
Attention as Inference via Fenchel Duality

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Attention has been widely adopted in many state-of-the-art deep learning models.
2 While the significant performance improvements it brings have attracted great
3 interest, attention is still poorly understood theoretically. This paper presents a new
4 perspective to understand attention by showing that it can be seen as a solver of a
5 family of estimation problems. In particular, we describe a convex optimization
6 problem that arises in a family of estimation tasks commonly appearing in the de-
7 sign of deep learning models. Rather than directly solving the convex optimization
8 problem, we solve its Fenchel dual and derive a closed-form approximation of the
9 optimal solution. Remarkably, the solution gives a generalized attention structure,
10 and its special case is equivalent to the popular dot-product attention adopted in
11 transformer networks. We show that T5 transformer has implicitly adopted the
12 general form of the solution by demonstrating that this expression unifies the word
13 mask and the positional encoding functions. Finally, we discuss how the proposed
14 attention structures can be integrated in practical models.

15 1 Introduction

16 Attention-based deep neural networks are now integrated into cutting-edge language models that
17 have revolutionized a broad range of tasks: machine translation [1, 15], sentiment classification [27],
18 image captioning [29] and unsupervised representation learning [5], etc. Especially, attention plays a
19 pivotal role in the construction of the transformer architecture [25], which has had a profound impact
20 on the deep learning field.

21 Despite great empirical success, the driving principles of attention are still poorly understood. This
22 lack of understanding impedes practitioners from confidently and appropriately using attention layers
23 and makes it challenging to develop new attention-based neural architectures.

24 In this paper, we offer a new perspective for understanding attention by showing that it is in fact
25 a solver for a certain type of optimization problem that corresponds to an inference task. We give
26 several examples, all of which can be characterized as follows: given 1) an unreliable estimate of
27 the mean of an unknown distribution p on \mathbb{R}^d and 2) a preference distribution u on \mathbb{R}^d encoding
28 beliefs on p 's selection, the inference task is to get a better estimate of p 's mean given its unreliable
29 estimate and u . We derive a convex optimization problem that is abstracted from the task and solve it
30 by instead solving its Fenchel dual [22, p.104]. Remarkably, the derived expression of the improved
31 estimate of p gives a generalized attention structure whose special case is equivalent to the popular
32 dot-product attention [15] that is also applied in the transformer network [25]. In addition, we show
33 that our generalized attention expression has been implicitly adopted by T5 transformer [19] as the
34 expression unifies the concept of word masks and its positional encoding functions. Extra examples
35 are given to show how the generalized attention structures can be used in practice.

36 2 Related Works

37 Since 2019, several authors have investigated the properties and working mechanism of attention.
 38 This series of works mainly addresses whether the attention mechanism can serve as a proxy of
 39 saliency [9, 18, 23, 24, 26, 28]. Most of these works obtain insights into the attention mechanism
 40 by performing empirical studies. The related methods include analyzing the behaviours of trained
 41 attention-based networks [4], or pruning a few heads, or analyzing the effects of altering the attention
 42 weights [18, 26], or a mixture of these [9, 24].

43 Apart from understanding attention empirically, some theoretical results presented by Brunner et
 44 al. [3] and Hahn [7] show that the self-attention layers are not identifiable. This implies there could
 45 exist multiple combinations of attention weights that can provide equally good final predictions.
 46 In particular, such non-uniqueness means that the use of attention may complicate interpretability.
 47 Another important approach to understand attention is to analyze its asymptotic behaviour when
 48 the number of heads and the network width approach infinity [8, 30]. In this limiting case, the
 49 entire network can be seen as a Gaussian process [13] and its behaviours can be characterized by
 50 closed-form expressions that are not available in the finite case.

51 Very recently (since 2021) several theoretical works have appeared that study attention outside the
 52 asymptotic regime. Lu et al. [14] set up a simple attention-based classification model and derive a
 53 closed-form relationship between the word’s embedding norm and the product of its key and the
 54 query. They empirically show that such relationship also exists in a more complicated and practical
 55 configuration. Ramsauer et al. [20] construct an equivalence relationship between attention and
 56 a newly proposed Hopfield network with continuous states. In particular, they show that the new
 57 Hopfield network’s update rule is equivalent to the attention mechanism used in transformers [25].

58 3 A Motivating Example

59 We first consider a seemingly unrelated example, to illustrate the key ingredients of this paper.

60 Assume a probability distribution p on \mathbb{R}^d has a spherical Gaussian prior $u \sim \mathcal{N}(\mu, I_d)$. Let h_p
 61 denote the mean of the unknown p . Given an unreliable observation b of h_p , what is the best guess
 62 of h_p ? To solve this problem, we may formulate the following optimization problem

$$p^* = \arg \min_p \frac{\alpha}{2} \left\| b - \int \mathbf{a} p(\mathbf{a}) \, d\mathbf{a} \right\|^2 + \mathcal{K}(p, u), \quad (1)$$

63 with $\alpha > 0$ responsible for the relative strength of the two terms, where $\mathcal{K}(p, u)$ denotes the KL
 64 divergence between p and u . The basic idea behind (1) is that: although b is not reliable, it should
 65 not be too far from $h_p = \int \mathbf{a} p(\mathbf{a}) \, d\mathbf{a}$. Also, as u encodes the preferred value of p , we add the
 66 KL divergence term to show preference for p that is close to u . As will be discussed later, such a
 67 formulation can be either obtained from the maximum likelihood principle or from the maximum
 68 entropy principle [10, 11]. In particular, Rioux et al [21] develop (1) for image de-blurring by applying
 69 Maximum Entropy on the Mean (MEM), an information-theoretic method due to Gamboa [6] but not
 70 yet widely known in machine learning.

