

Deep Neural Networks Evolve Human-like Attention Distribution during Goal-directed Reading Comprehension

Jiajie Zou

Zhejiang university

Nai Ding (✉ ding_nai@zju.edu.cn)

College of Biomedical Engineering and Instrument Sciences, Zhejiang University

Article

Keywords: deep neural networks, reading comprehension

Posted Date: August 26th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-813994/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Deep Neural Networks Evolve Human-like Attention Distribution**
2 **during Goal-directed Reading Comprehension**

3
4 Jiajie Zou², Nai Ding^{1,2*},

5
6 ¹Zhejiang lab; College of Biomedical Engineering and Instrument Sciences,
7 Zhejiang University, Hangzhou 311121, China

8 ²Key Laboratory for Biomedical Engineering of Ministry of Education,
9 Zhejiang University, Hangzhou 310027, China

10
11
12
13 ***Corresponding author:**

14 Nai Ding,

15 Email: ding_nai@zju.edu.cn

16 Zhejiang lab; College of Biomedical Engineering and Instrument Sciences,
17 Zhejiang University, Hangzhou 311121, China

18 **Abstract**

19 Attention is a key mechanism for information selection in both biological brains and
20 many state-of-the-art deep neural networks (DNNs). Here, we investigate whether
21 humans and DNNs allocate attention in comparable ways when seeking information
22 in a text passage to answer a question. We analyze 3 transformer-based DNNs that
23 reach human-level performance when trained to perform the reading comprehension
24 task. We find that the DNN attention distribution quantitatively resembles human
25 attention distribution measured by eye tracking: Human readers fixate longer on
26 words that are more relevant to the question-answering task, demonstrating that
27 attention is modulated by the top-down reading goal, on top of lower-level visual
28 layout and textual features. Further analyses reveal that the attention weights in DNNs
29 are also influenced by both the top-down reading goal and lower-level textual
30 features, with the shallow layers more strongly influenced by lower-level textual
31 features and the deep layers attending more to task-relevant words. Additionally, deep
32 layers' attention to task-relevant words gradually emerges when pre-trained DNN
33 models are fine-tuned to perform the reading comprehension task, which coincides
34 with the improvement in task performance. These results demonstrate that DNNs can
35 naturally evolve human-like attention distribution through task optimization. The
36 results suggest that human attention during goal-directed reading comprehension is a
37 consequence of task optimization and the attention weights in DNN are of biological
38 significance.

39

40

41 **Introduction**

42 Artificial intelligence (AI) and cognitive science separately investigate how machines
43 and brains solve complex information processing problems, such as language
44 comprehension and visual object recognition. As the artificial neural network approach
45 in AI was in part inspired by biological neural networks, there is continuing interest in
46 comparing the performance of artificial and biological neural networks¹⁻¹¹. If human-
47 or animal-like behaviors emerge in artificial neural networks, it indicates that the
48 computations implemented in artificial neural networks can serve as a possible model
49 for human/animal cognition¹²⁻¹⁵. Indeed, it has been found that artificial neurons in
50 modern DNNs can evolve receptive field properties that are comparable to those
51 measured from animal visual cortices¹⁶, and DNN models that have properties more
52 consistent with biological neural systems tend to perform better at information
53 processing tasks^{10,11}.

54
55 Traditional artificial neural networks mainly mimic lower-level or biophysical
56 properties of neurons, while the new generations of DNN models also attempt to mimic
57 high-level cognitive functions, e.g., attention. Attention mechanisms have greatly
58 improved the performance of DNNs and have become a necessary component in state-
59 of-the-art DNN models, especially in the field of natural language processing (NLP)<sup>17-
60 21</sup>. Recent studies have shown that the attention mechanism in DNN can play a wide
61 variety of roles in language processing, e.g., to extract task relevant information^{22,23} and
62 to analyze syntactic dependencies and semantic co-reference²⁴⁻²⁶. The attention
63 mechanism in DNNs, however, is not designed to quantitatively simulate human

64 attention, and few studies have systematically compared human and DNN attention
65 during the same language processing task (see Bolotova et al. 2020²⁷ for a notable
66 exception). Therefore, it remains unclear to what extent the attention mechanisms in
67 DNN language models are comparable to human attention and whether the attention
68 mechanisms in DNN can serve as a model for human attention.

