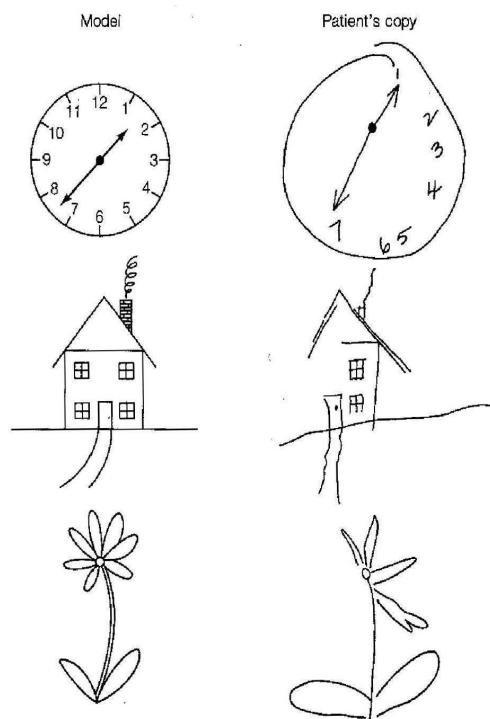


図 23 From Bloom & Lazerson (1988) Fig. 17-6

60 ポズナーとコーヘン

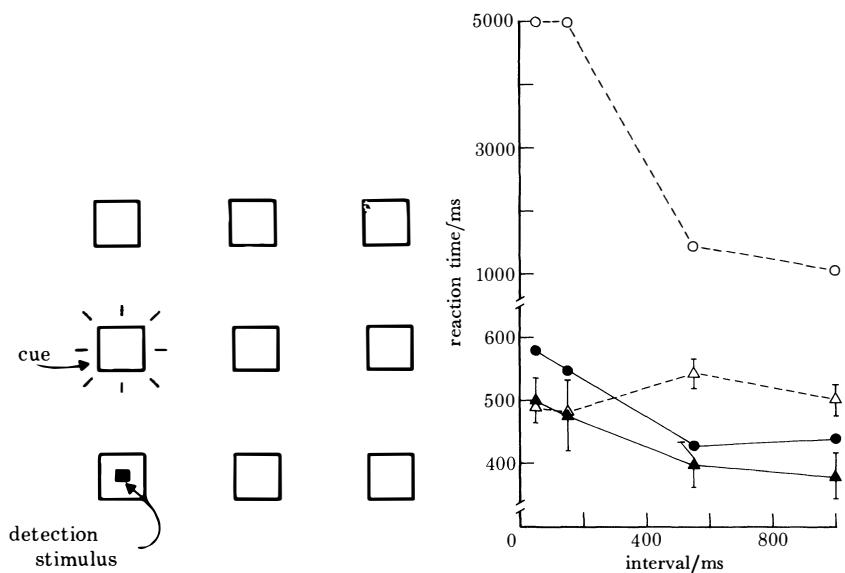


図 24 From Posner (1980) Fig. 1, Fig.6: 右頭頂葉障害を呈した患者 (R.S.) の結果。円:ターゲットが左視野提示、三角:ターゲット右視野提示。白点線:非有効手がかり、黒実線:有効手がかり。横軸は ISI。縦軸は反応時間中央値

61 特徴統合理論 (FIT)

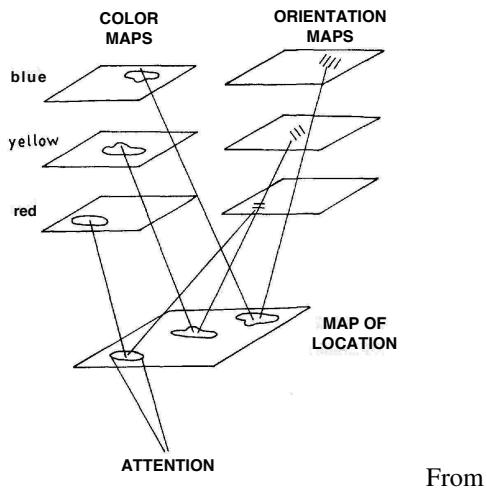


図 25 Treisman & Souther (1985) Fig. 9

62 探索非対称性 search asymmetry

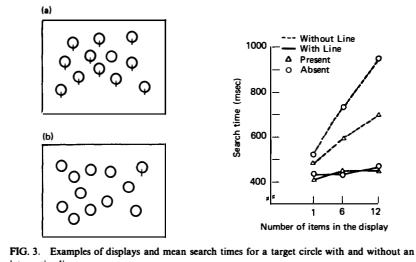


図 26 From Treisman (1988) Fig. 3

上図右の結果は横軸に同時に提示された刺激の個数であり、縦軸は反応時間です。線分特徴が存在する刺激 (Q) が目標となるか、存在しない (O) が目標となるかによって反応時間に差が認められます。結果は点線、すあんわち特徴が存在しない目標を探索する条件、点線で描画、では同時に提示された刺激数が増加するに従って反応時間が増大します。一方、特徴が存在する目標を探索する条件では、同時提示された刺激の個数によらず反応時間は平坦になります。以下に同様な実験結果を示しました。

63 スポットライトメタファー

- スポットライトメタファー Crick (1984)

Attention can be likened to a spotlight that enhances the efficiency of detection of events within its beam.

