

Attention in Natural Language Processing

Andrea Galassi^{ID}, Marco Lippi^{ID}, and Paolo Torrioni^{ID}

Abstract—Attention is an increasingly popular mechanism used in a wide range of neural architectures. The mechanism itself has been realized in a variety of formats. However, because of the fast-paced advances in this domain, a systematic overview of attention is still missing. In this article, we define a unified model for attention architectures in natural language processing, with a focus on those designed to work with vector representations of the textual data. We propose a taxonomy of attention models according to four dimensions: the representation of the input, the compatibility function, the distribution function, and the multiplicity of the input and/or output. We present the examples of how prior information can be exploited in attention models and discuss ongoing research efforts and open challenges in the area, providing the first extensive categorization of the vast body of literature in this exciting domain.

Index Terms—Natural language processing (NLP), neural attention, neural networks, review, survey.

I. INTRODUCTION

IN MANY problems that involve the processing of natural language, the elements composing the source text are characterized by having each a different relevance to the task at hand. For instance, in aspect-based sentiment analysis, cue words, such as “good” or “bad,” could be relevant to some aspects under consideration, but not to others. In machine translation, some words in the source text could be irrelevant in the translation of the next word. In a visual question-answering task, background pixels could be irrelevant in answering a question regarding an object in the foreground but relevant to questions regarding the scenery.

Arguably, effective solutions to such problems should factor in a notion of relevance, so as to focus the computational resources on a restricted set of important elements. One possible approach would be to tailor solutions to the specific genre at hand, in order to better exploit known regularities of the input, by feature engineering. For example, in the argumentative analysis of persuasive essays, one could decide to give special emphasis to the final sentence. However, such an approach is not always viable, especially if the input is long or very information-rich, such as in text summarization,

where the output is the condensed version of a possibly lengthy text sequence. Another approach of increasing popularity amounts to machine learning the relevance of input elements. In that way, neural architectures could automatically weigh the relevance of any region of the input and take such a weight into account while performing the main task. The commonest solution to this problem is a mechanism known as attention.

Attention was first introduced in natural language processing (NLP) for machine translation tasks by Bahdanau *et al.* [2]. However, the idea of glimpses had already been proposed in computer vision by Larochelle and Hinton [3], following the observation that biological retinas fixate on relevant parts of the optic array, while resolution falls off rapidly with eccentricity. The term visual attention became especially popular after Mnih *et al.* [4] significantly outperformed the state of the art in several image classification tasks as well as in dynamic visual control problems such as object tracking due to an architecture that could adaptively select and then process a sequence of regions or locations at high resolution and use a progressively lower resolution for further pixels.

Besides offering a performance gain, the attention mechanism can also be used as a tool for interpreting the behavior of neural architectures, which are notoriously difficult to understand. Indeed, neural networks are subsymbolic architectures; therefore, the knowledge they gather is stored in numeric elements that do not provide any means of interpretation by themselves. It then becomes hard if not impossible to pinpoint the reasons behind the wrong output of a neural architecture. Interestingly, attention could provide a key to partially interpret and explain neural network behavior [5]–[9], even if it cannot be considered a reliable means of explanation [10], [11]. For instance, the weights computed by attention could point us to relevant information discarded by the neural network or to irrelevant elements of the input source that have been factored in and could explain a surprising output of the neural network.

Therefore, visual highlights of attention weights could be instrumental in analyzing the outcome of neural networks, and a number of specific tools have been devised for such a visualization [12], [13]. Fig. 1 shows an example of attention visualization in the context of aspect-based sentiment analysis.

For all these reasons, attention has become an increasingly common ingredient of neural architectures for NLP [14], [15]. Table I presents a nonexhaustive list of neural architectures where the introduction of an attention mechanism has brought about a significant gain. Works are grouped by the NLP tasks they address. The spectrum of tasks involved is remarkably broad. Besides NLP and computer vision [16]–[18], attentive models have been successfully adopted in many other different fields, such as speech recognition [19]–[21],

Manuscript received July 5, 2019; revised December 17, 2019 and April 17, 2020; accepted August 20, 2020. Date of publication September 10, 2020; date of current version October 6, 2021. This work was supported by Horizon 2020, project A14EU, under Grant 825619. (Corresponding author: Andrea Galassi.)

Andrea Galassi and Paolo Torrioni are with the Department of Computer Science and Engineering (DISI), University of Bologna, 40126 Bologna, Italy (e-mail: a.galassi@unibo.it; paolo.torrioni@unibo.it).

Marco Lippi is with the Department of Sciences and Methods for Engineering (DISMI), University of Modena and Reggio Emilia, 41121 Modena, Italy (e-mail: marco.lippi@unimore.it).

Color versions of one or more of the figures in this article are available online at <https://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2020.3019893

Task: Hotel location

you get what you pay for . not the **cleanest rooms** but bed was **clean** and so was **bathroom** . bring your own **towels** though as very **thin** . service was **excellent** , let us book in at 8:30am ! for **location and price** , **this ca n't be beaten** , but it is **cheap** for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

Task: Hotel cleanliness

you get what you pay for . **not the cleanest rooms but bed was clean and so was bathroom** . bring your own **towels** though as very **thin** . service was excellent , let us book in at 8:30am ! for location and price , this ca n't be beaten , but it is **cheap** for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

Task: Hotel service

you get what you pay for . not the cleanest rooms but bed was **clean** and so was **bathroom** . bring your own **towels** though as very **thin** . **service was excellent** , let us book in at 8:30am ! for location and price , this ca n't be beaten , but it is **cheap** for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

Fig. 1. Example of attention visualization for an aspect-based sentiment analysis task, from [1, Fig. 6]. Words are highlighted according to attention scores. Phrases in bold are the words considered relevant for the task or human rationales.

TABLE I

NONEXHAUSTIVE LIST OF WORKS THAT EXPLOIT ATTENTION, GROUPED BY THE TASK(S) ADDRESSED

Task Addressed	Related Works
Machine Translation Translation Quality Estimation	[2, 6, 8, 29–48, 48–50] [51]
Text Classification Abusive content detection	[7, 8, 10, 11, 52, 53] [54]
Text Summarization	[41, 55–58]
Language Modelling	[59–61]
Question Answering Question Answering over Knowledge Base	[10, 47, 59, 62–75] [76]
Morphology Pun Recognition	[77] [78]
Multimodal Tasks Image Captioning Visual Question Answering Task-oriented Language Grounding	[16, 79] [80–82] [83]
Information Extraction Coreference Resolution Named Entity Recognition Optical Character Recognition Correction	[84, 85] [51, 86] [87]
Semantic Entity Disambiguation Natural Language Inference Semantic Relatedness Semantic Role Labelling Sentence Similarity Textual Entailment Word Sense Disambiguation	[88] [8, 10, 47, 89–96] [93] [97, 98] [96] [75, 99, 100] [101]
Syntax Constituency Parsing Dependency Parsing	[102, 103] [51, 104, 105]
Sentiment Analysis Agreement/Disagreement Identification Argumentation Mining Emoji prediction Emotion Cause Analysis Emotion Classification	[1, 7, 93, 95, 100, 106–120] [121] [57, 122–125] [126] [127, 128] [115]

recommendation [22], [23], time-series analysis [24], [25], games [26], and mathematical problems [27], [28].

In NLP, after an initial exploration by a number of seminal papers [2], [59], a fast-paced development of new attention models and attentive architectures ensued, resulting in a highly diversified architectural landscape. Because of, and adding to,

the overall complexity, it is not unheard of different authors who have been independently following similar intuitions leading to the development of almost identical attention models. For instance, the concepts of inner attention [68] and word attention [41] are arguably one and the same. Unsurprisingly, the same terms have been introduced by different authors to define different concepts, thus further adding to the ambiguity in the literature. For example, the term context vector is used with different meanings by Bahdanau *et al.* [2], Yang *et al.* [52], and Wang *et al.* [129].

In this article, we offer a systematic overview of attention models developed for NLP. To this end, we provide a general model of attention for NLP tasks and use it to chart the major research activities in this area. We also introduce a taxonomy that describes the existing approaches along four dimensions: input representation, compatibility function, distribution function, and input/output multiplicity. To the best of our knowledge, this is the first taxonomy of attention models. Accordingly, we provide a succinct description of each attention model, compare the models with one another, and offer insights on their use. Moreover, we present the examples regarding the use of prior information in unison with attention, debate about the possible future uses of attention, and describe some interesting open challenges.

We restrict our analysis to attentive architectures designed to work with vector representation of data, as it typically is the case in NLP. Readers interested in attention models for tasks where data have a graphical representation may refer to Lee *et al.* [130].

What this survey does not offer is a comprehensive account of all the neural architectures for NLP (for an excellent overview, see [131]) or of all the neural architectures for NLP that uses an attention mechanism. This would be impossible and would rapidly become obsolete because of the sheer volume of new articles featuring architectures that increasingly rely on such a mechanism. Moreover, our purpose is to produce a synthesis and a critical outlook rather than a flat listing of research activities. For the same reason, we do not offer a quantitative evaluation of different types of attention mechanisms since such mechanisms are generally embedded in larger neural network architectures devised to address

specific tasks, and it would be pointless in many cases to attempt comparisons using different standards. Even for a single specific NLP task, a fair evaluation of different attention models would require experimentation with multiple neural architectures, extensive hyperparameter tuning, and validation over a variety of benchmarks. However, attention can be applied to a multiplicity of tasks, and there are no data sets that would meaningfully cover such a variety of tasks. An empirical evaluation is thus beyond the scope of this article. There are, however, a number of experimental studies focused on particular NLP tasks, including machine translation [37], [42], [48], [132], argumentation mining [125], text summarization [58], and sentiment analysis [7]. It is worthwhile remarking that, on several occasions, attention-based approaches enabled a dramatic development of entire research lines. In some cases, such a development has produced an immediate performance boost. This was the case, for example, with the transformer [36] for sequence-to-sequence annotation, as well as with BERT [60], currently among the most popular architectures for the creation of embeddings. In other cases, the impact of attention-based models was even greater, paving the way to radically new approaches for some tasks. This was the influence of Bahdanau *et al.*'s work [2] to the field of machine translation. Likewise, the expressive power of memory networks [59] significantly contributed to the idea of using deep networks for reasoning tasks.

This survey is structured as follows. In Section II, we define a general model of attention and describe its components. We use a well-known machine-translation architecture introduced by Bahdanau *et al.* [2] as an illustration and an instance of the general model. In Section III, we elaborate on the uses of attention in various NLP tasks. Section IV presents our taxonomy of attention models. Section V discusses how attention can be combined with knowledge about the task or the data. Section VI is devoted to open challenges, current trends, and future directions. Section VII concludes this article.

II. ATTENTION FUNCTION

The attention mechanism is a part of a neural architecture that enables to dynamically highlight relevant features of the input data, which, in NLP, is typically a sequence of textual elements. It can be applied directly to the raw input or to its higher level representation. The core idea behind attention is to compute a weight distribution on the input sequence, assigning higher values to more relevant elements.

To illustrate, we briefly describe a classic attention architecture, called RNNsearch [2]. We chose RNNsearch because of its historical significance and for its simplicity with respect to other structures that we will describe further on.

A. Example for Machine Translation and Alignment

RNNsearch uses attention for machine translation. The objective is to compute an output sequence \mathbf{y} that is a translation of an input sequence \mathbf{x} . The architecture consists of an encoder followed by a decoder, as shown in Fig. 2.

The encoder is a bidirectional recurrent neural network (BiRNN) [133] that computes an annotation term \mathbf{h}_i for every

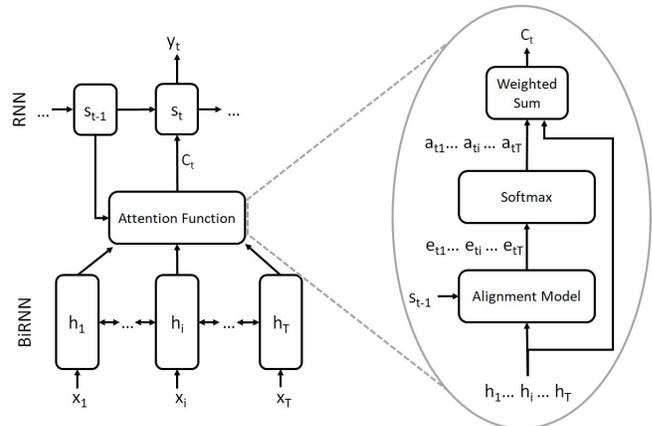


Fig. 2. Architecture of RNNsearch [2] (left) and its attention model (right).

input term x_i of \mathbf{x}

$$(\mathbf{h}_1, \dots, \mathbf{h}_T) = \text{BiRNN}(x_1, \dots, x_T). \quad (1)$$

The decoder consists of two cascading elements: the attention function and an RNN. At each time step t , the attention function produces an embedding \mathbf{c}_t of the input sequence, called a context vector. The subsequent RNN, characterized by a hidden state \mathbf{s}_t , computes from such an embedding a probability distribution over all possible output symbols, pointing to the most probable symbol y_t

$$p(y_t | y_1, \dots, y_{t-1}, \mathbf{x}) = \text{RNN}(\mathbf{c}_t). \quad (2)$$

The context vector is obtained as follows. At each time step t , the attention function takes as input the previous hidden state of the RNN \mathbf{s}_{t-1} and the annotations $\mathbf{h}_1, \dots, \mathbf{h}_T$. Such inputs are processed through an alignment model [see (3)] to obtain a set of scalar values e_{ti} that score the matching between the inputs around position i and the outputs around position t . These scores are then normalized through a softmax function, so as to obtain a set of weights a_{ti} [see (4)]

$$e_{ti} = f(\mathbf{s}_{t-1}, \mathbf{h}_i) \quad (3)$$

$$a_{ti} = \frac{\exp(e_{ti})}{\sum_{j=1}^T \exp(e_{tj})}. \quad (4)$$

Finally, the context vector \mathbf{c}_t is computed as a weighted sum of the annotations \mathbf{h}_i based on their weights a_{ti}

$$\mathbf{c}_t = \sum_i a_{ti} \mathbf{h}_i. \quad (5)$$

Quoting Bahdanau *et al.* [2], the use of attention “relieve[s] the encoder from the burden of having to encode all information in the source sentence into a fixed-length vector. With this new approach, the information can be spread throughout the sequence of annotations, which can be selectively retrieved by the decoder accordingly.”