71 After obtaining the minimizer p^* of (1), its mean $\int \mathbf{a} p^*(\mathbf{a}) \, d\mathbf{a}$ gives our estimate of h_p . Rioux et
 72 al. [21] prove, via Fenchel duality [22, p.104] that the minimizer p^* takes the form

$$p^*(\mathbf{a}) = \frac{u(\mathbf{a}) \exp\langle \mathbf{a}, \lambda^* \rangle}{\int u(\mathbf{a}') \exp\langle \mathbf{a}', \lambda^* \rangle \, d\mathbf{a}'}, \quad (2)$$

73 where

$$\lambda^* = \arg \max_{\lambda \in \mathbb{R}^d} \langle b, \lambda \rangle - \frac{1}{2\alpha} \|\lambda\|^2 - \log \int_{\mathbb{R}^d} u(\mathbf{a}) \exp\langle \mathbf{a}, \lambda \rangle \, d\mathbf{a}. \quad (3)$$

74 Note that $\int_{\mathbb{R}^d} u(\mathbf{a}) \exp\langle \mathbf{a}, \lambda \rangle \, d\mathbf{a} = \exp(\langle \mu, \lambda \rangle + \frac{1}{2} \|\lambda\|^2)$ as it is the moment generating function
 75 (MGF) of $u \sim \mathcal{N}(\mu, I_d)$. Substituting the expression into (3) followed by setting the derivative with
 76 respect to λ to zero yields $\lambda^* = \frac{\alpha}{\alpha+1} (b - \mu)$. By (2), $p^*(\mathbf{a}) \propto \exp(-\frac{1}{2} \|\mathbf{a} - \mu\|^2 + \langle \mathbf{a}, \lambda^* \rangle) \propto$
 77 $\exp(-\frac{1}{2} \|\mathbf{a} - (\mu + \lambda^*)\|^2)$. Substituting $\lambda^* = \frac{\alpha}{\alpha+1} (b - \mu)$ into it implies that p^* follows a Gaussian
 78 distribution $\mathcal{N}(\frac{1}{1+\alpha} \mu + \frac{\alpha}{1+\alpha} b, I_d)$. Thus, our estimate of h_p is $\frac{1}{1+\alpha} \mu + \frac{\alpha}{1+\alpha} b$.

79 In this paper, we focus on a similar optimization problem that estimates h_p assuming that u is instead
80 a discrete distribution. We show that such optimization problems naturally and frequently arise in
81 neural network designs. By solving the optimization problem, we derive a closed-form approximation
82 for the estimate of h_p , via Fenchel duality. The approximation then gives a generalized attention layer
83 structure as shown in Fig 1. A special case of it is equivalent to the familiar dot-product attention [15]
84 that is also adopted in transformers [25]. Moreover, we will show that T5 transformer [19] implicitly
85 adopts our generalized attention expression.

86 4 Setup of a Design Problem

87 Throughout the rest of the paper, we
88 consider a machine learning problem
89 in which the objective is to predict an
90 output quantity Y from a given input
91 X . Additionally, Y may include K
92 components, namely, be expressed as
93 $(Y^{(1)}, Y^{(2)}, \dots, Y^{(K)})$. To be more
94 concrete, we present a few example machine
95 learning problems and let them
96 run through our development.

97 **Example: Translation Problem.** In
98 this problem, the input X is a sentence,
99 or a sequence of words, in the source
100 language, and output Y is the sequence
101 of words in the target sentence, where
102 $Y^{(k)}$ denotes the k^{th} word.

103 **Example: Image Captioning.** In this problem, the input X is a raw image and output Y is the
104 sequence of words in the caption, where $Y^{(k)}$ denotes the k^{th} word.

105 **Example: Filling in the blanks task.** This task has been used to train the BERT model [5]. The
106 input X is a sequence of words with certain percentage of words masked. The output Y are the
107 predicted masked words, where $Y^{(k)}$ denotes the k^{th} masked one.

108 The objective of any of these problems and that we address in this paper is to learn a function F ,
109 mapping from the space of X to the space of Y so that $Y = F(X)$. We will denote by $F^{(k)}$ the part
110 of F responsible for predicting $Y^{(k)}$ (Fig 1a), namely, $Y^{(k)} = F^{(k)}(X)$. Although we here express
111 F as separate functions $(F^{(1)}, F^{(2)}, \dots, F^{(K)})$, we note that it is in fact possible that different $F^{(k)}$'s
112 share some component in common. We now focus on the design of $F^{(k)}$.

113 We restrict the architecture of $F^{(k)}$ to the form in Fig 1b with the main focus on the inference of $h^{(k)}$.
114 The extraction of feature $h^{(k)}$ is via two parallel modules $f_{\text{evd}}^{(k)}$ and $f_{\text{pref}}^{(k)}$ that directly operate on the
115 input X followed by a function $g^{(k)}$ (in Fig 1c), which we will design.

116 **The Design Problem** We describe the problem of designing g as follows.

117 Suppose that there is an unknown distribution $p^{(k)}$ on \mathbb{R}^d whose mean vector is $h^{(k)}$, namely,

$$h^{(k)} = \int_{\mathbb{R}^d} \mathbf{a} p^{(k)}(\mathbf{a}) \, d\mathbf{a}. \quad (4)$$

118 Let $u^{(k)}$ be another distribution on \mathbb{R}^d that is generated as the output of a network module $f_{\text{pref}}^{(k)}$. Here
119 $u^{(k)}$ is referred to as the preference distribution, which serves as a prior guess of $p^{(k)}$. Specifically
120 $u^{(k)}$ puts non-zero probability masses on M ‘‘template’’ vectors $\mathbf{t}_1^{(k)}, \mathbf{t}_2^{(k)}, \dots, \mathbf{t}_M^{(k)}$ in \mathbb{R}^d , and their
121 probabilities are respectively $u_1^{(k)}, u_2^{(k)}, \dots, u_M^{(k)}$ (which sum to 1). Collectively, we will denote the
122 set $\{\mathbf{t}_1^{(k)}, \mathbf{t}_2^{(k)}, \dots, \mathbf{t}_M^{(k)}\}$ of templates by $\mathbf{T}^{(k)}$.

123 The preference distribution $u^{(k)}$ is considered as a good approximation of $p^{(k)}$, in the sense that the
124 support of $p^{(k)}$ is contained in the set $\mathbf{T}^{(k)}$ of templates. Note that if \mathbb{R}^d is the word embedding space

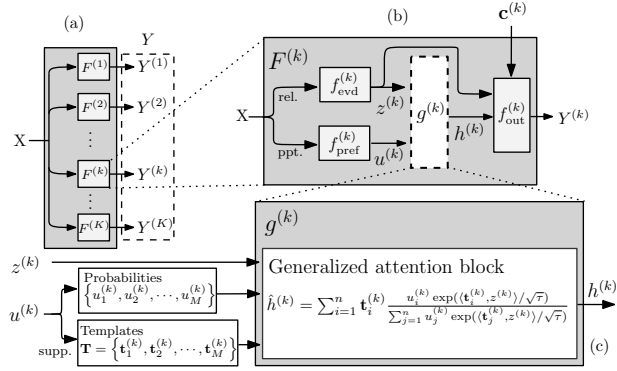


Figure 1: A conceptual graph of the deep learning model that we work with. The block $g^{(k)}$ is the one we will investigate.