Unlike when acuity is involved, the effect of the beam is not related to the fovea. When the fovea is unilluminated by attention, its ability to lead to detection is diminished, as would be the case with any other area of the visual system. Posner p.172

- Summerfield, Lepsius, Gitelman, Mesulam, & Nobre (2006) は AI の研究にも影響
- ネットワークの内部メモリから読み出す情報を選択するために注意機構
- 機械翻訳 Bahdanau et al. (2015), NTM Graves et al. (2016)
- コンテンツアドレス Hopfield (1982)
- BERT Devlin et al. (2018)

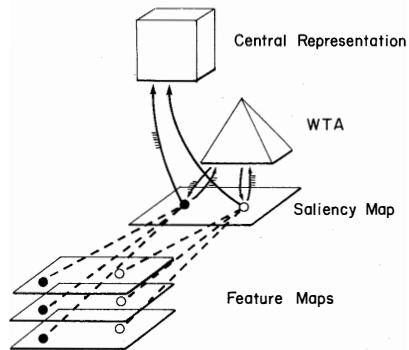


図 27 From Koch & Ullman (1985) Fig. 5

64 Inhibition of Return (IOR)

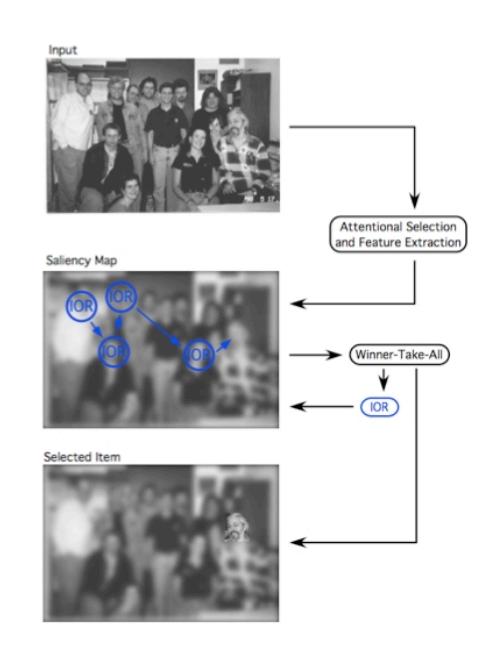


図 28 From http://www.scholarpedia.org/article/Inhibition_of_return

From The superior colliculus (SC) has been implicated as the neural substrate for IOR through four converging, but indirect, lines of evidence.

1. IOR is abnormal in patients with midbrain degeneration due to progressive supranuclear palsy (PSP).
2. It is preserved in patients with hemianopia, a condition in which only extrageniculate pathways are available to process visual information.
3. It is present in newborn infants, in whom the geniculostriate pathways are not yet developed.
4. It is generated asymmetrically in temporal and nasal visual fields, suggesting retinotectal mediation.

65 ガイド付き探索モデル Guided Search 2.0

最初にトップダウン注意を明示的に示した ガイド付き探索モデル Wolfe (1994)