B. Unified Attention Model

The characteristics of an attention model depend on the structure of the data whereupon they operate and on the desired output structure. The unified model we propose is based on

TABLE II
NOTATION

Symbol	Name	Definition
\mathbf{x}	Input sequence	Sequence of textual elements constituting the raw input.
\mathbf{K}	Keys	Matrix of d_k vectors (\mathbf{k}_i) of size n_k , whereupon attention weights are computed: $\mathbf{K} \in \mathbb{R}^{n_k \times d_k}$.
\mathbf{V}	Values	Matrix of d_k vectors (\mathbf{v}_i) of size n_v , whereupon attention is applied: $\mathbf{V} \in \mathbb{R}^{n_v \times d_k}$. Each \mathbf{v}_i and its corresponding \mathbf{k}_i offer two, possibly different, interpretations of the same entity.
\mathbf{q}	Query	Vector of size n_q , or sequence thereof, in which respect attention is computed: $\mathbf{q} \in \mathbb{R}^{n_q}$.
kaf vaf	qaf Annotation functions	Functions that encode the input sequence and query, producing \mathbf{K} , \mathbf{q} and \mathbf{V} respectively.
\mathbf{e}	Energy scores	Vector of size d_k , whose scalar elements (energy “scores”, e_i) represent the relevance of the corresponding \mathbf{k}_i , according to the compatibility function: $\mathbf{e} \in \mathbb{R}^{d_k}$.
\mathbf{a}	Attention weights	Vector of size d_k , whose scalar elements (attention “weights”, a_i) represent the relevance of the corresponding \mathbf{k}_i according to the attention model: $\mathbf{a} \in \mathbb{R}^{d_k}$.
f	Compatibility function	Function that evaluates the relevance of \mathbf{K} with respect to \mathbf{q} , returning a vector of energy scores: $\mathbf{e} = f(\mathbf{K}, \mathbf{q})$.
g	Distribution function	Function that computes the attention weights from the energy scores: $\mathbf{a} = g(\mathbf{e})$.
\mathbf{Z}	Weighted values	Matrix of d_k vectors (\mathbf{z}_i) of size n_v , representing the application of \mathbf{a} to \mathbf{V} : $\mathbf{Z} \in \mathbb{R}^{n_v \times d_k}$.
\mathbf{c}	Context vector	Vector of size n_v , offering a compact representation of \mathbf{Z} : $\mathbf{c} \in \mathbb{R}^{n_v}$.

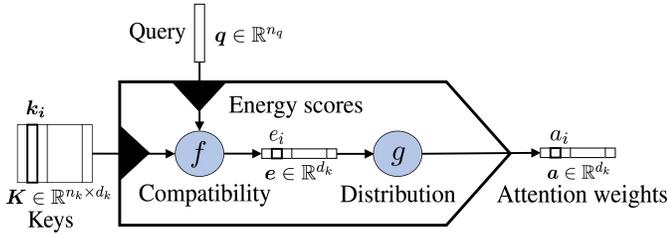


Fig. 3. Core attention model.

and extends the models proposed by Vaswani *et al.* [36] and Daniluk *et al.* [61]. It comprises a core part shared by almost the totality of the models found in the surveyed literature, as well as some additional components that, although not universally present, are still found in most literature models.

Fig. 3 shows the core attention model, which is part of the general model shown in Fig. 4. Table II lists the key terms and symbols. The core of the attention mechanism maps a sequence \mathbf{K} of d_k vectors \mathbf{k}_i , the keys, to a distribution \mathbf{a} of d_k weights a_i . \mathbf{K} encodes the data features whereupon attention is computed. For instance, \mathbf{K} may be word or character embeddings of a document, or the internal states of a recurrent architecture, as it happens with the annotation \mathbf{h}_i in RNNsearch. In some cases, \mathbf{K} could include multiple features or representations of the same object (e.g., both one-hot encoding and embedding of a word) or even—if the task calls for it—representations of entire documents.

More often than not, another input element \mathbf{q} , called query,¹ is used as a reference when computing the attention distribution. In that case, the attention mechanism will give emphasis to the input elements relevant to the task according to \mathbf{q} . If no query is defined, attention will give emphasis to the elements inherently relevant to the task at hand. In RNNsearch, for instance, \mathbf{q} is a single element, namely, the RNN hidden state \mathbf{s}_{t-1} . In other architectures, \mathbf{q} may represent different entities: embeddings of actual textual queries, contextual information,

background knowledge, and so on. It can also take the form of a matrix rather than a vector. For example, in their document attentive reader, Sordani *et al.* [67] made use of two query vectors.

From the keys and query, a vector \mathbf{e} of d_k energy scores e_i is computed through a compatibility function f

$$\mathbf{e} = f(\mathbf{q}, \mathbf{K}). \quad (6)$$

Function f corresponds to RNNsearch’s alignment model and to what other authors call energy function [43]. Energy scores are then transformed into attention weights using what we call a distribution function g

$$\mathbf{a} = g(\mathbf{e}). \quad (7)$$

Such weights are the outcome of the core attention mechanism. The commonest distribution function is the softmax function, as in RNNsearch, which normalizes all the scores to a probability distribution. Weights represent the relevance of each element to the given task, with respect to \mathbf{q} and \mathbf{K} .

The computation of these weights may already be sufficient for some tasks, such as the classification task addressed by Cui *et al.* [70]. Nevertheless, many tasks require the computation of new representation of the keys. In such cases, it is common to have another input element; a sequence \mathbf{V} of d_k vectors \mathbf{v}_i , the values, representing the data whereupon the attention computed from \mathbf{K} and \mathbf{q} is to be applied. Each element of \mathbf{V} corresponds to one and only one element of \mathbf{K} , and the two can be seen as different representations of the same data. Indeed, many architectures, including RNNsearch, do not distinguish between \mathbf{K} and \mathbf{V} . The distinction between keys and values was introduced by Daniluk *et al.* [61], who use different representations of the input for computing the attention distribution and the contextual information.

\mathbf{V} and \mathbf{a} are thus combined to obtain a new set \mathbf{Z} of weighted representations of \mathbf{V} [see (8)], which are then merged together so as to produce a compact representation of \mathbf{Z} usually called the context vector \mathbf{c} [see (9)].² The

¹The concept of “query” in attention models should not be confused with that used in tasks like question answering or information retrieval. In our model, the “query” is part of a general architecture and is task-independent.

²Although most authors use this terminology, we shall remark that Yang *et al.* [52], Wang *et al.* [129], and other authors used the term context vector to refer to other elements of the attention architecture.

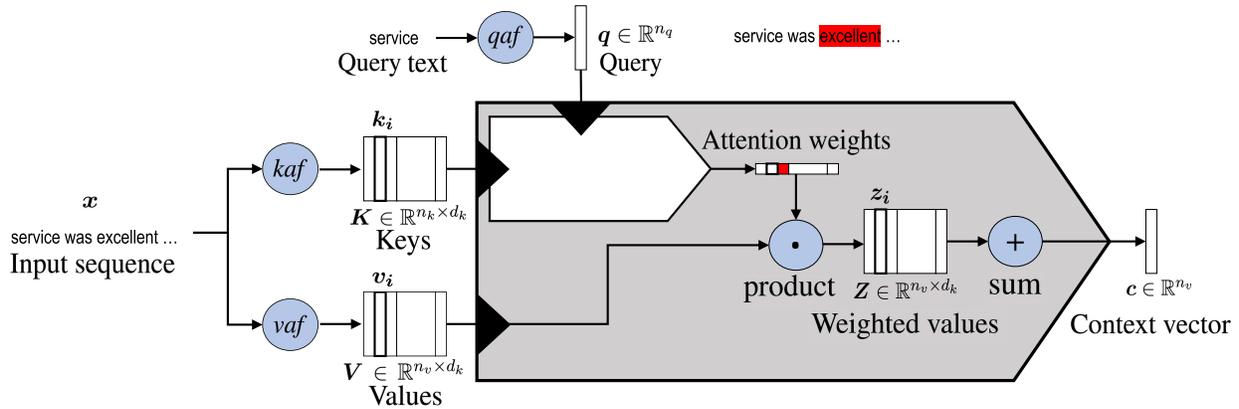


Fig. 4. General attention model.

commonest way of obtaining \mathbf{c} from \mathbf{Z} is by summation. However, alternatives have been proposed, including gating functions [93]. Nevertheless, \mathbf{c} will be mainly determined by values associated with higher attention weights

$$\mathbf{z}_i = a_i \mathbf{v}_i \quad (8)$$

$$\mathbf{c} = \sum_{i=1}^{d_k} \mathbf{z}_i. \quad (9)$$

What we described so far was a synthesis of the most frequent architectural choices made in the design of attentive architectures. Other options will be explored in Section IV-D.

C. Deterministic Versus Probabilistic Attention

Before we proceed, a brief remark about some relevant naming conventions is in order. The attention model described so far is sometimes described in the literature as a mapping with a probabilistic interpretation since the use of a softmax normalization allows one to interpret the attention weights as a probability distribution function (see [91]). Accordingly, some recent literature defines as deterministic attention models [134], [135] and those models where context words, whereupon attention is focused, are deterministically selected, for example, by using the constituency parse tree of the input sentence. However, other authors described the first model as deterministic (soft) attention, to contrast it with stochastic (hard) attention, where the probability distribution over weights is used to hardly sample a single input as context vector \mathbf{c} , in particular following a reinforcement learning strategy [16]. If the weights \mathbf{a} encode a probability distribution, the stochastic attention model differs from our general model for the absence of (8) and (9) that are replaced by the following equations:

$$\tilde{s} \sim \text{Multinoulli}(\{a_i\}) \quad (10)$$

$$\mathbf{c} = \sum_{i=1}^{d_k} \tilde{s}_i \mathbf{v}_i \quad \text{with } \tilde{s}_i \in \{0, 1\}. \quad (11)$$

To avoid any potential confusion, in the remainder of this article, we will abstain from characterizing the attention models as deterministic, probabilistic, or stochastic.

TABLE III

POSSIBLE USES OF ATTENTION AND EXAMPLES OF RELEVANT TASK

Use	Tasks
Feature selection	Multimodal tasks
Auxiliary task	Visual question answering Semantic role labelling
Contextual embedding creation	Machine Translation Sentiment Analysis Information Extraction
Sequence-to-sequence annotation	Machine Translation
Word selection	Dependency Parsing Cloze Question Answering
Multiple input processing	Question Answering

III. USES OF ATTENTION

Attention enables us to estimate the relevance of the input elements as well as to combine said elements into a compact representation—the context vector—that condenses the characteristics of the most relevant elements. Because the context vector is smaller than the original input, it requires fewer computational resources to be processed at later stages, yielding a computational gain.

We summarize possible uses of attention and the tasks in which they are relevant in Table III.

For tasks such as document classification, where usually there is only \mathbf{K} in input and no query, the attention mechanism can be seen as an instrument to encode the input into a compact form. The computation of such an embedding can be seen as a form of feature selection, and as such, it can be applied to any set of features sharing the same representation. This applies to cases where features come from different domains as in multimodal tasks [78] or from different levels of a neural architecture [38] or where they simply represent different aspects of the input document [136]. Similarly, attention can also be exploited as an auxiliary task during training so that specific features can be modeled via a multitask setting. This holds for several scenarios, such as visual question answering [137], domain classification for natural language understanding [138], and semantic role labeling [97].

When the generation of a text sequence is required, as in machine translation, attention enables us to make use of a dynamic representation of the input sequence, whereby the whole input does not have to be encoded into a single vector. At each time step, the encoding is tailored according to the task, and in particular, \mathbf{q} represents an embedding of the previous state of the decoder. More generally, the possibility to perform attention with respect to a query \mathbf{q} allows us to create representations of the input that depend on the task *context*, creating specialized embeddings. This is particularly useful in tasks, such as sentiment analysis and information extraction.

Since attention can create contextual representations of an element, it can also be used to build sequence-to-sequence annotators, without resorting to RNNs or convolutional neural networks (CNNs), as suggested by Vaswani *et al.* [36], who rely on an attention mechanism to obtain a whole encoder/decoder architecture.

Attention can also be used as a tool for selecting specific words. This could be the case, for example, in dependence parsing [97] and in cloze question-answering tasks [66], [70]. In the former case, attention can be applied to a sentence in order to predict dependences. In the latter, attention can be applied to a textual document or to a vocabulary to perform a classification among the words.

Finally, attention can come in handy when multiple interacting input sequences have to be considered in combination. In tasks such as question answering, where the input consists of two textual sequences—for instance, the question and the document or the question and the possible answers—an input encoding can be obtained by considering the mutual interactions between the elements of such sequences, rather than by applying a more rigid *a priori* defined model.

IV. TAXONOMY FOR ATTENTION MODELS

Attention models can be described on the basis of the following orthogonal dimensions: the nature of inputs (see Section IV-A), the compatibility function (see Section IV-B), the distribution function (see Section IV-C), and the number of distinct inputs/outputs, which we refer to as “multiplicity” (see Section IV-D). Moreover, attention modules can themselves be used inside larger attention models to obtain complex architectures such as hierarchical-input models (see Section IV-A2) or in some multiple-input coattention models (see Section IV-D2).

A. Input Representation

In NLP-related tasks, generally, \mathbf{K} and \mathbf{V} are representations of parts of documents, such as sequences of characters, words, or sentences. These components are usually embedded into continuous vector representations and then processed through key/value annotation functions (called *kaf* and *vaf* in Fig. 4), so as to obtain a hidden representation resulting in \mathbf{K} and \mathbf{V} . Typical annotation functions are RNN layers such as gated recurrent units (GRUs), long short-term memory networks (LSTMs), and CNNs. In this way, \mathbf{k}_i and \mathbf{v}_i represent an input element relative to its local context. If the layers in charge of

annotation are trained together with the attention model, they can learn to encode information useful to the attention model.