69

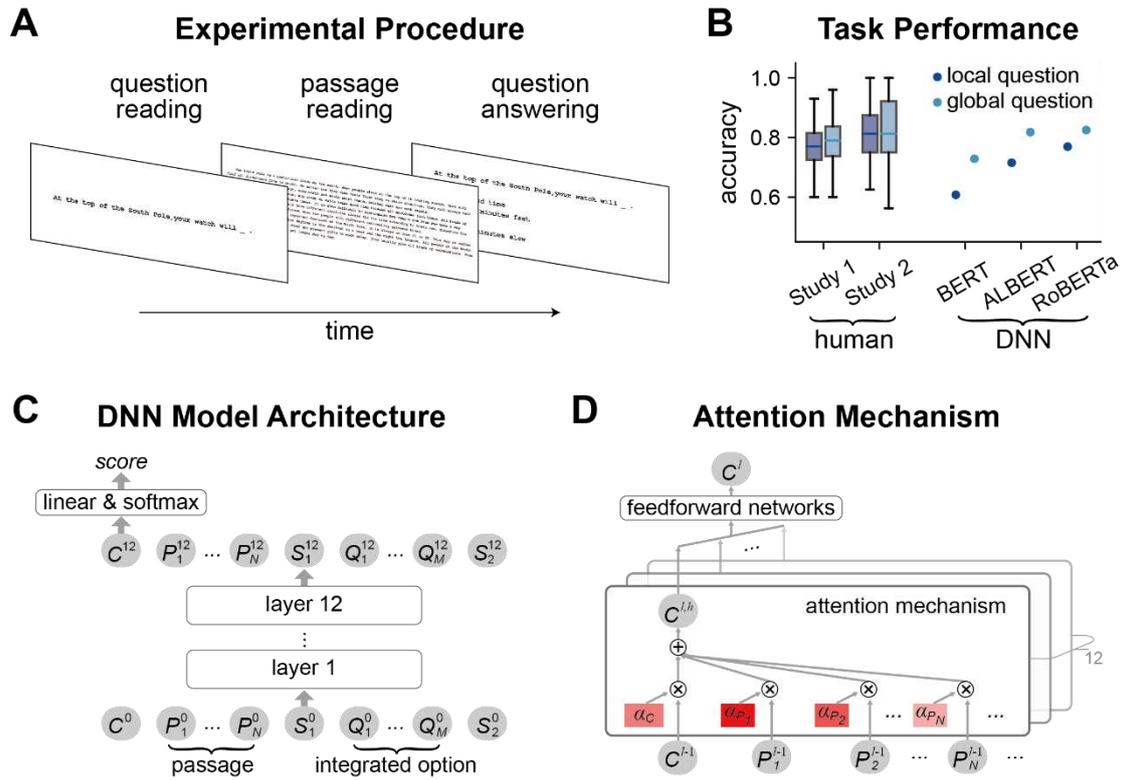
70 The human attention system has multiple components which contribute differently to
71 different tasks. For example, when freely viewing an image, attention is primarily
72 modulated by visual saliency, and this kind of attention is referred to as bottom-up
73 attention^{28,29}. When searching for a target object in a visual scene, however, viewers
74 attend more to possible locations of the target and objects with visual features more
75 consistent with the target³⁰. This form of attention - induced by the task - is called top-
76 down attention^{28,29}. In visual perception tasks that mainly engage bottom-up attention,
77 a large number of studies have shown that neural networks can be trained to model
78 human attention distribution measured through eye tracking^{31,32}. Recently, some
79 models have also been proposed to model top-down attention^{9,33-36}. In the domain of
80 language processing, computational models have been proposed to predict human
81 readers' eye movements when they read simple sentences without a specific purpose³⁷,
82 a task similar to free viewing. To our knowledge, however, no model has been proposed
83 to predict human attention when readers read a passage with a specific goal, e.g., to
84 answer a question, although goal-directed reading was the most common reading
85 behavior for adults³⁸.

86

87 Here we compare the attention distribution for humans and DNNs during a reading
88 comprehension task in which humans or DNNs have to answer a question by reading a
89 passage. We select this task since it is a benchmark task to test NLP algorithms^{39,40} and
90 also a common task to test human verbal ability, e.g., in exams such as SAT, GRE, and
91 TOEFL. This task is also suitable to investigate attention, since a passage contains an
92 enormous amount of information, but only a small portion of it is typically relevant to
93 answering a specific question, imposing a strong load for information selection. Finally,
94 state-of-the-art DNN models have recently achieved human-level performance on the
95 reading comprehension task for questions at the difficulty level corresponding to high
96 school exams in China^{18,19,21}.