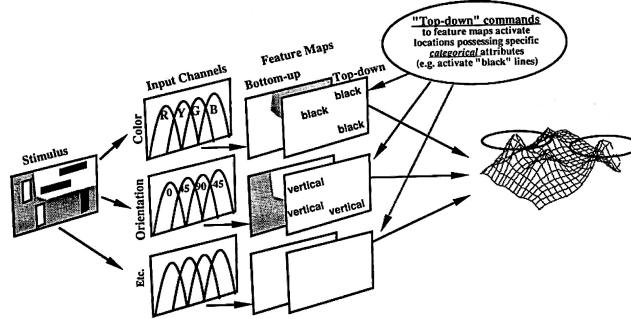


図 29 From Wolfe (1994) Fig. 2

Itti & Borji (2015) の総説論文からそれまでのモデルの概説図

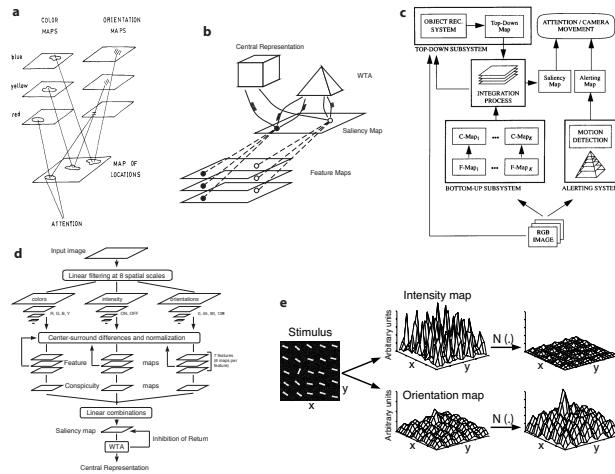


図 30 From Itti & Borji (2015) Fig. 2

66 Friston's attention

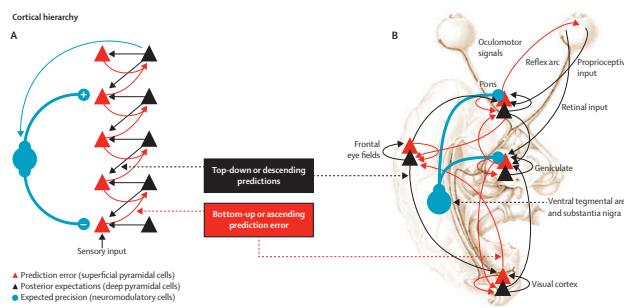


図 31 From Friston et al. (2014) Fig. 1

67 上丘 SC

From Olshausen, Anderson, & Esson (1993) Fig. 10a

- 灵長類の視覚系の動作は注意を伴う視線の移動により外界を認識
- すべての入力を並行して処理するのではなく、視覚的注意は場所や物体間の遷移 Koch & Ullman (1985); Moore & Zirnsak (2017); Posner & Petersen (1990)
- 情報の優先順位付け、取捨選択 Olshausen et al. (1993); Salinas & Abbott (1997)

68 リズム現象

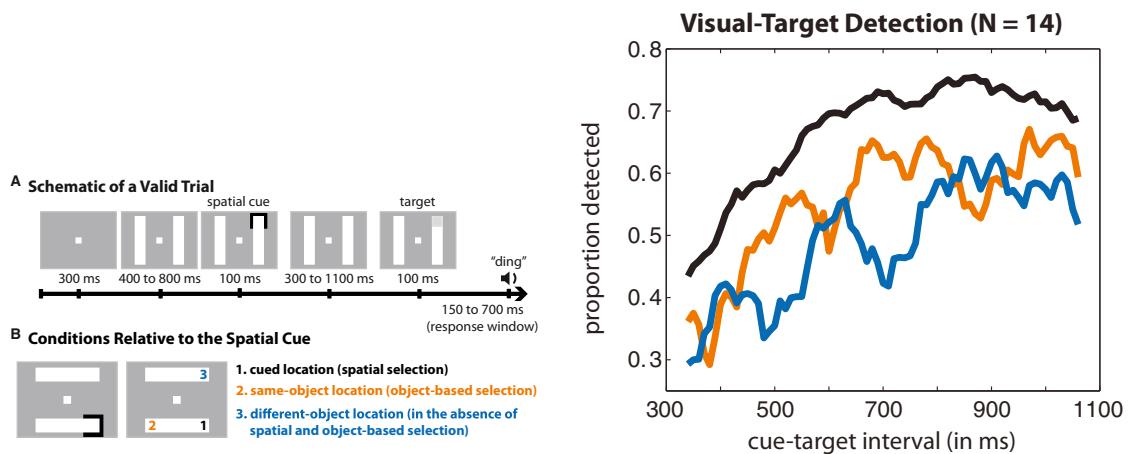


図 32 From Fiebelkorn et al. (2013) Fig. 1 and Fig. 2a

69 リズム現象 (2)

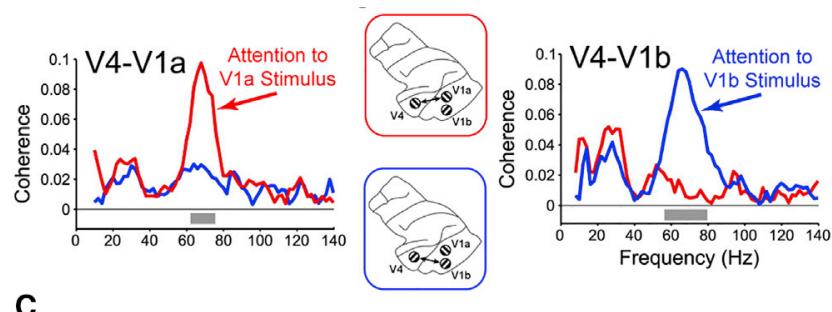


図 33 From Buschman & Kastner (2015) Fig. 3b

70 リズム現象 (3)