Alternatively, $\mathbf{k}_i/\mathbf{v}_i$ can be taken to represent each input element in isolation, rather than in context. For instance, they could be a one-hot encoding of words or characters or a pretrained word embedding. This results in an application of the attention mechanism directly to the raw inputs, which is a model known as inner attention [68]. Such a model has proven to be effective by several authors, who have exploited it in different fashions [36], [41], [54], [116]. The resulting architecture has a smaller number of layers and hyperparameters, which reduces the computational resources needed for training.

\mathbf{K} can also represent a single element of the input sequence. This is the case, for example, in the work by Bapna *et al.* [38], whose attention architecture operates on different encodings of the same element, obtained by the subsequent application of RNN layers. The context embeddings obtained for all components individually can then be concatenated, producing a new representation of the document that encodes the most relevant representation of each component for the given task. It would also be possible to aggregate each key or value with its neighbors, by computing their average or sum [128].

We have so far considered the input to be a sequence of characters, words, or sentences, which is usually the case. However, the input can also be other things, such as a juxtaposition of features or relevant aspects of the same textual element. For instance, Li *et al.* [56] and Zadeh *et al.* [78] considered the inputs composed of different sources, and in [136] and [139], the input represents different aspects of the same document. In that case, embeddings of the input can be collated together and fed into an attention model as multiple keys, as long as the embeddings share the same representation. This allows us to highlight the most relevant elements of the inputs and operate a feature selection, leading to a possible reduction of the dimensionality of the representation via the context embedding \mathbf{c} produced by the attention mechanism. Interestingly, Li *et al.* [120] proposed a model in which the textual input sequence is mixed with the output of the attention model. Their truncated history-attention model iterates the computation of attention on top of a bi-LSTM. At each step, the bi-LSTM hidden states are used as keys and values, while the context vectors computed in the previous steps act as a query.

We shall now explain in more detail two successful structures, which have become well-established building blocks of neural approaches for NLP, namely, self-attention and hierarchical-input architectures.

1) *Self-Attention*: We made a distinction between two input sources: the input sequence, represented by \mathbf{K} and \mathbf{V} , and the query, represented by \mathbf{q} . However, some architectures compute attention only based on the input sequence. These architectures are known as self-attentive or intraattentive models. We shall remark, however, that these terms are used to indicate many different approaches. The commonest one amounts to the application of multiple steps of attention to a vector \mathbf{K} , using the elements \mathbf{k}_t of the same vector as query at each step [18], [36]. At each step, the weights a_t^t

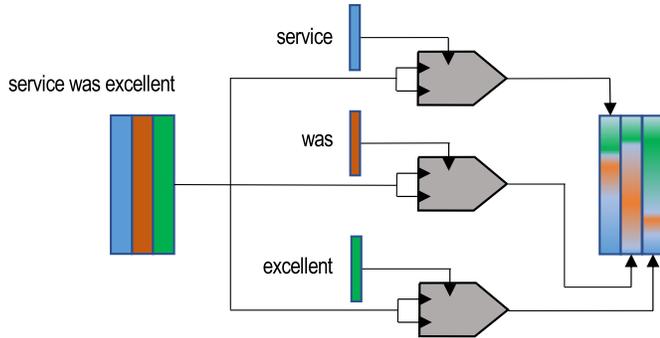


Fig. 5. Example of use of attention in a sequence-to-sequence model.

represent the relevance of k_i with respect to k_t , yielding d_K separate context embeddings, c^t , one per key. Attention could thus be used as a sequence-to-sequence model, as an alternative to CNNs or RNNs (see Fig. 5). In this way, each element of the new sequence may be influenced by elements of the whole input, incorporating contextual information without any locality boundaries. This is especially interesting since it could overcome a well-known shortcoming of RNNs: their limited ability of modeling long-range dependences [140]. For each element k_t , the resulting distribution of the weights a^t should give more emphasis to words that strongly relate to k_t . The analysis of these distributions will, therefore, provide information regarding the relationship between the elements inside the sequence. Modern text-sequence generation systems often rely on this approach. Another possibility is to construct a single query element q from the keys through a pooling operation. Furthermore, the same input sequence could be used both as keys K and query Q , applying a technique we will describe in Section IV-D2, known as coattention. Other self-attentive approaches, such as [52] and [100], are characterized by the complete absence of any query term q , which results in simplified compatibility functions (see Section IV-B).

2) *Hierarchical-Input Architectures*: In some tasks, portions of input data can be meaningfully grouped together into higher level structures, where hierarchical-input attention models can be exploited to subsequently apply multiple attention modules at different levels of the composition, as shown in Fig. 6.

Consider, for instance, data naturally associated with a two-level semantic structure, such as characters (the “micro” elements) forming words (the “macro” elements) or words forming sentences. Attention can be first applied to the representations of micro elements k_i , so as to build aggregate representations k_j of the macro elements, such as context vectors. Attention could then be applied again to the sequence of macroelement embeddings, in order to compute an embedding for the whole document D . With this model, attention first highlights the most relevant micro elements within each macro element and, then, the most relevant macro elements in the document. For instance, Yang *et al.* [52] applied attention first at word level, for each sentence in turn, to compute sentence embeddings. Then, they applied attention again on the sentence embeddings to obtain a document representation.

With reference to the model introduced in Section II, embeddings are computed for each sentence in D , and then, all such embeddings are used together as keys K to compute the document-level weights a and eventually D ’s context vector c . The hierarchy can be extended further. For instance, Wu *et al.* [141] added another layer on top, applying attention also at the document level.

If representations for both micro- and macro-level elements are available, one can compute attention on one level and then exploit the result as a key or query to compute attention on the other, yielding two different microrepresentation/macroe representation of D . In this way, attention enables us to identify the most relevant elements for the task at both levels. The attention-via-attention model by Zhao and Zhang [43] defines a hierarchy with characters at the micro level and words at the macro level. Both characters and words act as keys. Attention is first computed on word embeddings K_W , thus obtaining a document representation in the form of a context vector c_W , which in turn acts as a query q to guide the application of character-level attention to the keys (character embeddings) K_C , yielding a context vector c for D .

Ma *et al.* [113] identified a single “target” macro-object T as a set of words, which do not necessarily have to form a sequence in D , and then used such a macro-object as keys, K_T . The context vector c_T produced by a first application of the attention mechanism on K_T is then used as query q in a second application of the attention mechanism over D , with the keys being the document’s word embeddings K_W .

B. Compatibility Functions

The compatibility function is a crucial part of the attention architecture because it defines how keys and queries are matched or combined. In our presentation of compatibility functions, we will consider a data model where q and k_i are monodimensional vectors. For example, if K represents a document, each k_i may be the embedding of a sentence, a word, or a character. In such a model, q and k_i may have the same structure and, thus, the same size, although this is not always necessary. However, in some architectures, q can consist of a sequence of vectors or a matrix, a possibility we explore in Section IV-D2.

Some common compatibility functions are listed in Table IV. Two main approaches can be identified. The first one is to match and compare K and q . For instance, the idea behind the similarity attention proposed by Graves *et al.* [142] is that the most relevant keys are the most similar to the query. Accordingly, the authors present a model that relies on a similarity function (sim in Table IV) to compute the energy scores. For example, they rely on cosine similarity, a choice that suits cases where the query and the keys share the same semantic representation. A similar idea is followed by the widely used multiplicative or dot attention, where the dot product between q and K is computed. The complexity of this computation is $O(n_k d_k)$. A variation of this model is scaled multiplicative attention, where a scaling factor is introduced to improve performance with large keys [36]. General attention, proposed by Luong *et al.* [29],

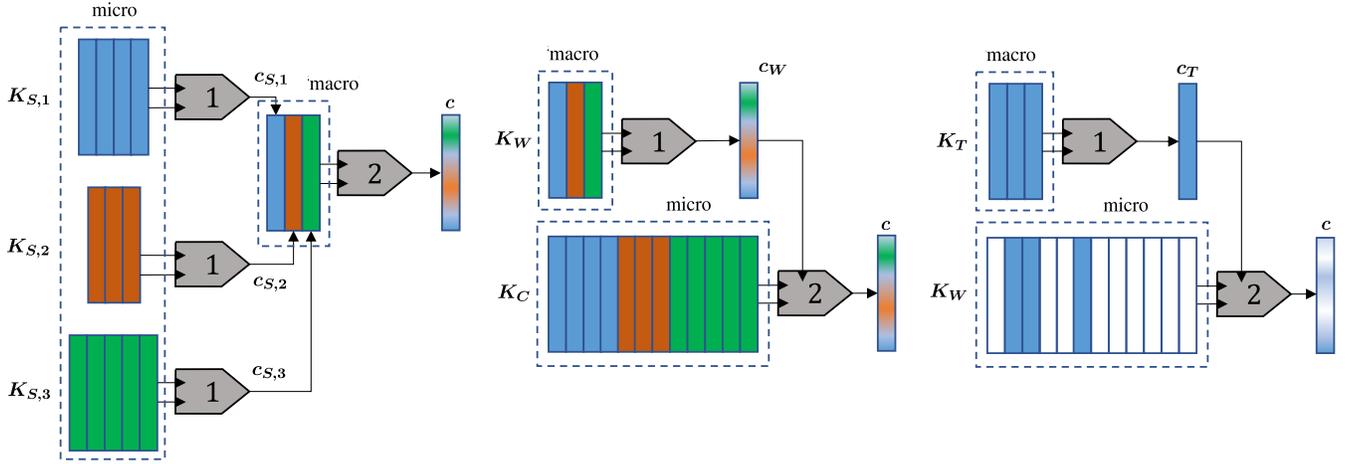


Fig. 6. Hierarchical input attention models defined by Zhao and Zhang [43] (center), Yang *et al.* [52] (left), and Ma *et al.* [113] (right). The number inside the attention shapes indicates the order of application. Different colors highlight different parts of the inputs.

TABLE IV

SUMMARY OF COMPATIBILITY FUNCTIONS FOUND IN THE LITERATURE. \mathbf{W} , \mathbf{W}_0 , \mathbf{W}_1, \dots , AND \mathbf{b} ARE LEARNABLE PARAMETERS

Name	Equation	Ref.
similarity	$f(\mathbf{q}, \mathbf{K}) = \text{sim}(\mathbf{q}, \mathbf{K})$	[142]
multiplicative or dot	$f(\mathbf{q}, \mathbf{K}) = \mathbf{q}^\top \mathbf{K}$	[29]
scaled multiplicative	$f(\mathbf{q}, \mathbf{K}) = \frac{\mathbf{q}^\top \mathbf{K}}{\sqrt{n_k}}$	[36]
general or bilinear	$f(\mathbf{q}, \mathbf{K}) = \mathbf{q}^\top \mathbf{W} \mathbf{K}$	[29]
biased general	$f(\mathbf{q}, \mathbf{K}) = \mathbf{K}^\top (\mathbf{W} \mathbf{q} + \mathbf{b})$	[67]
activated general	$f(\mathbf{q}, \mathbf{K}) = \text{act}(\mathbf{q}^\top \mathbf{W} \mathbf{K} + \mathbf{b})$	[111]
concat	$f(\mathbf{q}, \mathbf{K}) = \mathbf{w}_{imp}^\top \text{act}(\mathbf{W}[\mathbf{K}; \mathbf{q}] + \mathbf{b})$	[29]
additive	$f(\mathbf{q}, \mathbf{K}) = \mathbf{w}_{imp}^\top \text{act}(\mathbf{W}_1 \mathbf{K} + \mathbf{W}_2 \mathbf{q} + \mathbf{b})$	[2]
deep	$f(\mathbf{q}, \mathbf{K}) = \mathbf{w}_{imp}^\top \mathbf{E}^{(L-1)} + \mathbf{b}^L$ $\mathbf{E}^{(l)} = \text{act}(\mathbf{W}_l \mathbf{E}^{(l-1)} + \mathbf{b}^l)$ $\mathbf{E}^{(1)} = \text{act}(\mathbf{W}_1 \mathbf{K} + \mathbf{W}_0 \mathbf{q} + \mathbf{b}^1)$	[54]
convolution-based	$f(\mathbf{q}, \mathbf{K}) = [e_0; \dots; e_{d_k}]$ $e_j = \frac{1}{l} \sum_{i=j-l}^j e_{j,i}$ $e_{j,i} = \text{act}(\mathbf{w}_{imp}^\top [k_i; \dots; k_{i+l}] + \mathbf{b})$	[119]
location-based	$f(\mathbf{q}, \mathbf{K}) = f(\mathbf{q})$	[29]

extends this concept in order to accommodate keys and queries with different representations. To that end, it introduces a learnable matrix parameter \mathbf{W} . This parameter represents the transformation matrix that maps the query onto the vector

space of keys. This transformation increases the complexity of the operation to $O(n_q n_k d_k)$. In what could be called a biased general attention, Sordoni *et al.* [67] introduced a learnable bias, so as to consider some keys as relevant independently of the input. Activated general attention [111] employs a nonlinear activation function. In Table IV, *act* is a placeholder for a nonlinear activation function, such as hyperbolic tangent, tanh, rectifier linear unit, ReLU [143], or scaled exponential linear unit, SELU [144]. These approaches are particularly suitable in tasks where the concept of relevance of a key is known to be closely related to that of similarity to a query element. These include, for instance, tasks where specific keywords can be used as a query, such as abusive speech recognition and sentiment analysis.