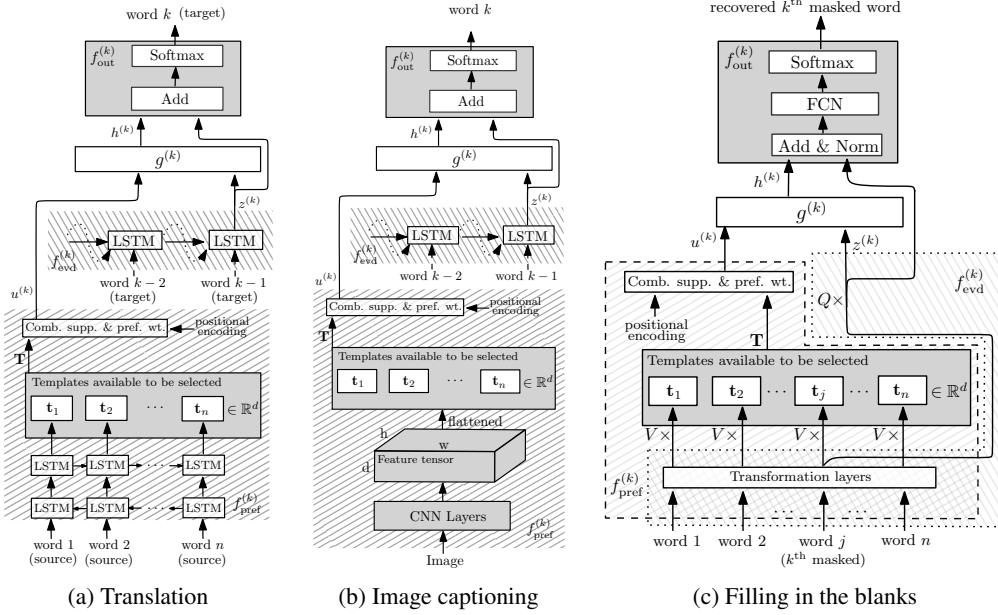


Figure 2: The model architectures of the three running examples. For the $f_{\text{evd}}^{(k)}$ in (a) and (b), the dashed links exist throughout the training and are replaced by the dotted ones in the generation stage.

125 for a large vocabulary, and if the size M of the template set $\mathbf{T}^{(k)}$ is relative small, then restricting the
 126 support of $p^{(k)}$ to within $\mathbf{T}^{(k)}$ imposes a strong constraint on $p^{(k)}$.

127 On the other hand, $u^{(k)}$ is not a sufficiently accurate approximation of $p^{(k)}$, in the sense that $u^{(k)}$ may
 128 assign probabilities to $\mathbf{T}^{(k)}$ somewhat differently. Such inaccuracy shifts the mean $\mu^{(k)}$ of $u^{(k)}$ from
 129 the mean $h^{(k)}$ of $p^{(k)}$. Suppose that there is another piece of information $z^{(k)} \in \mathbb{R}^d$ that is generated
 130 by another network module $f_{\text{evd}}^{(k)}$ and provides information regarding the mean shift. In particular, we
 131 assume that $z^{(k)}$ is a noisy version of the shift, more precisely,

$$z^{(k)} = h^{(k)} - \mu^{(k)} + \epsilon, \quad (5)$$

132 where $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ is the spherical Gaussian noise in \mathbb{R}^d with covariance $\sigma^2 \mathbf{I}$. We refer to $z^{(k)}$ as
 133 the evidence.

134 Then the design problem is to construct a function, or a network block, g , which infers the un-
 135 known distribution $p^{(k)}$ and hence its mean $h^{(k)}$ based on the evidence $z^{(k)}$ and the preference
 136 distribution $u^{(k)}$.

137 This formulation of the design problem might seem peculiar at the first glance, but we will show
 138 via examples (see Fig 2) that such a problem naturally arises in the construction of many machine
 139 learning models in practice.

140 **Example: Translation Problem.** For the translation problem, consider the model implementation
 141 plotted in Fig 2a that is similar to the one proposed in [1]. We will focus on the part of the model
 142 responsible for inferring the k^{th} word of the target sentence. In this model, $h^{(k)}$ corresponds to the
 143 constructed feature according to (4) that serves as an estimate of the context vector collecting the
 144 source sentence information. The estimated $h^{(k)}$ is then fed into a classifier $f_{\text{out}}^{(k)}$ to predict the k^{th}
 145 word. The preference distribution $u^{(k)}$ is generated by $f_{\text{pref}}^{(k)}$ which takes the source sentence words as
 146 inputs. In particular, the support of $u^{(k)}$ consists of the source sentence word embeddings \mathbf{T} (called
 147 annotations in [1]) which are pre-processed by two LSTM layers.¹ The preference weight for each
 148 template depends on some positional encoding functions, which, in principle, should assign higher

¹In this model, given input X , all $u^{(k)}$'s share the same support \mathbf{T} . The superscripts of the templates are then omitted to show their independence from k . Similar comments apply to implementations of the other two running examples.

149 weights to the templates appearing in the similar locations to the words we are inferring (that is, $h^{(k)}$
 150 is assumed to rely on the templates near \mathbf{t}_k more heavily).

151 Note that the inferred $p^{(k)}$'s support must be a subset of $u^{(k)}$'s as it is reasonable to assume that the
 152 target sentence words only depend on those appearing in the source sentence. Besides, although
 153 the preference weights specified by the positional encoding functions could provide some *a priori*
 154 information for the templates' weights in $p^{(k)}$, they cannot be accurate as their inferences do not
 155 consider the previously generated words $Y^{(i < t)}$. This results in the mean $\mu^{(k)}$ shifted from $h^{(k)}$,
 156 which is estimated by $z^{(k)} = f_{\text{evd}}^{(k)}$ that takes all the previously generated words $Y^{(i < t)}$ into account
 157 using another LSTM layer. Thus, $h^{(k)}$ and $p^{(k)}$ should not be far from $z^{(k)} + \mu^{(k)}$ and $u^{(k)}$,
 158 respectively.