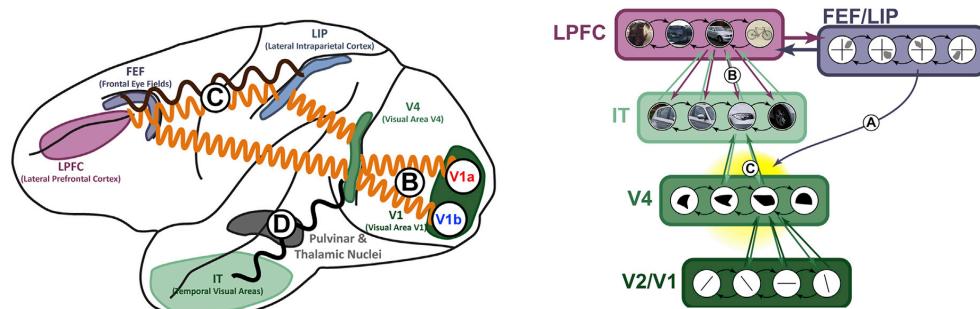


図 34 From Buschman & Kastner (2015) Fig. 3a, Fig. 6

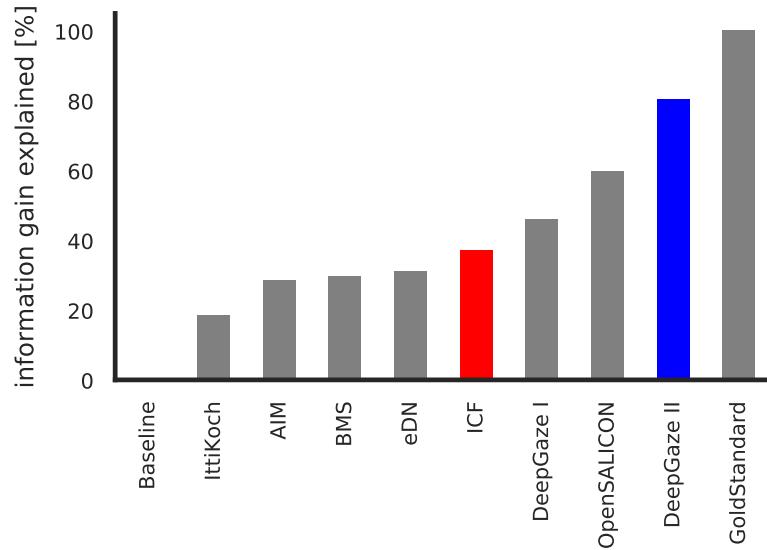


図 36 From Kümmeler et al. (2017) Fig. 2

71 DeepGaze II

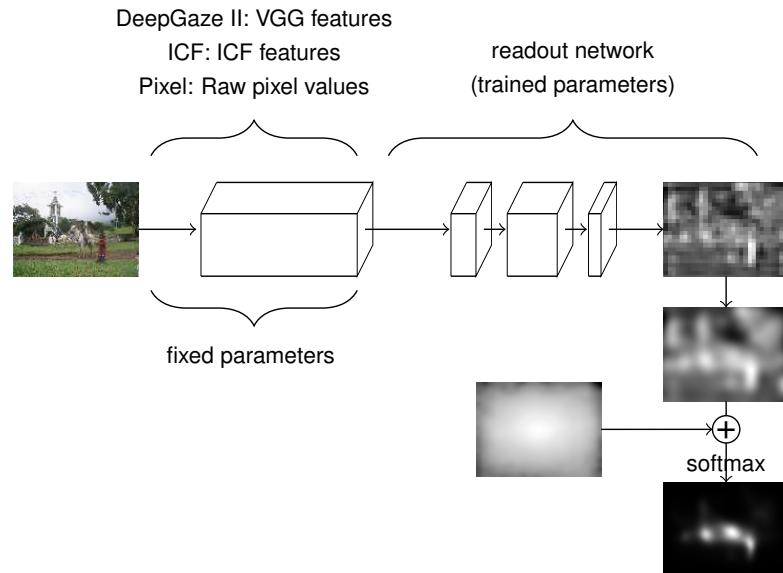


図 35 From Kümmeler et al. (2017) Fig. 2

72 DeepGaze II (2)

DeepGazeII より成績の良い最右の棒は人間の眼球運動データ

73 DeepGaze II (3)

Model	IG	IGE	AUC	sAUC	NSS
Centerbias	0.00	0.0	79.6	50.0	1.22
Pixel	0.13	10.7	81.2	60.2	1.38
IttiKoch [16]	0.23	18.6	82.3	64.1	1.41
AIM [6]	0.27	22.6	82.9	65.6	1.50
eDN [48]	0.38	31.1	83.8	68.7	1.61
ICF	0.45	37.2	84.4	70.1	1.74
DeepGaze I [32]	0.56	46.1	85.8	73.0	1.92
OpenSALICON [46]	0.73	59.7	86.4	74.2	2.14
DeepGaze II	0.98	80.3	88.3	77.7	2.48
Gold Standard	1.22	100.0	89.9	81.2	2.82

図 37 From Kümmerer et al. (2017) Fig. 3

IG: 情報ゲイン, IGE: 修正情報ゲイン, AUC: area under the ROC curve, sAUC: シャッフル精度, NSS: 正規化済キヤンパス顕在性 normalized scanpath saliency