97

98 With the reading comprehension task, we quantified the attention mechanisms in
99 humans using the fixation time and the attention mechanisms in DNNs using the
100 attention weight on each word. Note that both human eye fixations and the attention
101 weights in DNNs reflect intermediate processing steps instead of the outcome of
102 reading comprehension. We aim to investigate three closely related questions by
103 analyzing and comparing human and DNN attention. First, how is human and DNN
104 attention modulated by stimulus features and the top-down reading goal, i.e., the need
105 to answer a specific question? Second, do humans and DNNs show similar attention
106 distribution? Third, how does the DNN attention distribution evolve during training
107 and how does it relate to task performance?



108

109

110 **Fig. 1. Experimental procedure and DNN model.** (A) The experimental procedure in
 111 Study 1. In each trial, participants read a question first, and then read the corresponding
 112 passage, and finally proceed to read the question, coupled with 4 options, and answer
 113 it. (B) Performance of humans and DNN models on the reading comprehension task.
 114 (C) Architecture of the DNN models used for the reading comprehension task. The
 115 input to the models consists of all words in the passage and an integrated option, and
 116 also 3 special tokens, i.e., CLS, SEP₁, and SEP₂ (denoted as C, S₁, and S₂). The CLS
 117 token integrates information across words and is used to calculate a score that reflects
 118 how likely the option is the correct answer. The DNN model has 12 layers and has 12
 119 attention heads in each layer. (D) Illustration of the DNN attention mechanism in a layer.
 120 In the models, each word/token is represented by a vector, and information is integrated
 121 across words/tokens only in the self-attention module. For example, the vectorial
 122 representation of the CLS token is a weighted sum of the vectorial representations of
 123 all words and tokens. The attention weight for each word in the passage, i.e., α_{P_n} , is the
 124 DNN attention analyzed in this study. Output of the self-attention model, i.e., $C^{l,h}$, is
 125 further processed by feedforward networks and other operations that do not engage
 126 information integration across words.

127

Examples of Human Attention and Prediction of Different Features

A Local question ("At the top of the South Pole, your watch will __.")

human attention density prediction based on textual features

The South Pole is a particular place on the earth. When people stand at the top of it looking around, they will find all directions face to north. No matter how they make their first step in which direction, they will always walk towards the north. That's to say, only north and south exist there, neither east nor west exists.

At the top of the South Pole, any clock or watch keeps good time because all meridians join there. All kinds of local time are completely suitable there. It is even difficult to distinguish New Year's Eve from New Year's Day.

The explorers and scientists from different countries always fix the time according to their own. Therefore the time by their watches was different when the people with different nationality gathered there.

The Winter Solstice is an important festival at the South Pole. It is always on June 21 or 22. This day is called Midwinter Festival, on which the daytime is the shortest in a year and the night the longest. All people at the South Pole extend greetings to each other and present gifts to each other. They usually give all kinds of celebrations. From that day on, the daytime will get longer day by day.

The South Pole is a particular place on the earth. When people stand at the top of it looking around, they will find all directions face to north. No matter how they make their first step in which direction, they will always walk towards the north. That's to say, only north and south exist there, neither east nor west exists.

At the top of the South Pole, any clock or watch keeps good time because all meridians join there. All kinds of local time are completely suitable there. It is even difficult to distinguish New Year's Eve from New Year's Day.

The explorers and scientists from different countries always fix the time according to their own. Therefore the time by their watches was different when the people with different nationality gathered there.

The Winter Solstice is an important festival at the South Pole. It is always on June 21 or 22. This day is called Midwinter Festival, on which the daytime is the shortest in a year and the night the longest. All people at the South Pole extend greetings to each other and present gifts to each other. They usually give all kinds of celebrations. From that day on, the daytime will get longer day by day.

prediction based on layout features

prediction based on task relevance

The South Pole is a particular place on the earth. When people stand at the top of it looking around, they will find all directions face to north. No matter how they make their first step in which direction, they will always walk towards the north. That's to say, only north and south exist there, neither east nor west exists.

At the top of the South Pole, any clock or watch keeps good time because all meridians join there. All kinds of local time are completely suitable there. It is even difficult to distinguish New Year's Eve from New Year's Day.