159 **Example: Image Captioning.** The caption generation model presented in Fig 2b has a similar
 160 architecture reported in [29]. This model shares the designs of $f_{\text{evd}}^{(k)}$ and $f_{\text{out}}^{(k)}$ with the translation
 161 model while $f_{\text{pref}}^{(k)}$ instead extracts the templates from a raw image using a CNN network. In general,
 162 a word's position in the caption is independent of the location of the object it describes in the image.
 163 Therefore, in this model, all templates extracted by the CNN share the same preference weight.

164 As similar objects appear in an image would have similar features extracted by the CNN (for example,
 165 a zebra and a horse), allowing similar templates not in \mathbf{T} to participate in $h^{(k)}$'s estimation would
 166 possibly mix in information not contained in the raw image and harm the word inference accuracy.
 167 Therefore, we could improve the estimate of $h^{(k)}$ by choosing $p^{(k)}$ similar to $u^{(k)}$ in the sense that
 168 $p^{(k)}$'s support cannot contain elements not in $u^{(k)}$'s.

169 Intuitively, as the generation process proceeds, the context $h^{(k)}$ should be updated to provide relevant
 170 information in the image to facilitate the next word inference. Such change is governed by the
 171 caption's semantic evolution, which is captured by $z^{(k)} = f_{\text{evd}}^{(k)}$ that predicts the shift of the mean
 172 $\mu^{(k)}$ from $h^{(k)}$. For this reason, $\mu^{(k)} + z^{(k)}$ serves as an estimate of $h^{(k)}$ and should not be far away
 173 from it. Likewise, u^k should be close to $p^{(k)}$.

174 **Example: Filling in the blanks task.** For the filling-in-the-blank tasks, let us consider a model
 175 architecture plotted in Fig 2c that is similar to the one used in BERT [5]. We focus on the inference of
 176 the k^{th} masked word, which is assumed to be the j^{th} word of the input sentence. In this model, $f_{\text{pref}}^{(k)}$
 177 and $f_{\text{evd}}^{(k)}$ share the transformation layers (TL) that are commonly used in the NLP tasks to map one
 178 sequence of vector representations to another of the same length.² Taking the output sequence, $f_{\text{pref}}^{(k)}$
 179 applies a linear map V to each of its elements to form \mathbf{T} as the support of $u^{(k)}$ while the preference
 180 weights are specified by some positional encoding functions. At the same time, $z^{(k)} = f_{\text{evd}}^{(k)}$ estimates
 181 $h^{(k)}$'s shift from the mean $\mu^{(k)}$ due to the variation of the local information. For the same reasons
 182 discussed in the previous two examples, we need $\mu^{(k)} + z^{(k)}$ close to $h^{(k)}$ while $p^{(k)}$ is close to $u^{(k)}$.

183 Notably the formulation of the problem is based on the assumption that the network modules $f_{\text{evd}}^{(k)}$
 184 and $f_{\text{pref}}^{(k)}$ are fixed and generate $z^{(k)}$ and $u^{(k)}$ satisfying the above assumed properties. In reality,
 185 $f_{\text{evd}}^{(k)}$ and $f_{\text{pref}}^{(k)}$ are in fact obtained via training. However, we argue that if g is made to satisfy our
 186 design objective, then we can at least *interpret* $f_{\text{evd}}^{(k)}$ and $f_{\text{pref}}^{(k)}$ obtained from training as serving to
 187 produce $z^{(k)}$ and $u^{(k)}$ with our desired properties.

188 5 Formulation of an Optimization Problem

189 The discussion made in the previous section implies that the key optimization problem we are about
 190 to focus on should ensure

- 191 1. $h^{(k)}$ is not too far from $\mu^{(k)} + z^{(k)}$, where $h^{(k)}$ is constructed by $p^{(k)}$ according to (4) and
 192 $\mu^{(k)}$ is the mean of the preference distribution $u^{(k)}$.
- 193 2. $p^{(k)}$ is close to $u^{(k)}$ while $p^{(k)}$'s support must be a subset of $u^{(k)}$'s.

²Typical implementation of such layers include convolution layers, recurrent layers and self-attention layers.

194 These two desiderata prompt us to optimize:

$$\min_p \frac{\alpha}{2} \left\| \left(\mu^{(k)} + z^{(k)} \right) - \int_{\mathbb{R}^d} \mathbf{a} p(\mathbf{a}) \, d\mathbf{a} \right\|^2 + \mathcal{K}(p, u^{(k)}) \quad (6)$$

195 with $\alpha > 0$ responsible for the relative strength of the two terms, where $\mathcal{K}(p, u^{(k)})$ denotes the KL
 196 divergence from p to $u^{(k)}$. Remarkably, $\mathcal{K}(p, u^{(k)})$ has a finite value if and only if $p^{(k)}$ has non-zero
 197 values on the support of $u^{(k)}$. Thus, both requirements in the second desideratum are satisfied by
 198 using the KL divergence as a measure for the closeness of $p^{(k)}$ and $u^{(k)}$. Let $\tilde{p}^{(k)}$ be the minimizer of
 199 (6). The estimate of $h^{(k)}$ is

$$\hat{h}^{(k)} = \int_{\mathbb{R}^d} \mathbf{a} \tilde{p}^{(k)}(\mathbf{a}) \, d\mathbf{a}. \quad (7)$$

200 Naturally, this optimization problem can be derived from two different, though, related perspectives.³

201 **A maximum likelihood perspective.** The optimization problem in (6) can be derived using the
 202 maximum log likelihood method by treating the KL-divergence term as a regularizer. According to
 203 (5), the difference $(\mu^{(k)} + z^{(k)}) - h^{(k)}$ follows a Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. This implies the log
 204 likelihood function $\ell(z^{(k)}) \propto -\frac{1}{2\sigma^2} \left\| (\mu^{(k)} + z^{(k)}) - h^{(k)} \right\|^2$. Maximizing it with the KL-divergence
 205 term as a regularizer is the same as minimizing

$$\frac{1}{2\sigma^2} \left\| (\mu^{(k)} + z^{(k)}) - h^{(k)} \right\|^2 + \eta \mathcal{K}(p, u^{(k)}), \quad (8)$$

206 where $\eta > 0$ controls the strength of the regularization. Substituting (4) into (8) followed by
 207 rearrangement yields

$$\min_p \frac{1}{2\eta\sigma^2} \left\| (\mu^{(k)} + z^{(k)}) - \int_{\mathbb{R}^d} \mathbf{a} p(\mathbf{a}) \, d\mathbf{a} \right\|^2 + \mathcal{K}(p, u^{(k)}), \quad (9)$$

208 which is equivalent to (6) by setting $\alpha^{-1} = \eta\sigma^2$.