A different approach amounts to combining rather than comparing \mathbf{K} and \mathbf{q} , using them together to compute a joint representation, which is then multiplied by an importance vector³ \mathbf{w}_{imp} , which has to adhere to the same semantic of the new representation. Such a vector defines, in a way, relevance and could be an additional query element, as offered by Ma *et al.* [113] or a learnable parameter. In that case, we speculate that the analysis of a machine-learned importance vector could provide additional information on the model. One of the simplest models that follows this approach is the concat attention by Luong *et al.* [29], where a joint representation is given by juxtaposing keys and queries. Additive attention works similarly, except that the contribution of \mathbf{q} and \mathbf{K} can be computed separately. For example, Bahdanau *et al.* [2] precomputed the contribution of \mathbf{K} in order to reduce the computational footprint. The complexity of the computation, ignoring the application of the nonlinear function, is thus $O(n_w n_k d_k)$, where n_w indicates the size of \mathbf{w}_{imp} . Moreover, additive attention in principle could accommodate queries of different size. In additive attention and concat attention, the keys and the queries are fed into a single neural layer. We speak instead of deep attention if multiple layers are

³Part of our terminology. As previously noted, *wimp* is termed context vector by Yang *et al.* [52] and other authors.

employed [54]. Table IV shows a deep attention function with L levels of depth, $1 < l < L$. With an equal number of neurons for all levels and reasonably assuming L to be much smaller than any other parameter, the complexity becomes $O(n_w^L n_k d_k)$. These approaches are especially suitable when a representation of “relevant” elements is unavailable or it is available but encoded in a significantly different way from the way that keys are encoded. This may be the case, for instance, with tasks such as document classification and summarization.

A similar rationale lead to convolution-based attention [119]. It draws inspiration from biological models, whereby biological attention is produced through search templates and pattern matching. Here, the learned vector w_{imp} represents a convolutional filter, which embodies a specific relevance template. The filter is used in a pattern-matching process obtained by applying the convolution operation on key subsequences. Applying multiple convolutions with different filters enables a contextual selection of keys that match specific relevance templates, obtaining a specific energy score for each filter. Since each key belongs to multiple subsequences, each key yields multiple energy scores. Such scores are then aggregated in order to obtain a single value per key. Aggregation could be achieved by sum or average. Table IV illustrates convolution-based attention for a single filter of size l . The complexity of such an operation is $O(l^2 n_k d_k)$. These approaches are especially suitable when a representation of “relevant” elements is unavailable or it is available but encoded in a significantly different way from the way that keys are encoded. This may be the case, for instance, with tasks such as document classification and summarization.

Finally, in some models, the attention distribution ignores \mathbf{K} and only depends on \mathbf{q} . We then speak of location-based attention. The energy associated with each key is thus computed as a function of the key’s position, independently of its content [29]. Conversely, self-attention may be computed only based on \mathbf{K} , without any \mathbf{q} . The compatibility functions for self-attention, which are a special case of the more general functions, are omitted from Table IV.

For an empirical comparison between some compatibility functions (namely, additive and multiplicative attention) in the domain of machine translation, we suggest the reader to refer to Britz *et al.* [37].

C. Distribution Functions

Attention distribution maps energy scores to attention weights. The choice of the distribution function depends on the properties that the distribution is required to have—for instance, whether it is required to be a probability distribution, a set of probability scores, or a set of Boolean scores—on the need to enforce sparsity, and on the need to account for the keys’ positions.

One possible distribution function g is the logistic sigmoid, as proposed by Kim *et al.* [47]. In this way, each weight a_i is constrained between 0 and 1, thus ensuring that the values \mathbf{V}_i and their corresponding weighted counterparts \mathbf{Z}_i share the same boundaries. Such weights can thus be interpreted as probabilities that an element is relevant. The same range can

also be forced on the context vector’s elements c_i , by using a softmax function, as it is commonly done. In that case, the attention mechanism is called soft attention. Each attention weight can be interpreted as the probability that the corresponding element is the most relevant.

With sigmoid or softmax alike, all the key/value elements have a relevance, small as it may be. Yet, it can be argued that in some cases, some parts of the input are completely irrelevant, and if they were to be considered, they would likely introduce noise rather than contribute with useful information. In such cases, one could exploit attention distributions that altogether ignore some of the keys, thereby reducing the computational footprint. One option is the sparsemax distribution [91], which truncates to zero the scores under a certain threshold by exploiting the geometric properties of the probability simplex. This approach could be especially useful in those settings where a large number of elements are irrelevant, such as in document summarization or cloze question-answering tasks.

It is also possible to model structural dependences between the outputs. For example, structured attention networks [47] exploit neural conditional random fields to model the (conditional) attention distribution. The attention weights are thus computed as (normalized) marginals from the energy scores, which are treated as potentials.

In some tasks, such as machine translation or image captioning, the relevant features are found in a neighborhood of a certain position. In those cases, it could be helpful to focus the attention only on a specific portion of the input. If the position is known in advance, one can apply a positional mask, by adding or subtracting a given value from the energy scores before the application of the softmax [93]. Since the location may not be known in advance, the hard attention model by Xu *et al.* [16] considers the keys in a dynamically determined location. Such a solution is less expensive at inference time, but it is not differentiable. For that reason, it requires more advanced training techniques, such as reinforcement learning or variance reduction. Local attention [29] extends this idea while preserving differentiability. Guided by the intuition that in machine translation at each time step, only a small segment of the input can be considered relevant, and local attention considers only a small window of the keys at a time. The window has a fixed size, and the attention can be better focused on a precise location by combining the softmax distribution with a Gaussian distribution. The mean of the Gaussian distribution is dynamic, whereas its variance can either be fixed, as done by Luong *et al.* [29] or dynamic, as done by Yang *et al.* [39]. Selective attention [17] follows the same idea; using a grid of Gaussian filters, it considers only a patch of the keys, whose position, size, and resolution depend on dynamic parameters.

Shen *et al.* [94] combined soft and hard attention, by applying the former only on the elements filtered by the latter. More precisely, softmax is only applied among a subset of selected energy scores, whereas for the others, the weight is set to zero. The subset is determined according to a set of random variables, with each variable corresponding to a key. The probability associated with each variable is determined

through soft attention applied to the same set of keys. The proper “softness” of the distribution could depend not only on the task but also on the query. Lin *et al.* [44] defined a model whose distribution is controlled by a learnable, adaptive temperature parameter. When a “softer” attention is required, the temperature increases, producing a smoother distribution of weights. The opposite happens when a “harder” attention is needed.

Finally, the concept of locality can also be defined according to semantic rules, rather than the temporal position. This possibility will be further discussed in Section V.

D. Multiplicity

We shall now present variations of the general unified model where the attention mechanism is extended to accommodate multiple, possibly heterogeneous, inputs or outputs.

1) *Multiple Outputs*: Some applications suggest that the data could, and should, be interpreted in multiple ways. This can be the case when there is ambiguity in the data, stemming, for example, from words having multiple meanings or when addressing a multitask problem. For this reason, models have been defined that jointly compute not only one but multiple attention distributions over the same data. One possibility presented by Lin *et al.* [100] and also exploited by Du *et al.* [145] is to use additive attention (seen in Section IV-B) with an importance matrix, instead of a vector, $\mathbf{W}_{imp} \in \mathbb{R}^{n_k \times n_o}$, yielding an energy scores matrix where multiple scores are associated with each key. Such scores can be regarded as different models of relevance for the same values and can be used to create a context matrix $\mathbf{C} \in \mathbb{R}^{n_o \times n_o}$. Such embeddings can be concatenated together, creating a richer and more expressive representation of the values. Regularization penalties can be applied so as to enforce the differentiation between models of relevance (for example, the Frobenius norm). In multidimensional attention [93], where the importance matrix is a square matrix, attention can be computed featurewise. To that end, each weight $a_{i,j}$ is paired with the j th feature of the i th value $v_{i,j}$, and a featurewise product yields the new value z_i . Convolution-based attention [119] always produces multiple energy scores distributions according to the number of convolutional filters and the size of those filters. Another possibility explored by Vaswani *et al.* [36] is multihead attention. There, multiple linear projections of all the inputs (\mathbf{K} , \mathbf{V} , and \mathbf{q}) are performed according to learnable parameters, and multiple attention functions are computed in parallel. Usually, the processed context vectors are then merged together into a single embedding. A suitable regularization term is sometimes imposed so as to guarantee sufficient dissimilarity between attention elements. Li *et al.* [34] proposed three possibilities: regularization on the subspaces (the linear projections of \mathbf{V}), the attended positions (the sets of weights), or on the outputs (the context vectors). Multihead attention can be especially helpful when combined with nonsoft attention distribution since different heads can capture local and global contexts at the same time [39]. Finally, labelwise attention [126] computes a separate attention distribution for each class. This may improve the performance

as well as lead to a better interpretation of the data because it could help isolate data points that better describe each class. These techniques are not mutually exclusive. For example, multihead and multidimensional attention can be combined with one another [95].

2) *Multiple Inputs: Coattention*: Some architectures consider the query to be a sequence of d_q multidimensional elements, represented by a matrix $\mathbf{Q} \in \mathbb{R}^{n_q \times d_q}$, rather than by a plain vector. Examples of this setup are common in architectures designed for tasks where the query is a sentence, as in question answering, or a set of keywords, as in abusive speech detection. In those cases, it could be useful to find the most relevant query elements according to the task and the keys. A straightforward way of doing that would be to apply the attention mechanism to the query elements, thus treating \mathbf{Q} as keys and each \mathbf{k}_i as query, yielding two independent representations for \mathbf{K} and \mathbf{Q} . However, in that way, we would miss the information contained in the interactions between elements of \mathbf{K} and \mathbf{Q} . Alternatively, one could apply attention jointly on \mathbf{K} and \mathbf{Q} , which become the “inputs” of a coattention architecture [80]. Coattention models can be coarse-grained or fine-grained [112]. Coarse-grained models compute attention on each input, using an embedding of the other input as a query. Fine-grained models consider each element of an input with respect to each element of the other input. Furthermore, coattention can be performed sequentially or in parallel. In parallel models, the procedures to compute attention on \mathbf{K} and \mathbf{Q} symmetric, and thus, the two inputs are treated identically.

a) *Coarse-grained coattention*: Coarse-grained models use a compact representation of one input to compute attention on the other. In such models, the role of the inputs as keys and queries is no longer focal, and thus, a compact representation of \mathbf{K} may play as a query in parts of the architecture and vice versa. A sequential coarse-grained model proposed by Lu *et al.* [80] is alternating coattention, as shown in Fig. 7 (left), whereby attention is computed three times to obtain embeddings for \mathbf{K} and \mathbf{Q} . First, self-attention is computed on \mathbf{Q} . The resulting context vector is then used as a query to perform attention on \mathbf{K} . The result is another context vector \mathbf{C}_K , which is further used as a query as attention is again applied to \mathbf{Q} . This produces a final context vector, \mathbf{C}_Q . The architecture proposed by Sordoni *et al.* [67] can also be described using this model with a few adaptations. In particular, Sordoni *et al.* [67] omitted the last step and factor in an additional query element \mathbf{q} in the first two attention steps. An almost identical sequential architecture is used by Zhang *et al.* [86], who use \mathbf{q} only in the first attention step. A parallel coarse-grained model is shown in Fig. 7 (right). In such a model, proposed by Ma *et al.* [111], an average (*avg*) is initially computed on each input and then used as a query in the application of attention to generate the embedding of the other input. Sequential coattention offers a more elaborate computation of the final results, potentially allowing to discard all the irrelevant elements of \mathbf{Q} and \mathbf{K} , at the cost of a greater computational footprint. Parallel coattention can be optimized for better performance, at the expense of a “simpler” elaboration of the outputs. It is worthwhile noticing that the

averaging step in Ma *et al.* [111]’s model could be replaced by self-attention, in order to filter out irrelevant elements from Q and K at an early stage.

b) Fine-grained coattention: In fine-grained coattention models, the relevance (energy scores) associated with each key/query element pair $(\mathbf{k}_i/\mathbf{q}_j)$ is represented by the elements $E_{j,i}$ of a coattention matrix $\mathbf{E} \in \mathbb{R}^{d_q \times d_k}$ computed by a cocompatibility function.

Cocompatibility functions can be straightforward adaptations of any of the compatibility functions listed in Table IV. Alternatively, new functions can be defined. For example, Fan *et al.* [112] defined cocompatibility as a linear transformation of the concatenation between the elements and their product [see (12)]. In decomposable attention [89], the inputs are fed into neural networks, whose outputs are then multiplied [see (13)]. Delaying the product to after the processing by the neural networks reduces the number of inputs of such networks, yielding a reduction in the computational footprint. An alternative proposed by Tay *et al.* [75] exploits the Hermitian inner product. The elements of \mathbf{K} and \mathbf{Q} are first projected to a complex domain, then, the Hermitian product between them is computed, and finally, only the real part of the result is kept. Being the Hermitian product noncommutative, \mathbf{E} will depend on the roles played by the inputs as keys and queries

$$E_{j,i} = \mathbf{W}([\mathbf{k}_i; \mathbf{q}_j; \mathbf{k}_i \mathbf{q}_j]) \quad (12)$$

$$E_{j,i} = (\text{act}(\mathbf{W}_1 \mathbf{Q} + \mathbf{b}_1))^\top (\text{act}(\mathbf{W}_2 \mathbf{K} + \mathbf{b}_2)). \quad (13)$$

Because $E_{j,i}$ represent energy scores associated with $(\mathbf{k}_i/\mathbf{q}_j)$ pairs, computing the relevance of \mathbf{K} with respect to specific query elements, or, similarly, the relevance of \mathbf{Q} with respect to specific key elements, requires extracting information from \mathbf{E} using what we call an aggregation function. The output of such a function is a pair $\mathbf{a}_K/\mathbf{a}_Q$ of weight vectors.

The commonest aggregation functions are listed in Table V. A simple idea is the attention pooling parallel model adopted by dos Santos *et al.* [69]. It amounts to considering the highest score in each row or column of \mathbf{E} . By attention pooling, a key \mathbf{k}_i will be attributed a high attention weight if and only if it has a high coattention score with respect to at least one query element \mathbf{q}_j . Key attention scores are obtained through rowwise max pooling, whereas query attention scores are obtained through columnwise max pooling, as shown in Fig. 8 (left). Only considering the highest attention scores may be regarded as a conservative approach. Indeed, low-energy scores can only be obtained by keys whose coattention scores are all low and thus likely to be irrelevant to all query elements. Furthermore, the keys that are only relevant to a single query element may obtain the same or even higher energy score than keys that are relevant to all the query elements. The same reasoning applies to query attention scores. This is a suitable approach, for example, in tasks where the queries are keywords in disjunction, such as in abusive speech recognition. Conversely, the approaches follow to compute the energy score of an element by accounting for all the related attention scores, at the cost of a heavier computational footprint.