74 DeepGaze III

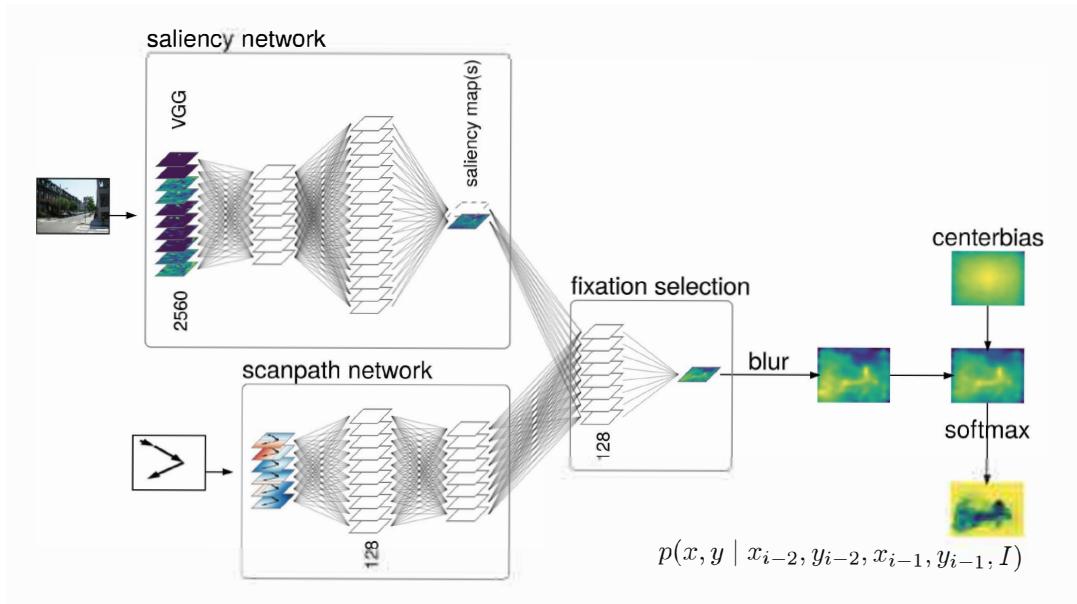


図 38 From Kümmerer et al. (2019) Fig. 1

75 ヘルムホルツマシン

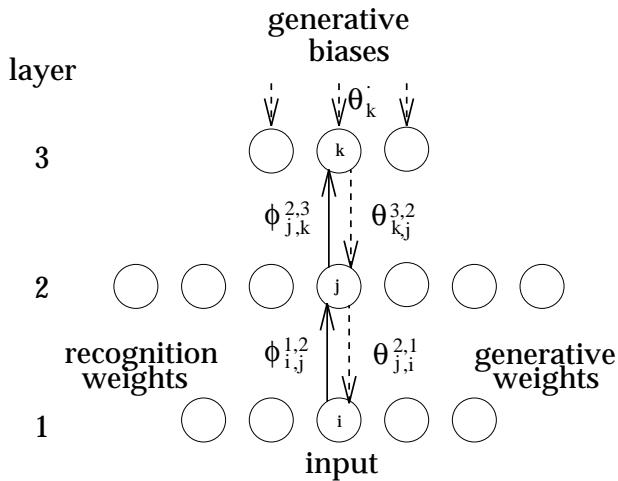


図 39 Dayan et al. (1995);Hinton et al. (1995)

76 ヘルムホルツマシン

$$\log p(d|\theta) = - \sum Q_a E_a - \sum Q_a \log Q_a + \sum Q_a \log \left(\frac{Q_a}{P_a} \right) = -F(d; \theta, Q) + \sum_a Q_a \log \left(\frac{Q_a}{P_a} \right) \quad (7)$$

$$q^{(l)} \left(\phi, \mathbf{s}^{(l-1)} \right) = \sigma \left(\sum s^{l-1} \phi^{(l-1,l)} \right) \quad (8)$$

$$Q_a(\phi, d) = \prod \prod \left[q^{(l)} \left(\phi, \mathbf{s}^{(l-1)} \right) \right]^{s^l} \left[1 - q^{(l)} \left(\phi, \mathbf{s}^{(l-1)} \right) \right]^{1-s} \quad (9)$$

$$p_j^{(l)} \left(\theta, \mathbf{s}^{(l+1)} \right) = \sigma \left(\sum s^{(l+1)} \theta^{(l+1)} \right) \quad (10)$$

$$p(\alpha|\theta) = \prod \prod \left[p_j^{(l)} \left(\theta, \mathbf{s}^{(l+1)} \right) \right] \quad (11)$$

77 モデル: ヘルムホルツマシン

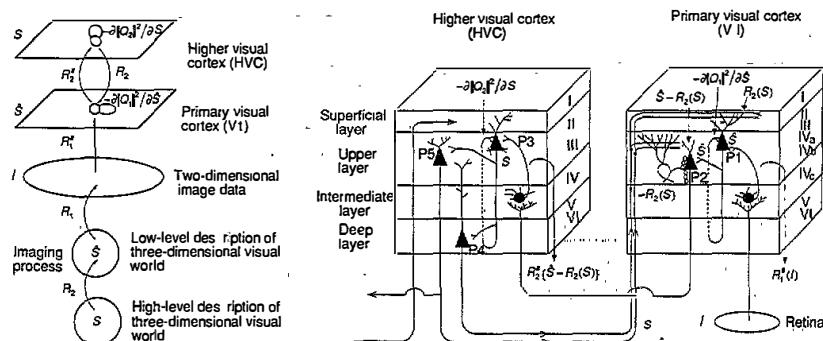


図 40 From Kawato et al. (1993) Fig. 1 より

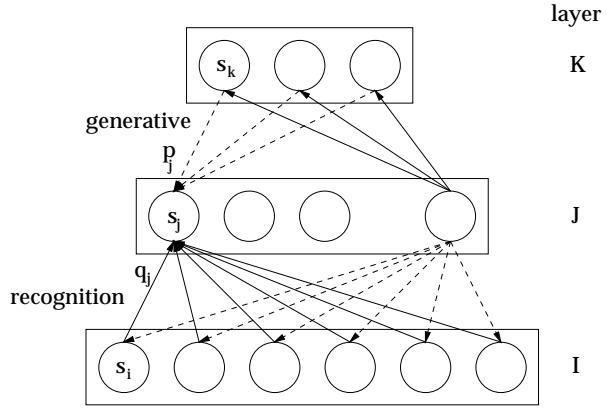


図 41 From Hinton et al. (1995) Fig. 1 より

- 上位層は下位層からの情報をサンプリング → 認識形成 **トップダウン**
- 下位層は上位層からの情報を受けとる → 情報再構成 **ボトムアップ**

ボトムアップ処理による認識とトップダウン処理による（こう見えるはずだという思い込みの）生成を n ($n = 2, \dots, 4$) 回繰り返す → **パレイドリア成立**