209 **A maximum entropy on the mean perspective** Consider a problem that seeks a distribution p such
 210 that the expectation $\int_{\mathbb{R}^d} \mathbf{a} p(\mathbf{a}) \, d\mathbf{a}$ is not far from $\mu^{(k)} + z^{(k)}$. In particular, we require

$$\left\| (\mu^{(k)} + z^{(k)}) - \int_{\mathbb{R}^d} \mathbf{a} p(\mathbf{a}) \, d\mathbf{a} \right\|^2 \leq \frac{1}{2\alpha}. \quad (10)$$

211 Note that, given $z^{(k)}$, there are infinitely many p 's that satisfy the constraints, which makes it difficult
 212 to pick a “best” p for later use. A technique known in information theory as the maximum entropy on
 213 the mean (MEM) [6, 21] solves this problem by picking the best guess of the ground truth p^* that
 214 simultaneously satisfies (10) and minimizes the KL divergence to the preference distribution $u^{(k)}$.
 215 That is,

$$\tilde{p}^{(k)} = \arg \min_p \mathcal{K}(p, u^{(k)}) \quad \text{subject to} \quad \left\| (\mu^{(k)} + z^{(k)}) - \int_{\mathbb{R}^d} \mathbf{a} p(\mathbf{a}) \, d\mathbf{a} \right\|^2 \leq \frac{1}{2\alpha}, \quad (11)$$

216 which is also the minimizer of (6) according to Equation (18) of [21] and Corollary 4.9 of [2].

217 6 Optimal Solution

218 Rioux et al. proved that the optimization problem stated in (6) has the following Fenchel dual (see
 219 Theorem 2 of [21]):

220 **Theorem 1.** *The dual of (6) is given by*

$$\max_{\lambda \in \mathbb{R}^d} \left\{ \left\langle \lambda, \mu^{(k)} + z^{(k)} \right\rangle - \frac{1}{2\alpha} \|\lambda\|^2 - \log \int_{\mathbb{R}^d} u^{(k)}(\mathbf{a}) \exp\langle \mathbf{a}, \lambda \rangle \, d\mathbf{a} \right\}. \quad (12)$$

221 *Given a maximizer λ^* of (12), one can recover the minimizer $\tilde{p}^{(k)}$ of (6) via*

$$\tilde{p}^{(k)}(\mathbf{a}) = \frac{u(\mathbf{a}) \exp\langle \mathbf{a}, \lambda^* \rangle}{\int_{\mathbb{R}^d} u(\mathbf{a}') \exp\langle \mathbf{a}', \lambda^* \rangle \, d\mathbf{a}'}. \quad (13)$$

³In fact, there is also a Bayesian perspective to derive the problem, which will be discussed in Appendix A.

222 By Theorem 1, the estimated $h^{(k)}$ defined in (7) can be re-written as

$$\hat{h}^{(k)} = \int_{\mathbb{R}^d} \mathbf{a} \tilde{p}^{(k)}(\mathbf{a}) \, d\mathbf{a} = \int_{\mathbb{R}^d} \mathbf{a} \frac{u^{(k)}(\mathbf{a}) \exp\langle \mathbf{a}, \lambda^* \rangle}{\int_{\mathbb{R}^d} u^{(k)}(\mathbf{a}') \exp\langle \mathbf{a}', \lambda^* \rangle \, d\mathbf{a}'} \, d\mathbf{a}, \quad (14)$$

223 where λ^* is a maximizer of (12).

224 In general, λ^* does not have a closed-form expression in terms of α , $u^{(k)}$ and $z^{(k)}$, and a standard
225 paradigm is to search for it using gradient ascent-based methods. In this paper, we will not search for
226 λ^* in this way; instead, we will derive a closed-form expression to approximate it. Remarkably, this
227 takes the form of the generalized attention presented in Fig 1.

228 Note that the integration in (12) equals $\mathbb{E}_{u^{(k)}}[\exp\langle W, \lambda \rangle]$, the expectation of the random variable
229 $\exp\langle W, \lambda \rangle$ where W has the probability distribution $u^{(k)}$. The expectation is just the moment
230 generating function (MGF) $M(\lambda)$ of W , and the value $\log M(\lambda)$ is called the cumulant of W [17,
231 p.26], which has an expansion [17, (2.4)]

$$\log M(\lambda) = \langle \mu^{(k)}, \lambda \rangle + \frac{1}{2} \langle \lambda, \Sigma^{(k)} \lambda \rangle + \mathcal{O}(\|\lambda\|^3), \quad (15)$$

232 where $\mu^{(k)} = \int_{\mathbb{R}^d} \mathbf{a} u^{(k)}(\mathbf{a}) \, d\mathbf{a}$ and $\Sigma^{(k)} = \int_{\mathbb{R}^d} (\mathbf{a} - \mu^{(k)}) (\mathbf{a} - \mu^{(k)})^T u^{(k)}(\mathbf{a}) \, d\mathbf{a}$ respectively
233 denote the expectation and the variance-covariance matrix of W .

234 Now we assume that α is small and we argue that this assumption is justified in practice. For instance,
235 in the translation task, all of words in the dictionary can serve as candidate templates, which could
236 be more than 10,000, but $u^{(k)}$ reduces this size to the length of the source sentence (usually less
237 than tens of words). The inference of $p^{(k)}$ should strongly anchor around this prior information;
238 consequently the information provided by $z^{(k)}$ should weigh less. On the other hand, $z^{(k)}$ can hardly
239 provide an accurate estimate of the mean shift, since the generation of $z^{(k)}$ is often ignorant of the
240 templates selected by $u^{(k)}$ (for example, in the example translation and image captioning models) or
241 generated by a low-capacity module (as in the example filling-in-the-blank model). For these reasons,
242 one should de-emphasize the constraint imposed by $z^{(k)}$ and hence choose a small α .