TABLE V

AGGREGATION FUNCTIONS. IN MOST CASES, \mathbf{a}_K AND \mathbf{a}_Q ARE OBTAINED BY APPLYING A DISTRIBUTION FUNCTION, SUCH AS THOSE SEEN IN SECTION IV-C, TO \mathbf{e}_K AND \mathbf{e}_Q , AND ARE THUS OMITTED FROM THIS TABLE IN THE INTEREST OF BREVITY. AS CUSTOMARY, *Act* IS A PLACEHOLDER FOR A GENERIC NONLINEAR ACTIVATION FUNCTION, WHEREAS *Dist* INDICATES A DISTRIBUTION FUNCTION SUCH AS SOFTMAX

Name	Equations	Ref.
<i>pooling</i>	$e_{K_i} = \max_{1 \leq j \leq d_q} (E_{j,i})$ $e_{Q_j} = \max_{1 \leq i \leq d_k} (E_{j,i})$	[69]
<i>perceptron</i>	$\mathbf{e}_K = \mathbf{W}_3^\top \text{act}(\mathbf{W}_1 \mathbf{K} + \mathbf{W}_2 \mathbf{Q} \mathbf{E})$ $\mathbf{e}_Q = \mathbf{W}_4^\top \text{act}(\mathbf{W}_2 \mathbf{K} + \mathbf{W}_1 \mathbf{Q} \mathbf{E}^\top)$	[80]
<i>linear transformation</i>	$\mathbf{e}_K = \mathbf{W}_1 \mathbf{E}$ $\mathbf{e}_Q = \mathbf{W}_2 \mathbf{E}$	[127]
<i>attention over attention</i>	$\mathbf{a}_K = \mathbf{M}_1 \cdot \mathbf{a}_Q$ $\mathbf{a}_Q = \text{average}_{1 \leq i \leq d_k} (\mathbf{M}_2, i)$ $\mathbf{M}_2, i = \text{dist}_{1 \leq j \leq d_q} (E_{j,i})$ $\mathbf{M}_1, j = \text{dist}_{1 \leq i \leq d_k} (E_{j,i})$	[70]
<i>perceptron with nested attention</i>	$\mathbf{e}_K = \mathbf{W}_3^\top \text{act}(\mathbf{W}_1 \mathbf{K} + (\mathbf{W}_2 \mathbf{Q}^\top) \mathbf{M}_2)$ $\mathbf{e}_Q = \mathbf{W}_4^\top \text{act}(\mathbf{W}_2 \mathbf{Q} + (\mathbf{W}_1 \mathbf{K}^\top) \mathbf{M}_1^\top)$ $\mathbf{M}_2, i = \text{dist}_{1 \leq j \leq d_q} (E_{j,i})$ $\mathbf{M}_1, j = \text{dist}_{1 \leq i \leq d_k} (E_{j,i})$	[88]

Lu *et al.* [80] used a multilayer perceptron in order to learn the mappings from \mathbf{E} to \mathbf{e}_K and \mathbf{e}_Q . In [127], the computation is even simpler since the final energy scores are a linear transformation of \mathbf{E} . Cui *et al.* [70] instead applied the nested model shown in Fig. 8 (right). First, two matrices \mathbf{M}_1 and \mathbf{M}_2 are computed by separately applying a rowwise and a columnwise softmax on \mathbf{E} . The idea is that each row of \mathbf{M}_1 represents the attention distribution over the document according to a specific query element—and it could already be used as such. Then, a rowwise average over \mathbf{M}_2 is computed so as to produce an attention distribution \mathbf{a}_Q over query elements. Finally, a weighted sum of \mathbf{M}_1 according to the relevance of query elements is computed through the dot product between \mathbf{M}_1 and \mathbf{a}_Q , obtaining the document’s attention distribution over the keys, \mathbf{a}_K . An alternative nested attention model is proposed by Nie *et al.* [88], whereby \mathbf{M}_1 and \mathbf{M}_2 are fed to a multilayer perceptron, as is done by Lu *et al.* [80]. Further improvements may be obtained by combining the results of multiple coattention models. Fan *et al.* [112], for instance, computed coarse-grained and fine-grained attention in parallel and combined the results into a single embedding.

V. COMBINING ATTENTION AND KNOWLEDGE

According to LeCun *et al.* [146], a major open challenge in artificial intelligence (AI) is combining connectionist (or subsymbolic) models, such as deep networks, with approaches based on symbolic knowledge representation, in order to

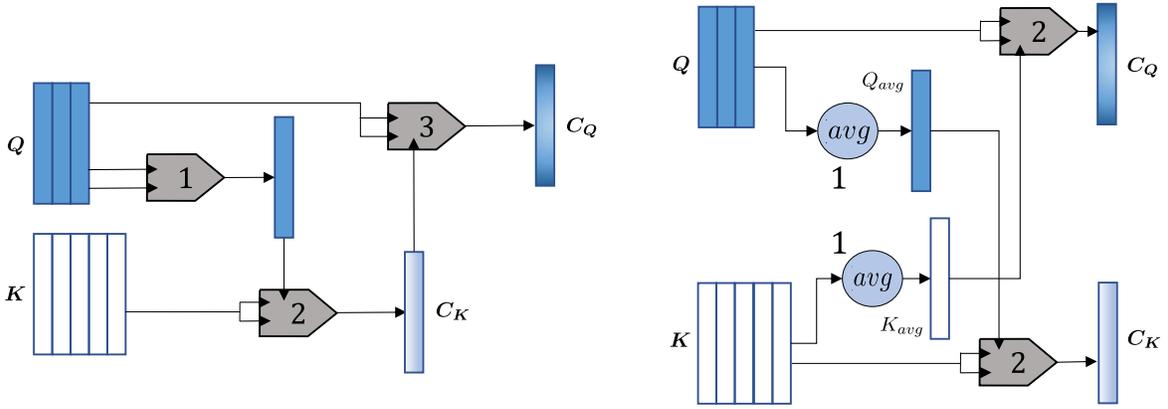


Fig. 7. Coarse-grained coattention by Lu *et al.* [80] (left) and Ma *et al.* [111] (right). The number inside the attention shapes (and besides the average operator) indicates the order in which they are applied. Different colors highlight different inputs.

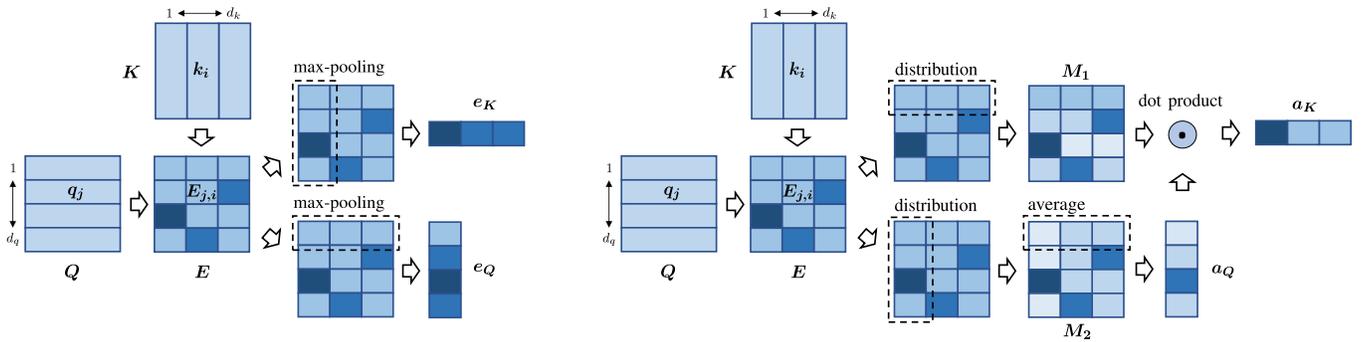


Fig. 8. Fine-grained coattention models presented by dos Santos *et al.* [69] (left) and by Cui *et al.* [70] (right). Dashed lines show how max-pooling/distribution functions are applied (columnwise or rowwise).

perform complex reasoning tasks. Throughout the last decade, filling the gap between these two families of AI methodologies has represented a major research avenue. Popular approaches include statistical relational learning [147], neural-symbolic learning [148], and the application of various deep learning architectures [149], such as memory networks [59], neural Turing machines [142], and several others.

From this perspective, attention can be seen both as an attempt to improve the interpretability of neural networks and as an opportunity to plug external knowledge into them. As a matter of fact, since the weights assigned by attention represent the relevance of the input with respect to the given task, in some contexts, it could be possible to exploit this information to isolate the most significant features that allow the deep network to make its predictions. On the other hand, any prior knowledge regarding the data, the domain, or the specific task, whenever available, could be exploited to generate information about the desired attention distribution, which could be encoded within the neural architecture.

In this section, we overview different techniques that can be used to inject this kind of knowledge in a neural network. We leave to Section VI further discussions on the open challenges regarding the combination of knowledge and attention.

A. Supervised Attention

In most of the works we surveyed, the attention model is trained with the rest of the neural architecture to perform a

specific task. Although trained alongside a supervised procedure, the attention model per se is trained in an unsupervised fashion⁴ to select useful information for the rest of the architecture. Nevertheless, in some cases, knowledge about the desired weight distribution could be available. Whether it is present in the data as a label or it is obtained as additional information through external tools, it can be exploited to perform a supervised training of the attention model.

1) *Preliminary Training*: One possibility is to use an external classifier. The weights learned by such a classifier are subsequently plugged into the attention model of a different architecture. We name this procedure as preliminary training. For example, Zhang *et al.* [53] first trained an attention model to represent the probability that a sentence contains relevant information. The relevance of a sentence is given by rationales [150], which are snippets of text that supports the corresponding document categorizations. In work by Long *et al.* [118], a model is preliminarily trained on eye-tracking data sets to estimate the reading time of words. Then, the reading time predicted by the model is used as an energy score in a neural model for sentiment analysis.

2) *Auxiliary Training*: Another possibility is to train the attention model without preliminary training, but by treating attention learning as an auxiliary task that is performed jointly with the main task. This procedure has led to good results

⁴Meaning that there is no target distribution for the attention model.

in many scenarios, including machine translation [30], [35], visual question answering [137], and domain classification for natural language understanding [138].

In some cases, this mechanism can also be exploited to have attention model-specific features. For example, since the linguistic information is useful for semantic role labeling, attention can be trained in a multitask setting to represent the syntactic structure of a sentence. Indeed, in LISA [97], a multilayer multiheaded architecture for semantic role labeling, one of the attention heads is trained to perform dependence parsing as an auxiliary task.

3) *Transfer Learning*: Furthermore, it is possible to perform transfer learning across different domains [1] or tasks [115]. By performing a preliminary training of an attentive architecture on a source domain to perform a source task, a mapping between the inputs and the distribution of weights will be learned. Then, when another attentive architecture is trained on the target domain to perform the target task, the pretrained model can be exploited. Indeed, the desired distribution can be obtained through the first architecture. Attention learning can, therefore, be treated as an auxiliary task as in the previously mentioned cases. The difference is that the distribution of the pretrained model is used as ground truth, instead of using data labels.

B. Attention Tracking

When attention is applied multiple times on the same data, as in sequence-to-sequence models, a useful piece of information could be how much relevance has been given to the input along different model iterations. Indeed, one may need to keep track of the weights that the attention model assigns to each input. For example, in machine translation, it is desirable to ensure that all the words of the original phrase are considered. One possibility to maintain this information is to use a suitable structure and provide it as an additional input to the attention model. Tu *et al.* [33] exploited a piece of symbolic information called coverage to keep track of the weight associated with the inputs. Every time attention is computed, such information is fed to the attention model as a query element, and it is updated according to the output of the attention itself. In [31], the representation is enhanced by making use of a subsymbolic representation for the coverage.

C. Modeling the Distribution Function by Exploiting Prior Knowledge

Another component of the attention model where prior knowledge can be exploited is the distribution function. For example, constraints can be applied to the computation of the new weights to enforce the boundaries on the weights assigned to the inputs. In [46] and [51], the coverage information is exploited by a constrained distribution function, regulating the amount of attention that the same word receives over time.

Prior knowledge could also be exploited to define or to infer a distance between the elements in the domain. Such domain-specific distance could then be considered in any position-based distribution function, instead of the positional distance. An example of distance could be derived by the

syntactical information. Chen *et al.* [40] and He *et al.* [109] used distribution functions that consider the distance between two words along the dependence graph of a sentence.

VI. CHALLENGES AND FUTURE DIRECTIONS

In this section, we discuss open challenges and possible applications of the attention mechanism in the analysis of neural networks, as a support of the training process and as an enabling tool for the integration of symbolic representations within neural architectures.

A. Attention for Deep Networks Investigation

Whether attention may or may not be considered as a mean to explain neural networks is currently an open debate. Some recent studies [10], [11] suggest that attention cannot be considered a reliable mean to explain or even interpret neural networks. Nonetheless, other works [6]–[9] advocate the use of attention weights as an analytic tool. Specifically, Jain and Wallace [10] proved that attention is not consistent with other explainability metrics and that it is easy to create local adversarial distributions (distributions that are similar to the trained model but produce a different outcome). Wiegrefe and Pinter [9] pushed the discussion further, providing experiments that demonstrate that creating an effective global adversarial attention model is much more difficult than creating a local one and that attention weights may contain information regarding feature importance. Their conclusion is that attention may indeed provide an explanation of a model, if by explanation, we mean a plausible, but not necessarily faithful, reconstruction of the decision-making process, as suggested by Rudin [151] and Riedl [152].