78 定式化

思い込みの印象 α と入力画像 d を用いて %% の記述長は、単なる前隠れ層ユニットの記述損失であり

$$C(\alpha, d) = C(\alpha) + C(d|\alpha) = \sum_{\ell \in L} \sum_{j \in \ell} C(s_j^\alpha) + \sum_i C(s_i^d | \alpha) \quad (12)$$

上式を用いて結合係数の更新を行う

$$\Delta w_{kj} = \epsilon s_k^\alpha \left(s_j^\alpha - p_j^\alpha \right), \quad (13)$$

$$C(d) = \sum_\alpha Q(\alpha|d) C(\alpha, d) - \left[- \sum_\alpha Q(\alpha|d) \log Q(\alpha|d) \right]. \quad (14)$$

$$p(\alpha|d) = \frac{e^{-C(\alpha,d)}}{\sum_\beta e^{-C(\beta,d)}} \quad (15)$$

$$\Delta s_{j,t+1} = \epsilon s_{j,t}^\gamma (s_{j,t}^\gamma - q_{j,t}^\gamma) \quad (16)$$

全体の良い表象が得られるまで、すなわち下位層の活性を再構築するように複数回繰り返す

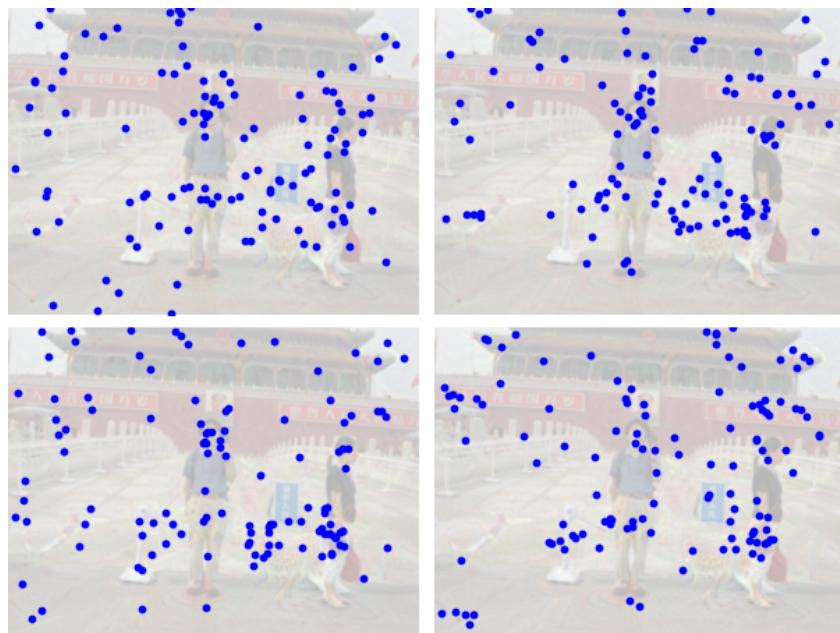
79 計算例



80 計算例



81 計算例 (2) 眼球運動のサンプリング



引用文献

- Bahdanau, D., Cho, K., & Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In Y. Bengio & Y. LeCun (Eds.), *Proceedings in the International Conference on Learning Representations (ICLR)*. San Diego, CA, USA.
- Bichot, N. P., Heard, M. T., DeGennaro, E. M., & Desimone, R. (2015). A source for feature-based attention in the prefrontal cortex. *Neuron*, 88, 832-844.
- Bloom, F. E., & Lazerson, A. (1988). *Brain, mind, and behavior* (2nd ed.). New York, NY: Freeman.
- Borji, A., & Itti, L. (2013). State-of-the-art in visual attention modeling. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 35, 185-207.
- Broadbent, D. E. (1958). *Perception and communication*. Oxford, UK: Pergamon.
- Buschman, T. J., & Kastner, S. (2015). From behavior to neural dynamics: An integrated theory of attention. *Neuron*, 88, 127-144.
- Chen, Y., Kalantidis, Y., Li, J., Yan, S., & Feng, J. (2018). a^2 -nets: Double attention networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances in neural information processing systems 31* (pp. 352-361). Curran Associates, Inc.
- Cordonnier, J.-B., Loukas, A., & Jaggi, M. (2020). On the relationship between self-attention and convolutional layers. *arXiv preprint, [cs.LG]*.
- Crick, F. (1984). Function of the thalamic reticular complex: the search light hypothesis. *Proceedings of the National Academy of Sciences*, 81, 4586-4590.
- Dayan, P., Hinton, G. E., Neal, R. M., & Zemel, R. S. (1995). The helmholtz machine. *Neural Computation*, 7, 889-904.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint*.
- Duncan, J., & Humphreys, G. W. (1989). Visual search and stimulus similarity. *Psychological Review*, 96, 433-458.
- Duncan, J., & Humphreys, G. W. (1992). Beyond the search surface: Visual search and attentional engagement. *Journal of Experimental Psychology: Human Perception and Performance*, 18, 578-588.
- Elman, J. L. (1990). Finding structure in time. *Cognitive Science*, 14, 179-211.
- Eriksen, C. W., & St.James, J. D. (1986). Visual attention within and around the field of focal attention: A zoom lens model. *Perception and Psychophysics*, 40, 225-240.
- Fiebelkorn, I. C., Saalmann, Y. B., & Kastner, S. (2013). Rhythmic sampling within and between objects despite sustained attention at a cued location. *Current Biology*, 23, 2553-2558.