243 When α is picked to be small enough (see (12)), the optimization of λ gets a large penalty on its L2
244 norm and thus, $\|\lambda^*\|$ is close to zero. Then, by (15), we have

$$\log \int_{\mathbb{R}^d} u^{(k)}(\mathbf{a}) \exp\langle \mathbf{a}, \lambda^* \rangle \, d\mathbf{a} = \log M(\lambda^*) \approx \langle \mu^{(k)}, \lambda^* \rangle + \frac{1}{2} \langle \lambda^*, \Sigma^{(k)} \lambda^* \rangle. \quad (16)$$

245 Substituting (16) into (12) followed by setting the derivative with respect to λ to zero yields

$$\lambda^* = \alpha (I_d + \alpha \Sigma^{(k)})^{-1} z^{(k)}, \quad (17)$$

246 where I_d denotes the $d \times d$ identity matrix. As α is assumed close to zero, (17) is further reduced to

$$\lambda^* = \alpha z^{(k)}. \quad (18)$$

247 Plugging the expression into (14) gives the result stated as follows:

248 **Theorem 2.** *For a small enough $\alpha > 0$, the estimated $h^{(k)}$ defined in (7) can be approximated by*

$$\hat{h}^{(k)} = \int_{\mathbb{R}^d} \mathbf{a} \frac{u^{(k)}(\mathbf{a}) \exp(\alpha \langle \mathbf{a}, z^{(k)} \rangle)}{\int_{\mathbb{R}^d} u^{(k)}(\mathbf{a}') \exp(\alpha \langle \mathbf{a}', z^{(k)} \rangle) \, d\mathbf{a}'} \, d\mathbf{a}. \quad (19)$$

249 For the case that $u^{(k)}$ is a discrete distribution with support $\{\mathbf{t}_1^{(k)}, \mathbf{t}_2^{(k)}, \dots, \mathbf{t}_n^{(k)}\}$ and the preference
250 probability $\{u_1^{(k)}, u_2^{(k)}, \dots, u_n^{(k)}\}$, (19) becomes simply

$$\hat{h}^{(k)} = \sum_{i=1}^n \mathbf{t}_i \frac{u_i^{(k)} \exp(\alpha \langle \mathbf{t}_i, z^{(k)} \rangle)}{\sum_{j=1}^n u_j^{(k)} \exp(\alpha \langle \mathbf{t}_j, z^{(k)} \rangle)}. \quad (20)$$

251 In Fig 3, we set $d = 2$ and visualize the approximation of $h^{(k)}$ for different selections of α . We can
252 observe that, as α decreases, (20) outputs a better approximation of $\hat{h}^{(k)}$.

253 Let $\alpha = \tau^{-\frac{1}{2}}$, we rewrite Theorem 2 as follows for later reference.

254 **Corollary 1.** *For a sufficiently large τ , the best guess of $h^{(k)}$ defined in (7) with $\alpha = \tau^{-\frac{1}{2}}$ equals*

$$\hat{h}^{(k)} = \sum_{i=1}^n \mathbf{t}_i \frac{u_i^{(k)} \exp(\langle \mathbf{t}_i, z^{(k)} \rangle / \sqrt{\tau})}{\sum_{j=1}^n u_j^{(k)} \exp(\langle \mathbf{t}_j, z^{(k)} \rangle / \sqrt{\tau})}. \quad (21)$$

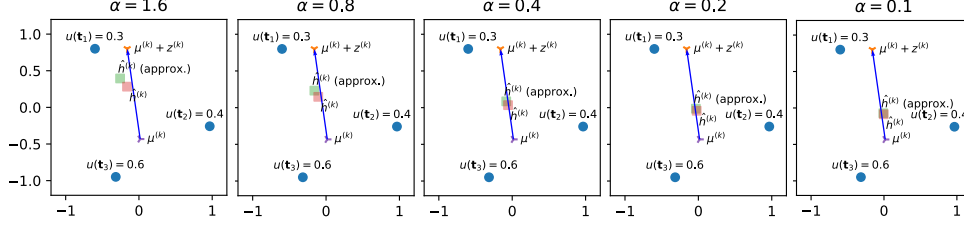


Figure 3: The approximation of $\hat{h}^{(k)}$ for different choices of α . The dots in light blue compose the support of discrete $u^{(k)}$ with the preference weights labelled above. The dark blue arrow starting from the mean $\mu^{(k)}$ of $u^{(k)}$ denotes the evidence $z^{(k)}$. The red square marks the $\hat{h}^{(k)}$ constructed by (14) with the λ^* maximizes (12), while the green one marks the $\hat{h}^{(k)}$ approximated by (20). As we can observe, (20) gives a precise approximation of $\hat{h}^{(k)}$ when α is sufficiently small.

255 7 Discussion

256 In Section 6, we derived an alternative expression of $\hat{h}^{(k)}$ defined in (7) by solving the Fenchel dual of
 257 the optimization problem stated in (6). Although the expression is not in closed form, as we are only
 258 interested in the case when α is small, a closed-form approximation of $\hat{h}^{(k)}$ is derived in Theorem 2
 259 and reduced to the form stated in (21) when considering a discrete distribution $u^{(k)}$.

260 As we pointed out, the block $g^{(k)}$ in Fig 2a, Fig 2b and Fig 2c is expected to find the inferred $\tilde{p}^{(k)}$
 261 minimizing (6) followed by plugging it into (7) to construct $\hat{h}^{(k)}$. Therefore, one can complete
 262 the architecture designs of the three running examples by replacing $g^{(k)}$ with a network layer
 263 implementing (21), namely, the structure in Figure 1 (c).

264 **The relationship between the optimal solution and the attention models.** Remarkably, the expres-
 265 sion stated in (21) gives a generalized attention block. By setting $u_i^{(k)} = \frac{1}{n}$ for all i , the expression
 266 is equivalent to the well known dot-product attention [15], which is also applied in the transformer
 267 network [25]. The equivalence of the expression of $\hat{h}^{(k)}$ and the dot-product attention layer tells us:
 268 (a) by applying a dot-product attention layer in a model, we essentially ask the model to perform
 269 an optimization task defined in (6) and construct the output according to (7). (b) the derivation of
 270 $h^{(k)}$ depends on two relatively independent pieces of information: a preference distribution given the
 271 global information and an estimate of the output’s deviation from the preference distribution’s mean
 272 according to some local information. This suggests that the design of attention-based model can be
 273 decomposed into two parts that respectively estimate these two values.

274 **The model consisting of a stack of attention layers.** Although our discussion focuses on the case
 275 that contains a single attention layer, any attention layer \mathcal{L} in an attention stack fits our frameworks
 276 (see Fig 1). In particular, all the attention layers closer to the input X than \mathcal{L} can be grouped into the
 277 functions $f_{\text{pref}}^{(k)}$ or $f_{\text{evd}}^{(k)}$. For those layers that take the current layer’s output as input, we can group
 278 them into $f_{\text{out}}^{(k)}$, where $\mathbf{c}^{(k)}$ may contain the outputs of other attention layers working in parallel.