In the context of a multilayer neural architecture, it is fair to assume that the deepest levels will compute the most abstract features [146], [153]. Therefore, the application of attention to deep networks could enable the selection of higher level features, thus providing hints to understand which complex features are relevant for a given task. Following this line of inquiry in the computer vision domain, Zhang *et al.* [18] showed that the application of attention to middle-to-high level feature sets leads to better performance in image generation. The visualization of the self-attention weights has revealed that higher weights are not attributed to proximate image regions, but rather to those regions whose color or texture is most similar to that of the query image point. Moreover, the spatial distribution does not follow a specific pattern, but instead, it changes, modeling a region that corresponds to the object depicted in the image. Identifying abstract features in an input text might be less immediate than in an image, where the analysis process is greatly aided by visual intuition. Yet, it may be interesting to test the effects of the application of attention at different levels and to assess whether its weights correspond to specific high-level features. For example, Vaswani *et al.* [36] analyzed the possible relation between attention weights and syntactic predictions, Voita *et al.* [49] did the same with anaphora resolutions, and Clark *et al.* [6] investigated the correlation with many linguistic features. Voita *et al.* [50] analyzed the behavior of the heads of a multihead model,

discovering that different heads develop different behaviors, which can be related to specific position or specific syntactical element. Yang *et al.* [39] seemed to confirm that the deeper levels of neural architectures capture nonlocal aspects of the textual input. They studied the application of locality at different depths of an attentive deep architecture and showed that its introduction is especially beneficial when it is applied to the layers that are closer to the inputs. Moreover, when the application of locality is based on a variable-size window, higher layers tend to have a broader window.

A popular way of investigating whether an architecture has learned high-level features amounts to using the same architecture to perform other tasks, as it happens with transfer learning. This setting has been adopted outside the context of attention, for example, by Shi *et al.* [154], who perform syntactic predictions by using the hidden representations learned with machine translation. In a similar way, attention weights could be used as input features in a different model, so as to assess whether they can select relevant information for a different learning task. This is what happens, for example, in attention distillation, where a student network is trained to penalize the most confusing features according to a teacher network, producing an efficient and robust model in the task of machine reading comprehension [155]. Similarly, in a transfer learning setting, attentional heterogeneous transfer [156] has been exploited in heterolingual text classification to selectively filter input features from heterogeneous sources.

B. Attention for Outlier Detection and Sample Weighing

Another possible use of attention may be for outlier detection. In tasks such as classification or the creation of a representative embedding of a specific class, attention could be applied over all the samples belonging to that task. In doing so, the samples associated with small weights could be regarded as outliers with respect to their class. The same principle could be potentially applied to each data point in a training set, independently of its class. The computation of a weight for each sample could be interpreted as assessing the relevance of that specific data point for a specific task. In principle, assigning such samples a low weight and excluding them from the learning could improve a model's robustness to noisy input. Moreover, a dynamic computation of these weights during training would result in a dynamic selection of different training data in different training phases. While attention-less adaptive data selection strategies have already proven to be useful for efficiently obtaining more effective models [117], to the best of our knowledge, no attention-based approach has been experimented to date.

C. Attention Analysis for Model Evaluation

The impact of attention is greatest when all the irrelevant elements are excluded from the input sequence, and the importance of the relevant elements is properly balanced. A seemingly uniform distribution of the attention weights could be interpreted as a sign that the attention model has been unable to identify the more useful elements. This, in turn, may be due to the data that do not contain useful information for

the task at hand or it may be ascribed to the poor ability of the model to discriminate information. Nevertheless, the attention model would be unable to find relevant information in the specific input sequence, which may lead to errors. The analysis of the distribution of the attention weights may, therefore, be a tool for measuring an architecture's confidence in performing a task on a given input. We speculate that a high entropy in the distribution or the presence of weights above a certain threshold may be correlated with a higher probability of success of the neural model. These may, therefore, be used as indicators, to assess the uncertainty of the architecture, as well as to improve its interpretability. Clearly, this information would be useful to the user, to better understand the model and the data, but it may also be exploited by more complex systems. For example, Heo *et al.* [157] proposed to exploit the uncertainty of their stochastic predictive model to avoid making risky predictions in healthcare tasks.

In the context of an architecture that relies on multiple strategies to perform its task, such as a hybrid model that relies on both symbolic and subsymbolic information, the uncertainty of the neural model can be used as a parameter in the merging strategy. Other contexts in which this information may be relevant are multitask learning and reinforcement learning. Examples of exploitation of the uncertainty of the model, although in contexts other than attention and NLP, can be found in works by Poggi and Mattocchia [158], Kendall *et al.* [159], and Blundell *et al.* [160].

D. Unsupervised Learning With Attention

To properly exploit unsupervised learning is widely recognized as one of the most important long-term challenges of AI [146]. As already mentioned in Section V, attention is typically trained in a supervised architecture, although without a direct supervision on the attention weights. Nevertheless, a few works have recently attempted to exploit attention within purely unsupervised models. We believe this to be a promising research direction, as the learning process of humans is indeed largely unsupervised.

For example, in work by He *et al.* [161], attention is exploited in a model for aspect extraction in sentiment analysis, with the aim to remove words that are irrelevant for the sentiment and to ensure more coherence of the predicted aspects. In work by Zhang and Wu [162], attention is used within autoencoders in a question-retrieval task. The main idea is to generate semantic representations of questions, and self-attention is exploited during the encoding and decoding phases, with the objective to reconstruct the input sequences, as in traditional autoencoders. Following a similar idea, Zhang *et al.* [163] exploited bidimensional attention-based recursive autoencoders for bilingual phrase embeddings, whereas Tian and Fang [164] exploited attentive autoencoders to build sentence representations and performed topic modeling on short texts. Yang *et al.* [165] instead adopted an attention-driven approach to unsupervised sentiment modification in order to explicitly separate sentiment words from content words.

In computer vision, attention alignment has been proposed for unsupervised domain adaptation, with the aim to align the

attention patterns of networks trained on the source and target domain, respectively [166]. We believe that this could be an interesting scenario also for NLP.

E. Neural-Symbolic Learning and Reasoning

Recently, attention mechanisms started to be integrated within some neural-symbolic models, whose application to NLP scenarios is still at an early stage. For instance, in the context of neural logic programming [167], they have been exploited for reasoning over knowledge graphs, in order to combine parameter and structure learning of first-order logic rules. They have also been used in logic attention networks [168] to aggregate information coming from graph neighbors with both rule- and network-based attention weights. Moreover, prior knowledge has also been exploited by Shen *et al.* [169] to enable the attention mechanism to learn the knowledge representation of entities for ranking question–answer pairs.

Neural architectures exploiting attention performed well also in reasoning tasks that are also addressed with symbolic approaches, such as textual entailment [99]. For instance, Hudson and Manning [170] recently proposed a new architecture for complex reasoning problems, with NLP usually being one of the target sub-tasks, as in the case of visual question answering. In such an architecture, attention is used within several parts of the model, for example, over question words or to capture long-range dependences with self-attention.

An attempt to introduce constraints in the form of logical statements within neural networks has been proposed in [171] where rules governing attention are used to enforce word alignment in tasks, such as machine comprehension and natural language inference.

VII. CONCLUSION

Attention models have become ubiquitous in NLP applications. Attention can be applied to different input parts, different representations of the same data, or different features, to obtain a compact representation of the data as well as to highlight the relevant information. The selection is performed through a distribution function, which may consider locality in different dimensions, such as space, time, or even semantics. Attention can be used to compare the input data with a query element based on measures of similarity or significance. It can also autonomously learn what is to be considered relevant, by creating a representation encoding what the important data should be similar to. Integrating attention in neural architectures may thus yield a significant performance gain. Moreover, attention can be used as a tool for investigating the behavior of the network.

In this article, we have introduced a taxonomy of attention models, which enabled us to systematically chart a vast portion of the approaches in the literature and compare them to one another. To the best of our knowledge, this is the first systematic, comprehensive taxonomy of attention models for NLP.

We have also discussed the possible role of attention in addressing fundamental AI challenges. In particular, we have

shown how attention can be instrumental in injecting knowledge into the neural model, so as to represent specific features, or to exploit previously acquired knowledge, as in transfer learning settings. We speculate that this could pave the way to new challenging research avenues, where attention could be exploited to enforce the combination of subsymbolic models with symbolic knowledge representations to perform reasoning tasks or to address natural language understanding. Recent results also suggest that attention could be a key ingredient of unsupervised learning architectures, by guiding and focusing the training process where no supervision is given in advance.

REFERENCES

- [1] Y. Bao, S. Chang, M. Yu, and R. Barzilay, “Deriving machine attention from human rationales,” in *Proc. EMNLP*, 2018, pp. 1903–1913.
- [2] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” in *Proc. ICLR*, 2015, pp. 1–15.
- [3] H. Larochelle and G. E. Hinton, “Learning to combine foveal glimpses with a third-order Boltzmann machine,” in *Proc. NIPS*, 2010, pp. 1243–1251.
- [4] V. Mnih, N. Heess, A. Graves, and K. Kavukcuoglu, “Recurrent models of visual attention,” in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2014, pp. 2204–2212.
- [5] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and D. Pedreschi, “A survey of methods for explaining black box models,” *ACM Comput. Surv.*, vol. 51, no. 5, pp. 93:1–93:42, Aug. 2018.
- [6] K. Clark, U. Khandelwal, O. Levy, and C. D. Manning, “What does BERT look at? An analysis of BERT’s attention,” in *BlackboxNLP@ACL*. Florence, Italy: ACL, Aug. 2019, pp. 276–286.
- [7] G. Letarte, F. Paradis, P. Giguère, and F. Laviolette, “Importance of self-attention for sentiment analysis,” in *Proc. BlackboxNLP@EMNLP*, 2018, pp. 267–275.
- [8] S. Vashishth, S. Upadhyay, G. S. Tomar, and M. Faruqui, “Attention interpretability across NLP tasks,” 2019, *arXiv:1909.11218*. [Online]. Available: <http://arxiv.org/abs/1909.11218>
- [9] S. Wiegrefe and Y. Pinter, “Attention is not not explanation,” in *Proc. EMNLP/IJCNLP*, 2019, pp. 11–20.
- [10] S. Jain and B. C. Wallace, “Attention is not explanation,” in *Proc. NAACL-HLT*, 2019, pp. 3543–3556.
- [11] S. Serrano and N. A. Smith, “Is attention interpretable?” *CoRR*, vol. abs/1906.03731, 2019.
- [12] S. Liu, T. Li, Z. Li, V. Srikumar, V. Pascucci, and P.-T. Bremer, “Visual interrogation of attention-based models for natural language inference and machine comprehension,” in *Proc. EMNLP Syst. Demonstrations*, 2018, pp. 36–41.
- [13] J. Lee, J.-H. Shin, and J.-S. Kim, “Interactive visualization and manipulation of attention-based neural machine translation,” in *Proc. EMNLP Syst. Demonstrations*, 2017, pp. 121–126.
- [14] A. Gatt and E. Kraemer, “Survey of the state of the art in natural language generation: Core tasks, applications and evaluation,” *J. Artif. Intell. Res.*, vol. 61, pp. 65–170, Jan. 2018.
- [15] T. Young, D. Hazarika, S. Poria, and E. Cambria, “Recent trends in deep learning based natural language processing [review article],” *IEEE Comput. Intell. Mag.*, vol. 13, no. 3, pp. 55–75, Aug. 2018.
- [16] K. Xu *et al.*, “Show, attend and tell: Neural image caption generation with visual attention,” in *Proc. ICML*, vol. 37, 2015, pp. 2048–2057.
- [17] K. Gregor, I. Danihelka, A. Graves, D. Rezende, and D. Wierstra, “Draw: A recurrent neural network for image generation,” in *Proc. ICML*, in Proceedings of Machine Learning Research, vol. 37, 2015, pp. 1462–1471.
- [18] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, “Self-attention generative adversarial networks,” in *Proc. ICML*, vol. 97, 2019, pp. 7354–7363.
- [19] J. Chorowski, D. Bahdanau, K. Cho, and Y. Bengio, “End-to-end continuous speech recognition using attention-based recurrent NN: First results,” *CoRR*, vol. abs/1412.1602, 2014.
- [20] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Mar. 2016, pp. 4960–4964.
- [21] M. Sperber, J. Niehues, G. Neubig, S. Stüker, and A. Waibel, “Self-attentional acoustic models,” in *Proc. Interspeech*, Sep. 2018, pp. 3723–3727.