- Friston, K. J., Stephan, K. E., Montague, R., & Dolan, R. J. (2014). Computational psychiatry: the brain as a phantastic organ. *The Lancet Psychiatry*, 1, 148-158.
- Gers, F. A., Schmidhuber, J., & Cummins, F. (1999). Learning to forget: Continual prediction with LSTM. In *Artificial Neural Networks ICANN 99. Ninth International Conference on* (Vol. 2, pp. 850-855). Edinburgh, Scotland.
- Graves, A., Wayne, G., & Danihelka, I. (2014). Neural Turing machines. *arXiv:1410.5401*.
- Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Agnieszka Grabska-Barwińska, Colmenarejo, S. G., Grefenstette, E., Ramalho, T., Agapiou, J., Badia, A. P., Hermann, K. M., Zwols, Y., Ostrovski, G., Cain, A., King, H., Summerfield, C., Blunsom, P., Kavukcuoglu, K., & Hassabis, D. (2016). Hybrid computing using a neural network with dynamic external memory. *Nature*, 538, 471-476.
- Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R., & Schmidhuber, J. (2015). LSTM: A search space odyssey. *arXiv:1503.04069*.
- Heilman, K. M., & Valenstein, E. (1979). Mechanisms underlying hemispatial neglect. *The Annals of Neurology*, 5, 166-170.
- Hewitt, J., & Manning, C. D. (2019). A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)* (pp. 4129-4138). Minneapolis, Minnesota: Association for Computational Linguistics.
- Hinton, G. E., Dayan, P., Frey, B. J., & Neal, R. M. (1995). The "wake-sleep" algorithm for unsupervised neural networks. *Science*, 268, 1158-1161.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9, 1735-1780.
- Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79, 2554-2558.
- Itti, L., & Borji, A. (2014). Computational models: Bottom-up and top-down aspects. In A. C. Nobre & S. Kastner (Eds.), *The oxford handbook of attention* (p. 1122-1158). Oxford University Press.
- Itti, L., & Borji, A. (2015). Computational models of attention. *arXiv preprint*.
- Itti, L., & Koch, C. (2001). Computational modelling of visual attention. *Nature Reviews Neuroscience*, 2, 1-11.
- Itti, L., Koch, C., & Niebur, E. (1998). A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20, 1254-1259.
- Jones, K. S. (1972). A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation*, 28, 11-21.
- Kawato, M., Hayakawa, H., & Inui, T. (1993). A forward-inverse optics model of reciprocal connections between visual cortical areas. *Network: Computation in Neural Systems*, 4, 415-422.
- Kim, J., Kim, M., Kang, H., & Lee, K. (2019). U-GAT-IT: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. *arXiv preprint, [cs.CV]*.
- Kimura, A., Yonetani, R., & Hirayama, T. (2013). Computational models of human visual attention and their implementations: A survey. *IEICE Transactions of Information & Systems*, E96-D, 562-578.
- Knudsen, E. I. (2007). Fundamental components of attention. *Annual Review of Neuroscience*, 30, 57-78.
- Koch, C., & Ullman, S. (1985). Shifts in selective visual attention: towards the underlying neural circuitry. *Human Neurobiology*, 4, 219-227.
- Krauzlis, R. J., Lovejoy, L. P., & Zénon, A. (2013). Superior colliculus and visual spatial attention. *Annual Review of Neuroscience*, 36.