279 **T5 transformer implicitly adopts the generalized attention structure.** We now show that T5
 280 transformer [19] can be seen as a realization of the generalized attention in (21), where the preference
 281 weights $u^{(k)}$ unifies the concepts of word masks and T5’s positional encoding functions. Consider
 282 the running example: filing in the blanks, with the preference distribution

$$u^{(k)}(\mathbf{t}_i) = \begin{cases} 0 & \text{if the } i^{\text{th}} \text{ word is masked} \\ \exp(b_{j-i})/Z & \text{otherwise,} \end{cases} \quad (22)$$

283 where Z is a normalizing constant and b_{j-i} is a trainable scalar that only depends on the relative
 284 position of word i and word j (which is the k^{th} masked word that we are inferring). Substituting
 285 such $u^{(k)}$ into (21) with $\tau = d$ yields

$$\hat{h}^{(k)} = \sum_{i=1}^n \mathbf{t}_i \frac{\exp\left(\frac{\langle \mathbf{t}_i, \mathbf{z}^{(k)} \rangle}{\sqrt{d}} + b_{j-i} + \mathbf{1}_{\text{masked}}(i)\right)}{\sum_{l=1}^n \exp\left(\frac{\langle \mathbf{t}_l, \mathbf{z}^{(k)} \rangle}{\sqrt{d}} + b_{j-l} + \mathbf{1}_{\text{masked}}(l)\right)}, \quad (23)$$

286 where $\mathbf{1}_{\text{masked}}(i)$ is an indicator function that equals $-\infty$ if word i is masked and zero otherwise.
287 The expression in (23) has the same structure as that adopted in T5 transformer, where the indicator
288 function serves as the mask function to prevent the model from assigning weights to the masked
289 words. In this way, the concepts of word masks and the positional encoding functions are unified
290 by $u^{(k)}$ in (22). Conversely, T5 transformer is a realization of the generalized attention with the
291 preference weights $u^{(k)}$ specified in (22).

292 **Generalized attention structures suggested by the optimal solution.** While T5 transformer has
293 implicitly adopted the generalized attention, (21) hints further generalizations could be made. For
294 instance, in T5 transformer, the function outputting template’s preference weights only considers the
295 word masks and the word’s relative positions. This function could be generalized to also consider the
296 input sentence contexts, and the output weights encode the importance of each word before giving
297 the local information stored in $z^{(k)}$. The same idea could be applied to the image captioning example
298 to replace the uniform preference weights. By adding a neural network taking the input image to
299 generate non-uniform preference weights, we devise a mechanism to estimate the importance of each
300 part of the image before the caption generation. In this way, the newly added network collects global
301 information from the image to propose a preference distribution, which could be updated locally
302 based on current generation stage encoded in $z^{(k)}$.

303 Moreover, although we mainly focus on the case when $u^{(k)}$ is discrete, we want to emphasize that the
304 analysis performed in Section 6 also covers continuous $u^{(k)}$. This hints that a continuous attention
305 mechanism could also be implemented, which might prove to be useful in some applications.

306 **Limitations and other comments.** The approximations performed in (15) and (18) have implicitly
307 assumed that random variable W following distribution $u^{(k)}$ has bounded moments. For a discrete $u^{(k)}$
308 with fixed support $\mathbf{T} = \{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_n\}$, all the moments are bounded and we can always pick a small
309 enough α (or equivalently large enough scaling factor τ in Cor 1) to make the approximation meet
310 our requirements. A concern may arise as the support \mathbf{T} in our running examples are supplied by
311 some neural layers, which could output templates of increasing norms as the training evolves. This
312 problem could be alleviated by adding norm regularization or using normalized templates instead.

313 8 Conclusion

314 This paper presented a new perspective to understand the attention mechanism by showing that it can
315 be treated as realizing a solver of a family of inference tasks. These tasks involve improving the noisy
316 estimate of a distribution p ’s mean by a preference distribution that encodes some beliefs of p ’s value.
317 We have used three running examples with the typical model architectures to show that such tasks
318 naturally exist in neural network design. We then abstracted a convex optimization problem from
319 these tasks and derived a closed-form approximation of the optimal solution by solving the problem’s
320 Fenchel dual. We find that the closed-form approximation can be seen as a generalized attention
321 layer and show that one of its special cases is equivalent to the dot-product attention adopted in
322 transformers. We further performed an analysis on the general form and showed that T5 transformer
323 implicitly adopts the generalized attention structure with attention weights unifying the concepts of
324 the word masks and the positional encoding functions.

325 This paper is the first work that presents a principled justification for the design of attention modules
326 in neural networks. The generalized attention structure presented in this paper potentially opens a door
327 to a wide design space. For example, the preference weights need not be derived from the positional
328 encoding functions; they could integrate a variety of information provided by other components of
329 the network. Additionally, this research might have pointed to new ways to analyze the functioning
330 of a neural network component, namely, via isolating the component from the complex network
331 structure and asking: is there a “local problem” that is solved by the design of this component?

332 **Potential negative societal impacts.** This paper presents a new perspective to understand attention
333 and derived a generalized attention structure. Our work is foundational, which we believe does not
334 have direct negative societal impacts. Due to the very wide range of applications of attention, such as
335 self-driving [12] and healthcare [16], our work may have unexpected negative impacts on these areas.