- [22] S. Wang, L. Hu, L. Cao, X. Huang, D. Lian, and W. Liu, "Attention-based transactional context embedding for next-item recommendation," in *Proc. AAAI*, 2018, pp. 2532–2539.
- [23] H. Ying *et al.*, "Sequential recommender system based on hierarchical attention networks," in *Proc. IJCAI*, Jul. 2018, pp. 3926–3932.
- [24] D. T. Tran, A. Iosifidis, J. Kannianen, and M. Gabbouj, "Temporal attention-augmented bilinear network for financial time-series data analysis," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 5, pp. 1407–1418, May 2019.
- [25] H. Song, D. Rajan, J. J. Thiagarajan, and A. Spanias, "Attend and diagnose: Clinical time series analysis using attention models," in *Proc. AAAI*, 2018, pp. 4091–4098.
- [26] J. Choi, B.-J. Lee, and B.-T. Zhang, "Multi-focus attention network for efficient deep reinforcement learning," *CoRR*, vol. abs/1712.04603, 2017.
- [27] O. Vinyals, M. Fortunato, and N. Jaitly, "Pointer networks," in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2015, pp. 2692–2700.
- [28] W. Kool, H. van Hoof, and M. Welling, "Attention, learn to solve routing problems!" in *Proc. ICLR*, 2019, pp. 1–25.
- [29] M. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," *CoRR*, vol. abs/1508.04025, 2015.
- [30] H. Mi, Z. Wang, and A. Ittycheriah, "Supervised attentions for neural machine translation," in *Proc. EMNLP*, 2016, pp. 2283–2288.
- [31] H. Mi, B. Sankaran, Z. Wang, and A. Ittycheriah, "Coverage embedding models for neural machine translation," in *Proc. EMNLP*, 2016, pp. 955–960.
- [32] F. Meng, Z. Lu, H. Li, and Q. Liu, "Interactive attention for neural machine translation," in *Proc. COLING*, 2016, pp. 2174–2185.
- [33] Z. Tu, Z. Lu, Y. Liu, X. Liu, and H. Li, "Modeling coverage for neural machine translation," in *Proc. ACL*, 2016, pp. 76–85.
- [34] J. Li, Z. Tu, B. Yang, M. R. Lyu, and T. Zhang, "Multi-head attention with disagreement regularization," in *Proc. EMNLP*, 2018, pp. 2897–2903.
- [35] L. Liu, M. Utiyama, A. Finch, and E. Sumita, "Neural machine translation with supervised attention," in *Proc. COLING*, 2016, pp. 3093–3102.
- [36] A. Vaswani *et al.*, "Attention is all you need," in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2017, pp. 5998–6008.
- [37] D. Britz, A. Goldie, M.-T. Luong, and Q. Le, "Massive exploration of neural machine translation architectures," in *Proc. EMNLP*, 2017, pp. 1442–1451.
- [38] A. Bapna, M. Chen, O. Firat, Y. Cao, and Y. Wu, "Training deeper neural machine translation models with transparent attention," in *Proc. EMNLP*, 2018, pp. 3028–3033.
- [39] B. Yang, Z. Tu, D. F. Wong, F. Meng, L. S. Chao, and T. Zhang, "Modeling localness for self-attention networks," in *Proc. EMNLP*, 2018, pp. 4449–4458.
- [40] K. Chen, R. Wang, M. Utiyama, E. Sumita, and T. Zhao, "Syntax-directed attention for neural machine translation," in *Proc. AAAI*, 2018, pp. 4792–4799.
- [41] L. Wu, F. Tian, L. Zhao, J. Lai, and T. Liu, "Word attention for sequence to sequence text understanding," in *Proc. AAAI*, 2018, pp. 5578–5585.
- [42] T. Domhan, "How much attention do you need? A granular analysis of neural machine translation architectures," in *Proc. ACL*, 2018, pp. 1799–1808.
- [43] S. Zhao and Z. Zhang, "Attention-via-attention neural machine translation," in *Proc. AAAI*, 2018, pp. 563–570.
- [44] J. Lin, X. Sun, X. Ren, M. Li, and Q. Su, "Learning when to concentrate or divert attention: Self-adaptive attention temperature for neural machine translation," in *Proc. EMNLP*, 2018, pp. 2985–2990.
- [45] J. Lin, S. Ma, Q. Su, and X. Sun, "Decoding-history-based adaptive control of attention for neural machine translation," *CoRR*, vol. abs/1802.01812, 2018.
- [46] C. Malaviya, P. Ferreira, and A. F. T. Martins, "Sparse and constrained attention for neural machine translation," in *Proc. ACL*, 2018, pp. 370–376.
- [47] Y. Kim, C. Denton, L. Hoang, and A. M. Rush, "Structured attention networks," in *Proc. ICLR*, 2017, pp. 1–21.
- [48] P. Michel, O. Levy, and G. Neubig, "Are sixteen heads really better than one?" in *Proc. NeurIPS*, Vancouver, BC, Canada, Dec. 2019, pp. 14014–14024.
- [49] E. Voita, P. Serdyukov, R. Sennrich, and I. Titov, "Context-aware neural machine translation learns anaphora resolution," in *Proc. ACL*, 2018, pp. 1264–1274.
- [50] E. Voita, D. Talbot, F. Moiseev, R. Sennrich, and I. Titov, "Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned," in *Proc. ACL*, 2019, pp. 5797–5808.
- [51] A. F. T. Martins and J. Kreutzer, "Learning what's easy: Fully differentiable neural easy-first taggers," in *Proc. EMNLP*, 2017, pp. 349–362.
- [52] Z. Yang, D. Yang, C. Dyer, X. He, A. Smola, and E. Hovy, "Hierarchical attention networks for document classification," in *Proc. HLT-NAACL*, 2016, pp. 1480–1489.
- [53] Y. Zhang, I. Marshall, and B. C. Wallace, "Rationale-augmented convolutional neural networks for text classification," in *Proc. EMNLP*, 2016, pp. 795–804.
- [54] J. Pavlopoulos, P. Malakasiotis, and I. Androutsopoulos, "Deeper attention to abusive user content moderation," in *Proc. EMNLP*, 2017, pp. 1125–1135.
- [55] A. M. Rush, S. Chopra, and J. Weston, "A neural attention model for abstractive sentence summarization," in *Proc. EMNLP*, 2015, pp. 379–389.
- [56] H. Li, J. Zhu, T. Liu, J. Zhang, and C. Zong, "Multi-modal sentence summarization with modality attention and image filtering," in *Proc. IJCAI*, 2018, pp. 4152–4158.
- [57] A. Lauscher, G. Glavaš, S. P. Ponzetto, and K. Eckert, "Investigating the role of argumentation in the rhetorical analysis of scientific publications with neural multi-task learning models," in *Proc. EMNLP*, 2018, pp. 3326–3338.
- [58] I. C. T. Ngoko, A. Mukherjee, and B. Kabaso, "Abstractive text summarization using recurrent neural networks: Systematic literature review," in *Proc. ICICKM*, vol. 13, 2018, pp. 435–439.
- [59] S. Sukhbaatar, A. Szlam, J. Weston, and R. Fergus, "End-to-end memory networks," in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2015, pp. 2440–2448.
- [60] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [61] M. Daniluk, T. Rocktäschel, J. Welbl, and S. Riedel, "Frustratingly short attention spans in neural language modeling," in *Proc. ICLR*, 2017, pp. 1–10.
- [62] K. M. Hermann *et al.*, "Teaching machines to read and comprehend," in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2015, pp. 1693–1701.
- [63] J. Cheng, L. Dong, and M. Lapata, "Long short-term memory-networks for machine reading," in *Proc. EMNLP*, 2016, pp. 551–561.
- [64] F. Hill, A. Bordes, S. Chopra, and J. Weston, "The goldilocks principle: Reading children's books with explicit memory representations," in *Proc. ICLR*, 2016, pp. 1–13.
- [65] Y. Cui, T. Liu, Z. Chen, S. Wang, and G. Hu, "Consensus attention-based neural networks for Chinese reading comprehension," in *Proc. COLING*, 2016, pp. 1777–1786.
- [66] R. Kadlec, M. Schmid, O. Bajgar, and J. Kleindienst, "Text understanding with the attention sum reader network," in *Proc. ACL*, 2016, pp. 908–918.
- [67] A. Sordani, P. Bachman, and Y. Bengio, "Iterative alternating neural attention for machine reading," *CoRR*, vol. abs/1606.02245, 2016.
- [68] B. Wang, K. Liu, and J. Zhao, "Inner attention based recurrent neural networks for answer selection," in *Proc. ACL*, 2016, pp. 1288–1297.
- [69] C. N. dos Santos, M. Tan, B. Xiang, and B. Zhou, "Attentive pooling networks," *CoRR*, vol. abs/1602.03609, 2016.
- [70] Y. Cui, Z. Chen, S. Wei, S. Wang, T. Liu, and G. Hu, "Attention-over-attention neural networks for reading comprehension," in *Proc. ACL*, 2017, pp. 593–602.
- [71] B. Dhingra, H. Liu, Z. Yang, W. Cohen, and R. Salakhutdinov, "Gated-attention readers for text comprehension," in *Proc. ACL*, 2017, pp. 1832–1846.
- [72] Y. Shen, P.-S. Huang, J. Gao, and W. Chen, "ReasoNet: Learning to stop reading in machine comprehension," in *Proc. KDD*, 2017, pp. 1047–1055.
- [73] S. Kundu and H. T. Ng, "A question-focused multi-factor attention network for question answering," in *Proc. AAAI*, 2018, pp. 5828–5835.
- [74] H. Zhu, F. Wei, B. Qin, and T. Liu, "Hierarchical attention flow for multiple-choice reading comprehension," in *Proc. AAAI*, 2018, pp. 6077–6085.
- [75] Y. Tay, A. T. Luu, and S. C. Hui, "Hermitian co-attention networks for text matching in asymmetrical domains," in *Proc. IJCAI*, Jul. 2018, pp. 4425–4431.
- [76] Y. Hao *et al.*, "An end-to-end model for question answering over knowledge base with cross-attention combining global knowledge," in *Proc. ACL*, 2017, pp. 221–231.
- [77] Y. Diao *et al.*, "WECA: A WordNet-encoded collocation-attention network for homographic pun recognition," in *Proc. EMNLP*, 2018, pp. 2507–2516.

- [78] A. Zadeh, P. P. Liang, S. Poria, P. Vij, E. Cambria, and L. Morency, “Multi-attention recurrent network for human communication comprehension,” in *Proc. AAAI*, 2018, pp. 5642–5649.
- [79] H. Chen, G. Ding, Z. Lin, Y. Guo, and J. Han, “Attend to knowledge: Memory-enhanced attention network for image captioning,” in *Advances in Brain Inspired Cognitive Systems* Cham, Switzerland: Springer, 2018, pp. 161–171.
- [80] J. Lu, J. Yang, D. Batra, and D. Parikh, “Hierarchical question-image co-attention for visual question answering,” in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2016, pp. 289–297.
- [81] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, “Stacked attention networks for image question answering,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 21–29.
- [82] J. Song, P. Zeng, L. Gao, and H. T. Shen, “From pixels to objects: Cubic visual attention for visual question answering,” in *Proc. IJCAI*, Jul. 2018, pp. 906–912.
- [83] D. S. Chaplot, K. M. Sathyendra, R. K. Pasumarthi, D. Rajagopal, and R. Salakhutdinov, “Gated-attention architectures for task-oriented language grounding,” in *Proc. AAAI*, 2018, pp. 2819–2826.
- [84] R. Zhang, C. N. dos Santos, M. Yasunaga, B. Xiang, and D. Radev, “Neural coreference resolution with deep biaffine attention by joint mention detection and mention clustering,” in *Proc. ACL*, 2018, pp. 102–107.
- [85] G. Kundu, A. Sil, R. Florian, and W. Hamza, “Neural cross-lingual coreference resolution and its application to entity linking,” in *Proc. ACL*, 2018, pp. 395–400.
- [86] Q. Zhang, J. Fu, X. Liu, and X. Huang, “Adaptive co-attention network for named entity recognition in tweets,” in *Proc. AAAI*, 2018, pp. 5674–5681.
- [87] R. Dong and D. Smith, “Multi-input attention for unsupervised OCR correction,” in *Proc. ACL*, 2018, pp. 2363–2372.
- [88] F. Nie, Y. Cao, J. Wang, C. Lin, and R. Pan, “Mention and entity description co-attention for entity disambiguation,” in *Proc. AAAI*, 2018, pp. 5908–5915.
- [89] A. Parikh, O. Täckström, D. Das, and J. Uszkoreit, “A decomposable attention model for natural language inference,” in *Proc. EMNLP*, 2016, pp. 2249–2255.
- [90] Y. Tay, A. T. Luu, and S. C. Hui, “Compare, compress and propagate: Enhancing neural architectures with alignment factorization for natural language inference,” in *Proc. EMNLP*, 2018, pp. 1565–1575.
- [91] A. Martins and R. Astudillo, “From softmax to sparsemax: A sparse model of attention and multi-label classification,” in *Proc. ICML*, vol. 48, 2016, pp. 1614–1623.
- [92] S. Wang and J. Jiang, “Learning natural language inference with LSTM,” in *Proc. HLT-NAACL*, 2016, pp. 1442–1451.
- [93] T. Shen, T. Zhou, G. Long, J. Jiang, S. Pan, and C. Zhang, “Disan: Directional self-attention network for RNN/CNN-free language understanding,” in *Proc. AAAI*, 2018, pp. 5446–5455.
- [94] T. Shen, T. Zhou, G. Long, J. Jiang, S. Wang, and C. Zhang, “Reinforced self-attention network: A hybrid of hard and soft attention for sequence modeling,” in *Proc. IJCAI*, Jul. 2018, pp. 4345–4352.
- [95] Q. Chen, Z.-H. Ling, and X. Zhu, “Enhancing sentence embedding with generalized pooling,” in *Proc. COLING*, 2018, pp. 1815–1826.
- [96] I. Lopez-Gazpio, M. Maritxalar, M. Lapata, and E. Agirre, “Word n-gram attention models for sentence similarity and inference,” *Expert Syst. Appl.*, vol. 132, pp. 1–11, Oct. 2019. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0957417419302842>
- [97] E. Strubell, P. Verga, D. Andor, D. Weiss, and A. McCallum, “Linguistically-informed self-attention for semantic role labeling,” in *Proc. EMNLP*, 2018, pp. 5027–5038.
- [98] Z. Tan, M. Wang, J. Xie, Y. Chen, and X. Shi, “Deep semantic role labeling with self-attention,” in *Proc. AAAI*, 2018, pp. 4929–4936.
- [99] T. Rocktäschel, E. Grefenstette, K. M. Hermann, T. Kociský, and P. Blunsom, “Reasoning about entailment with neural attention,” in *Proc. ICLR*, 2016, pp. 1–9.
- [100] Z. Lin *et al.*, “A structured self-attentive sentence embedding,” in *Proc. ICLR*, 2017, pp. 1–15.
- [101] F. Luo, T. Liu, Q. Xia, B. Chang, and Z. Sui, “Incorporating glosses into neural word sense disambiguation,” in *Proc. ACL*, 2018, pp. 2473–2482.
- [102] O. Vinyals, L. U. Kaiser, T. Koo, S. Petrov, I. Sutskever, and G. Hinton, “Grammar as a foreign language,” in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2015, pp. 2773–2781.
- [103] N. Kitaev and D. Klein, “Constituency parsing with a self-attentive encoder,” in *Proc. ACL*, 2018, pp. 2676–2686.
- [104] R. Kohita, H. Noji, and Y. Matsumoto, “Dynamic feature selection with attention in incremental parsing,” in *Proc. COLING*, 2018, pp. 785–794.
- [105] T. Dozat and C. D. Manning, “Deep biaffine attention for neural dependency parsing,” in *Proc. ICLR*, 2017, pp. 1–8.
- [106] D. Tang, B. Qin, and T. Liu, “Aspect level sentiment classification with deep memory network,” in *Proc. EMNLP*, 2016, pp. 214–224.
- [107] J. Cheng, S. Zhao, J. Zhang, I. King, X. Zhang, and H. Wang, “Aspect-level sentiment classification with heat (hierarchical attention) network,” in *Proc. CIKM*, 2017, pp. 97–106.
- [108] Y. Wang, M. Huang, X. Zhu, and L. Zhao, “Attention-based LSTM for aspect-level sentiment classification,” in *Proc. EMNLP*, 2016, pp. 606–615.
- [109] R. He, W. S. Lee, H. T. Ng, and D. Dahlmeier, “Effective attention modeling for aspect-level sentiment classification,” in *Proc. COLING*, 2018, pp. 1121–1131.
- [110] P. Zhu and T. Qian, “Enhanced aspect level sentiment classification with auxiliary memory,” in *Proc. COLING*, 2018, pp. 1077–1087.
- [111] D. Ma, S. Li, X. Zhang, and H. Wang, “Interactive attention networks for aspect-level sentiment classification,” in *Proc. IJCAI*, 2017, pp. 4068–4074.
- [112] F. Fan, Y. Feng, and D. Zhao, “Multi-grained attention network for aspect-level sentiment classification,” in *Proc. EMNLP*, 2018, pp. 3433–3442.
- [113] Y. Ma, H. Peng, and E. Cambria, “Targeted aspect-based sentiment analysis via embedding commonsense knowledge into an attentive LSTM,” in *Proc. AAAI*, 2018, pp. 5876–5883.
- [114] Y. Tay, A. T. Luu, S. C. Hui, and J. Su, “Attentive gated lexicon reader with contrastive contextual co-attention for sentiment classification,” in *Proc. EMNLP*, 2018, pp. 3443–3453.
- [115] J. Yu, L. Marujo, J. Jiang, P. Karuturi, and W. Brendel, “Improving multi-label emotion classification via sentiment classification with dual attention transfer network,” in *Proc. EMNLP*, 2018, pp. 1097–1102.
- [116] Z. Li, Y. Wei, Y. Zhang, and Q. Yang, “Hierarchical attention transfer network for cross-domain sentiment classification,” in *Proc. AAAI*, 2018, pp. 5852–5859.
- [117] Y. Fan, F. Tian, T. Qin, J. Bian, and T.-Y. Liu, “Learning what data to learn,” *CoRR*, 2017, pp. 1–7.
- [118] Y. Long, R. Xiang, Q. Lu, C.-R. Huang, and M. Li, “Improving attention model based on cognition grounded data for sentiment analysis,” *IEEE Trans. Affect. Comput.*, early access, Mar. 4, 2019, doi: [10.1109/TAFFC.2019.2903056](https://doi.org/10.1109/TAFFC.2019.2903056).
- [119] J. Du, L. Gui, Y. He, R. Xu, and X. Wang, “Convolution-based neural attention with applications to sentiment classification,” *IEEE Access*, vol. 7, pp. 27983–27992, 2019.
- [120] X. Li, L. Bing, P. Li, W. Lam, and Z. Yang, “Aspect term extraction with history attention and selective transformation,” in *Proc. IJCAI*, 2018, pp. 4194–4200.
- [121] D. Chen, J. Du, L. Bing, and R. Xu, “Hybrid neural attention for agreement/disagreement inference in online debates,” in *Proc. EMNLP*, 2018, pp. 665–670.
- [122] I. Habernal and I. Gurevych, “What makes a convincing argument? Empirical analysis and detecting attributes of convincings in Web argumentation,” in *Proc. EMNLP*, 2016, pp. 1214–1223.
- [123] C. Stab, T. Miller, B. Schiller, P. Rai, and I. Gurevych, “Cross-topic argument mining from heterogeneous sources,” in *Proc. EMNLP*, 2018, pp. 3664–3674.
- [124] M. Spliethöver, J. Klaff, and H. Heuer, “Is it worth the attention? A comparative evaluation of attention layers for argument unit segmentation,” in *ArgMining@ACL*. Florence, Italy: ACL, Aug. 2019, pp. 74–82.
- [125] M. Spliethöver, J. Klaff, and H. Heuer, “Is it worth the attention? A comparative evaluation of attention layers for argument unit segmentation,” in *Proc. 6th Workshop Argument Mining*. Florence, Italy: Association for Computational Linguistics, Aug. 2019, pp. 74–82. [Online]. Available: <https://www.aclweb.org/anthology/W19-4509>
- [126] F. Barbieri, L. E. Anke, J. Camacho-Collados, S. Schockaert, and H. Saggion, “Interpretable emoji prediction via label-wise attention LSTMs,” in *Proc. EMNLP*, 2018, pp. 4766–4771.
- [127] X. Li, K. Song, S. Feng, D. Wang, and Y. Zhang, “A co-attention neural network model for emotion cause analysis with emotional context awareness,” in *Proc. EMNLP*, 2018, pp. 4752–4757.
- [128] L. Gui, J. Hu, Y. He, R. Xu, Q. Lu, and J. Du, “A question answering approach for emotion cause extraction,” in *Proc. EMNLP*, 2017, pp. 1593–1602.