- Kümmerer, M., Wallis, T. S., & Bethge, M. (2019). DeepGaze III: Using deep learning to probe interactions between scene content and scanpath history in fixation selection. In *Proceedings of Cognitive Computational Neuroscience* (p. 542–545). Berlin, Germany.
- Kümmerer, M., Wallis, T. S. A., Gatys, L. A., & Bethge, M. (2017). Understanding low- and high-level contributions to fixation prediction. In *The IEEE International Conference on Computer Vision (ICCV)* (pp. 4789–4798). Venice, Italy.
- Lample, G., & Conneau, A. (2019). Cross-lingual language model pretraining. *arXiv preprint, 1901.07291v1 [cs.CL]*.
- Liu, X., He, P., Chen, W., & Gao, J. (2019). Multi-task deep neural networks for natural language understanding. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (pp. 4487–4496). Florence, Italy: Association for Computational Linguistics.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A robustly optimized bert pretraining approach. *arXiv preprint*.
- Luong, M.-T., Pham, H., & Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. *arXiv preprint, cs.CL, 1508.04025*.
- Manning, C. D., & Schütze, H. (1999). *Foundations of statistical natural language processing*. Cambridge, MA, USA: MIT press.
- Mikolov, T., Karafiat, M., Burget, L., Černocký, J. H., & Khudanpur, S. (2010). Recurrent neural network based language model. In T. Kobayashi, K. Hirose, & S. Nakamura (Eds.), *Proceedings of INTERSPEECH2010* (pp. 1045–1048). Makuhari, JAPAN. (First proposal of dynamic evaluation)
- Mikolov, T., Kombrink, S., Burget, L., Černocký, J. H., & Khudanpur, S. (2011). Extensions of recurrent neural network language model. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Prague, Czech Republic.
- Milanese, R., Wechsler, H., Gill, S., Bost, J.-M., & Pun, T. (1994). Integration of bottom-up integration of bottom-up and top-down cues for visual attention using non-linear relaxation. In *The proceedings of CVPR, IEEE – Institute of Electrical and Electronics Engineers* (p. 781-785). Dallas Texas, USA.
- Miller, E. K., & Cohen, J. D. (2001). An integrative theory of prefrontal cortex function. *Annual Review of Neuroscience, 24*.
- Mishra, N., Rohaninejad, M., Chen, X., & Abbeel, P. (2018). A simple neural attentive meta-learner. *arXiv preprint, [cs.AI]*.
- Mnih, V., Heess, N., Graves, A., & Kavukcuoglu, K. (2014). Recurrent models of visual attention. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, & K. Weinberger (Eds.), *Advances in neural information processing systems 27* (pp. 2204–2212). Curran Associates, Inc.
- Monosov, I. E., & Thompson, K. G. (2009). Frontal eye field activity enhances object identification during covert visual search. *Journal of Neurophysiology, 102*, 3656–3672.
- Moore, T., & Zirmsak, M. (2017). Neural mechanisms of selective visual attention. *Annual Review of Psychology, 68*, 47–72.
- Olshausen, B. A., Anderson, C. H., & Essen, D. C. V. (1993). A neurobiological model of visual attention and invariant pattern recognition based on dynamic routing of information. *The Journal of Neuroscience, 13*, 4700–4719.
- Petersen, S. E., & Posner, M. I. (2012). The attention system of the human brain: 20 years after. *Annual Review of Neuroscience, 35*, 73–89.
- Posner, M. I. (1980). Orienting of attention. *Quarterly Journal of Experimental Psychology, 32*, 3–25.
- Posner, M. I., & Petersen, S. E. (1990). The attention system of the human brain. *Annual Review of Neuroscience, 13*, 25–42.
- Ramachandran, P., Parmar, N., Vaswani, A., Bello, I., Levskaya, A., & Shlens, J. (2019). Stand-alone self-attention in vision models. *arXiv preprint, [cs.CV]*.
- Salinas, E., & Abbott, L. F. (1997). Invariant visual responses from attentional gain fields. *Journal of Neurophysiology, 77*, 3267–3272.
- Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2020). Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint*.
- Sperry, R. W. (1961). Cerebral organization and behavior. *Science, 133*, 1749–1757.
- Sperry, R. W. (1968). Hemisphere disconnection and unity in conscious awareness. *American Psychologist, 28*, 723–733.
- Summerfield, J. J., Lepsius, J., Gitelman, D. R., Mesulam, M. M., & Nobre, A. C. (2006). Orienting attention based on long-term memory experience. *Neuron, 49*, 905–916.
- Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, & K. Weinberger (Eds.), *Advances in Neural Information Processing Systems (NIPS)* (Vol. 27, pp. 3104–3112). Montreal, BC, Canada.
- Treisman, A. (1964). Selective attention in man. *British Medical Bulletin, 20*, 12–16.
- Treisman, A. (1988). Feature and objects: The fourteenth bartlett memorial lecture. *The quarterly Journal of Experimental Psychology, 40A*, 201–237.
- Treisman, A., & Gelade, G. (1980). A feature integration theory of attention. *Cognitive Psychology, 12*, 97–136.
- Treisman, A., & Souther, J. (1985). Search asymmetry: A diagnostic for preattentive processing of separable features. *Journal of Experimental Psychology: General, 114*, 285–310.
- Treisman, A. M. (1969). Strategies and models of selective attention. *Psychological Review, 76*, 282–299.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., & Kaiser, L. (2017). Attention is all you need. *arXiv preprint, [cs.CL]*.
- Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. In *Computer Vision and Pattern Recognition (CVPR)*. Boston, MA, USA.
- Wang, F., Jiang, M., Qian, C., Yang, S., & Li, C. (2017). Residual attention network for image classification. In *Proceedings of International Conference of Computer Vision (ICCV), IEEE International Conference*. Venice, Italy.
- Wang, W., & Shen, J. (2018). Deep visual attention prediction. *IEEE Transactions on Image Processing, 27*, 2368–2378.
- Wang, X., Girshick, R., Gupta, A., & He, K. (2018). Non-local neural networks. *arXiv preprint, [cs.CV]*.
- Wardak, C., Olivier, E., & Duhamel, J.-R. (2004). A deficit in covert attention after parietal cortex inactivation in the monkey. *Neuron, 42*, 501–508.
- Wolfe, J. M. (1994). Guided search 2.0 a revised model of visual search. *Psychonomic Bulletin and Review, 1*, 202–238.
- Xu, K., Ba, J. L., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R. S., & Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. *arXiv:1502.03044*.
- Zhang, H., Goodfellow, I., Metaxas, D., & Odena, A. (2019). Self-attention generative adversarial networks. *arXiv preprint, [stat.ML]*.