References

- 336
- 337 [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural Machine Translation by
338 Jointly Learning to Align and Translate. In *ICLR 2015*, pages 1–15, 2014.
- 339 [2] J. M. Borwein and A. S. Lewis. Partially finite convex programming, part i: Quasi relative
340 interiors and duality theory. *Mathematical Programming*, 57(1):15–48, 1992.
- 341 [3] Gino Brunner, Yang Liu, Damián Pascual, Oliver Richter, Massimiliano Ciaramita, and Roger
342 Wattenhofer. On Identifiability in Transformers. In *International Conference on Learning
343 Representations (ICLR)*, 2019.
- 344 [4] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What Does BERT
345 Look at? An Analysis of BERT’s Attention. *arXiv preprint arXiv:1906.04341*, 2019.
- 346 [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
347 deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Confer-
348 ence of the North American Chapter of the Association for Computational Linguistics: Human
349 Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis,
350 Minnesota, June 2019. Association for Computational Linguistics.
- 351 [6] Fabrice Gamboa. Methode du maximum d’entropie sur la moyenne et applications. *Phd Thesis*,
352 1989.
- 353 [7] Michael Hahn. Theoretical Limitations of Self-Attention in Neural Sequence Models. *Transac-
354 tions of the Association for Computational Linguistics*, 8:156–171, 2020.
- 355 [8] Jiri Hron, Yasaman Bahri, and Jascha Sohl-dickstein Roman. Infinite attention : NNGP and
356 NTK for deep attention networks. In *International Conference on Machine Learning (ICML)*,
357 2020.
- 358 [9] Sarthak Jain and Byron C. Wallace. Attention is not Explanation. In *Annual Conference of the
359 North American Chapter of the Association for Computational Linguistics: Human Language
360 Technologies (NAACL-HLT)*, 2019.
- 361 [10] E. T. Jaynes. Information theory and statistical mechanics. *Phys. Rev.*, 106:620–630, May 1957.
- 362 [11] E. T. Jaynes. Information theory and statistical mechanics. ii. *Phys. Rev.*, 108:171–190, Oct
363 1957.
- 364 [12] Jinkyu Kim and John Canny. Interpretable learning for self-driving cars by visualizing causal
365 attention. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2961–
366 2969, 2017.
- 367 [13] Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and
368 Yasaman Bahri. Deep neural networks as gaussian processes. In *International Conference on
369 Learning Representations (ICLR)*, 2018.
- 370 [14] Haoye Lu, Yongyi Mao, and Amiya Nayak. On the dynamics of training attention models. In
371 *International Conference on Learning Representations*, 2021.
- 372 [15] Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based
373 neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Nat-
374 ural Language Processing*, pages 1412–1421, Lisbon, Portugal, September 2015. Association
375 for Computational Linguistics.
- 376 [16] Fenglong Ma, Radha Chitta, Jing Zhou, Quanzeng You, Tong Sun, and Jing Gao. Dipole:
377 Diagnosis prediction in healthcare via attention-based bidirectional recurrent neural networks.
378 In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery
379 and Data Mining, KDD ’17*, page 1903–1911, New York, NY, USA, 2017. Association for
380 Computing Machinery.
- 381 [17] P. McCullagh. *Tensor Methods in Statistics : Monographs on Statistics and Applied Probability*.
382 Chapman and Hall/CRC, Boca Raton, FL, first edition. edition, 1987.

- 383 [18] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In
384 *Neural Information Processing Systems (NIPS)*, volume 32. Curran Associates, Inc., 2019.
- 385 [19] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
386 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified
387 text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- 388 [20] Hubert Ramsauer, Bernhard Schöfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Lukas
389 Gruber, Markus Holzleitner, Thomas Adler, David Kreil, Michael K Kopp, Günter Klambauer,
390 Johannes Brandstetter, and Sepp Hochreiter. Hopfield networks is all you need. In *International
391 Conference on Learning Representations*, 2021.
- 392 [21] Gabriel Rioux, Rustom Choksi, Tim Hoheisel, Pierre Marechal, and Christopher Scarvelis. The
393 maximum entropy on the mean method for image deblurring. *Inverse Problems*, oct 2020.
- 394 [22] R. Tyrrell Rockafellar. *Convex analysis*. Princeton mathematical series ; 28. Princeton University
395 Press, Princeton, N.J, 1970.
- 396 [23] Sofia Serrano and Noah A. Smith. Is attention interpretable? In *Annual Meeting of the
397 Association for Computational Linguistics (ACL)*, 2020.
- 398 [24] Shikhar Vashishth, Shyam Upadhyay, Gaurav Singh Tomar, and Manaal Faruqui. Attention
399 Interpretability Across NLP Tasks. In *International Conference on Learning Representations
400 (ICLR)*, 2020.
- 401 [25] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
402 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg,
403 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural
404 Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 405 [26] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head
406 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of
407 the 57th Annual Meeting of the Association for Computational Linguistics*, July 2019.
- 408 [27] Yequan Wang, Minlie Huang, Xiaoyan Zhu, and Li Zhao. Attention-based LSTM for aspect-
409 level sentiment classification. In *Empirical Methods in Natural Language Processing (EMNLP)*,
410 pages 606–615, Austin, Texas, November 2016. Association for Computational Linguistics.
- 411 [28] Sarah Wiegrefe and Yuval Pinter. Attention is not not explanation. In *2019 Conference on
412 Empirical Methods in Natural Language Processing and 9th International Joint Conference on
413 Natural Language Processing (EMNLP-IJCNLP)*, 2019.
- 414 [29] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show,
415 Attend and Tell: Neural Image Caption Generation with Visual Attention. In *Proc. 32nd Int.
416 Conf. Mach. Learn.*, pages 257–261, 2015.
- 417 [30] Greg Yang. Tensor Programs I : Wide Feedforward or Recurrent Neural Networks of Any
418 Architecture are Gaussian Processes. In *Neural Information Processing Systems (NIPS)*, 2019.

419 Checklist

- 420 1. For all authors...
- 421 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
422 contributions and scope? [Yes]
- 423 (b) Did you describe the limitations of your work? [Yes] *We have created a sub-section,
424 called **Limitations and other comments**, in Section 7 to discuss the potential limita-
425 tions of our work.*
- 426 (c) Did you discuss any potential negative societal impacts of your work? [Yes] *We have
427 discussed this at the end of the paper.*
- 428 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
429 them? [Yes]

- 430 2. If you are including theoretical results...
- 431 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 432 (b) Did you include complete proofs of all theoretical results? [Yes]
- 433 3. If you ran experiments...
- 434 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
- 435 imental results (either in the supplemental material or as a URL)? [N/A] *Our work*
- 436 *presents a theoretical framework and does not contain any experimental study.*
- 437 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 438 were chosen)? [N/A]
- 439 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 440 ments multiple times)? [N/A]
- 441 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 442 of GPUs, internal cluster, or cloud provider)? [N/A]
- 443 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 444 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 445 (b) Did you mention the license of the assets? [N/A]
- 446 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 447
- 448 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 449 using/curating? [N/A]
- 450 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 451 information or offensive content? [N/A]
- 452 5. If you used crowdsourcing or conducted research with human subjects...
- 453 (a) Did you include the full text of instructions given to participants and screenshots, if
- 454 applicable? [N/A]
- 455 (b) Did you describe any potential participant risks, with links to Institutional Review
- 456 Board (IRB) approvals, if applicable? [N/A]
- 457 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 458 spent on participant compensation? [N/A]