- [129] X. Wang, Y. Hua, E. Kodirov, G. Hu, and N. M. Robertson, "Deep metric learning by online soft mining and class-aware attention," in *Proc. AAAI*, 2019, pp. 5361–5368, doi: 10.1609/aaai.v33i01.33015361.
- [130] J. B. Lee, R. A. Rossi, S. Kim, N. K. Ahmed, and E. Koh, "Attention models in graphs: A survey," *ACM Trans. Knowl. Discovery Data*, vol. 13, no. 6, pp. 62:1–62:25, Nov. 2019.
- [131] Y. Goldberg, *Neural Network Methods for Natural Language Processing* (Synthesis Lectures on Human Language Technologies), vol. 10, no. 1. San Rafael, CA, USA: Morgan & Claypool, 2017.
- [132] G. Tang, M. Müller, A. Rios, and R. Sennrich, "Why self-attention? A targeted evaluation of neural machine translation architectures," in *Proc. EMNLP*, 2018, pp. 4263–4272.
- [133] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *IEEE Trans. Signal Process.*, vol. 45, no. 11, pp. 2673–2681, Nov. 1997.
- [134] C. Ma, L. Liu, A. Tamura, T. Zhao, and E. Sumita, "Deterministic attention for sequence-to-sequence constituent parsing," in *Proc. AAAI*, 2017, pp. 3237–3243.
- [135] L. Liu, M. Zhu, and S. Shi, "Improving sequence-to-sequence constituency parsing," in *Proc. AAAI*, 2018, pp. 4873–4880.
- [136] S. Maharjan, M. Montes, F. A. González, and T. Solorio, "A genre-aware attention model to improve the likability prediction of books," in *Proc. EMNLP*, 2018, pp. 3381–3391.
- [137] T. Qiao, J. Dong, and D. Xu, "Exploring human-like attention supervision in visual question answering," in *Proc. AAAI*, 2018, pp. 7300–7307.
- [138] J.-K. Kim and Y.-B. Kim, "Supervised domain enablement attention for personalized domain classification," in *Proc. EMNLP*, 2018, pp. 894–899.
- [139] D. Kiela, C. Wang, and K. Cho, "Dynamic meta-embeddings for improved sentence representations," in *Proc. EMNLP*, 2018, pp. 1466–1477.
- [140] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Trans. Neural Netw.*, vol. 5, no. 2, pp. 157–166, Mar. 1994.
- [141] C. Wu, F. Wu, J. Liu, and Y. Huang, "Hierarchical user and item representation with three-tier attention for recommendation," in *Proc. NAACL-HLT*, 2019, pp. 1818–1826.
- [142] A. Graves, G. Wayne, and I. Danihelka, "Neural Turing machines," *CoRR*, vol. abs/1410.5401, 2014.
- [143] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proc. AISTATS*, vol. 15, 2011, pp. 315–323.
- [144] G. Klambauer, T. Unterthiner, A. Mayr, and S. Hochreiter, "Self-normalizing neural networks," in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2017, pp. 971–980.
- [145] J. Du, J. Han, A. Way, and D. Wan, "Multi-level structured self-attentions for distantly supervised relation extraction," in *Proc. EMNLP*, 2018, pp. 2216–2225.
- [146] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, p. 436, 2015.
- [147] L. Getoor and B. Taskar, *Introduction to Statistical Relational Learning*, vol. 1. Cambridge, MA, USA: MIT Press, 2007.
- [148] A. S. D. Garcez, K. B. Broda, and D. M. Gabbay, *Neural-Symbolic Learning Systems: Foundations and Applications*. Berlin, Germany: Springer, 2012.
- [149] M. Lippi, "Reasoning with deep learning: An open challenge," in *Proc. CEUR Workshop*, vol. 1802, 2017, pp. 38–43.
- [150] O. Zaidan, J. Eisner, and C. Piatko, "Using 'annotator rationales' to improve machine learning for text categorization," in *Proc. HLT-NAACL*, 2007, pp. 260–267.
- [151] C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Mach. Intell.*, vol. 1, pp. 206–215, May 2019.
- [152] M. O. Riedl, "Human-centered artificial intelligence and machine learning," *Hum. Behav. Emerg. Technol.*, vol. 1, no. 1, pp. 33–36, 2019.
- [153] Q. V. Le, "Building high-level features using large scale unsupervised learning," in *Proc. IEEE ICASSP*, May 2013, pp. 8595–8598.
- [154] X. Shi, I. Padhi, and K. Knight, "Does string-based neural MT learn source syntax?" in *Proc. EMNLP*, 2016, pp. 1526–1534.
- [155] M. Hu *et al.*, "Attention-guided answer distillation for machine reading comprehension," in *Proc. EMNLP*, 2018, pp. 2077–2086.
- [156] S. Moon and J. G. Carbonell, "Completely heterogeneous transfer learning with attention—What and what not to transfer," in *Proc. IJCAI*, 2017, pp. 2508–2514.
- [157] J. Heo *et al.*, "Uncertainty-aware attention for reliable interpretation and prediction," in *Proc. NIPS*. Red Hook, NY, USA: Curran Associates, 2018, pp. 917–926.
- [158] M. Poggi and S. Mattocchia, "Learning a general-purpose confidence measure based on O(1) features and a smarter aggregation strategy for semi global matching," in *Proc. 3DV*, 2016, pp. 509–518.
- [159] R. Cipolla, Y. Gal, and A. Kendall, "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 7482–7491.
- [160] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural networks," in *Proc. ICML*, 2015, pp. 1613–1622.
- [161] R. He, W. S. Lee, H. T. Ng, and D. Dahlmeier, "An unsupervised neural attention model for aspect extraction," in *Proc. ACL*, vol. 1, 2017, pp. 388–397.
- [162] M. Zhang and Y. Wu, "An unsupervised model with attention autoencoders for question retrieval," in *Proc. AAAI*, 2018, pp. 4978–4986.
- [163] B. Zhang, D. Xiong, and J. Su, "Battara: Bidimensional attention-based recursive autoencoders for learning bilingual phrase embeddings," in *Proc. AAAI*, 2017, pp. 3372–3378.
- [164] T. Tian and Z. F. Fang, "Attention-based autoencoder topic model for short texts," *Procedia Comput. Sci.*, vol. 151, pp. 1134–1139, Jan. 2019.
- [165] P. Yang, J. Lin, J. Xu, J. Xie, Q. Su, and X. Sun, "Specificity-driven cascading approach for unsupervised sentiment modification," in *Proc. EMNLP-IJCNLP*, 2019, pp. 5511–5520.
- [166] G. Kang, L. Zheng, Y. Yan, and Y. Yang, "Deep adversarial attention alignment for unsupervised domain adaptation: The benefit of target expectation maximization," in *Proc. ECCV*, 2018, pp. 401–416.
- [167] F. Yang, Z. Yang, and W. W. Cohen, "Differentiable learning of logical rules for knowledge base reasoning," in *Proc. NIPS*, 2017, pp. 2319–2328.
- [168] P. Wang, J. Han, C. Li, and R. Pan, "Logic attention based neighborhood aggregation for inductive knowledge graph embedding," in *Proc. AAAI*, vol. 33, 2019, pp. 7152–7159.
- [169] Y. Shen *et al.*, "Knowledge-aware attentive neural network for ranking question answer pairs," in *Proc. SIGIR*, 2018, pp. 901–904.
- [170] D. A. Hudson and C. D. Manning, "Compositional attention networks for machine reasoning," in *Proc. ICLR*, 2018, pp. 1–20.
- [171] T. Li and V. Srikumar, "Augmenting neural networks with first-order logic," in *Proc. ACL*, 2019, pp. 292–302.



Andrea Galassi received the master's degree in computer engineering from the Department of Computer Science and Engineering, University of Bologna, Bologna, Italy, in 2017, where he is currently pursuing the Ph.D. degree.

His research activity concerns artificial intelligence and machine learning, focusing on argumentation mining and related natural language processing (NLP) tasks. Other research interests involve deep learning applications to games and constraint satisfaction problems (CSPs).



Marco Lippi held positions at the University of Florence, Florence, Italy, University of Siena, Siena, Italy, and the University of Bologna, Bologna, Italy. He was a Visiting Scholar with Université Pierre et Marie Curie (UPMC), Paris, France. He is currently an Associate Professor with the Department of Sciences and Methods for Engineering, University of Modena and Reggio Emilia, Modena, Italy. His work focuses on machine learning and artificial intelligence, with applications to several areas, including argumentation mining, legal informatics, and medicine.

Prof. Lippi received the "E. Caianello" Prize for the Best Italian Ph.D. Thesis in the field of neural networks in 2012.



Paolo Torroni has been an Associate Professor with the University of Bologna, Bologna, Italy, since 2015. He has edited over 20 books and special issues and authored over 150 articles in computational logics, multiagent systems, argumentation, and natural language processing. His main research interest includes artificial intelligence.

Prof. Torroni serves as an Associate Editor for *Fundamenta Informaticae* and *Intelligenza Artificiale*.