The explorers and scientists from different countries always fix the time according to their own. Therefore the time by their watches was different when the people with different nationality gathered there.

The Winter Solstice is an important festival at the South Pole. It is always on June 21 or 22. This day is called Midwinter Festival, on which the daytime is the shortest in a year and the night the longest. All people at the South Pole extend greetings to each other and present gifts to each other. They usually give all kinds of celebrations. From that day on, the daytime will get longer day by day.

The South Pole is a particular place on the earth. When people stand at the top of it looking around, they will find all directions face to north. No matter how they make their first step in which direction, they will always walk towards the north. That's to say, only north and south exist there, neither east nor west exists.

At the top of the South Pole, any clock or watch keeps good time because all meridians join there. All kinds of local time are completely suitable there. It is even difficult to distinguish New Year's Eve from New Year's Day.

The explorers and scientists from different countries always fix the time according to their own. Therefore the time by their watches was different when the people with different nationality gathered there.

The Winter Solstice is an important festival at the South Pole. It is always on June 21 or 22. This day is called Midwinter Festival, on which the daytime is the shortest in a year and the night the longest. All people at the South Pole extend greetings to each other and present gifts to each other. They usually give all kinds of celebrations. From that day on, the daytime will get longer day by day.

prediction based on DNN attention

DNN attention in the last layer

The South Pole is a particular place on the earth. When people stand at the top of it looking around, they will find all directions face to north. No matter how they make their first step in which direction, they will always walk towards the north. That's to say, only north and south exist there, neither east nor west exists.

At the top of the South Pole, any clock or watch keeps good time because all meridians join there. All kinds of local time are completely suitable there. It is even difficult to distinguish New Year's Eve from New Year's Day.

The explorers and scientists from different countries always fix the time according to their own. Therefore the time by their watches was different when the people with different nationality gathered there.

The Winter Solstice is an important festival at the South Pole. It is always on June 21 or 22. This day is called Midwinter Festival, on which the daytime is the shortest in a year and the night the longest. All people at the South Pole extend greetings to each other and present gifts to each other. They usually give all kinds of celebrations. From that day on, the daytime will get longer day by day.

The South Pole is a particular place on the earth. When people stand at the top of it looking around, they will find all directions face to north. No matter how they make their first step in which direction, they will always walk towards the north. That's to say, only north and south exist there, neither east nor west exists.

At the top of the South Pole, any clock or watch keeps good time because all meridians join there. All kinds of local time are completely suitable there. It is even difficult to distinguish New Year's Eve from New Year's Day.

The explorers and scientists from different countries always fix the time according to their own. Therefore the time by their watches was different when the people with different nationality gathered there.

The Winter Solstice is an important festival at the South Pole. It is always on June 21 or 22. This day is called Midwinter Festival, on which the daytime is the shortest in a year and the night the longest. All people at the South Pole extend greetings to each other and present gifts to each other. They usually give all kinds of celebrations. From that day on, the daytime will get longer day by day.

attention density min max

B Global question ("What is the passage mainly about?")

human attention density prediction based on textual features

All the oceans of the world will be dead in the future unless action is taken at once. How can this happen?

We have already seen that people allow all sorts of waste products to flow into the sea. It is almost impossible to measure how much waste water and industrial waste in our oceans, but we can find out how much oil is poured into them legally and illegally. It is illegal to pour oil into the sea close to the shore, but when a ship is many miles out to sea, there are no such rules. Scientists have different ideas as to the amount of oil putting into the sea, but the lowest figure for oil poured in European waters alone is nearly 200 000 tons every year. Some people say the figure could be ten times as high.

It is not only our coasts that suffer from oil pollution. Many shell fish, for example, now have high amounts of poisonous substances. Next time you have shell fish to eat, how can you be sure that they are free from oil pollution? You cannot see the effects and you cannot taste them, either. It is really quite a problem.

All the oceans of the world will be dead in the future unless action is taken at once. How can this happen?

We have already seen that people allow all sorts of waste products to flow into the sea. It is almost impossible to measure how much waste water and industrial waste in our oceans, but we can find out how much oil is poured into them legally and illegally. It is illegal to pour oil into the sea close to the shore, but when a ship is many miles out to sea, there are no such rules. Scientists have different ideas as to the amount of oil putting into the sea, but the lowest figure for oil poured in European waters alone is nearly 200 000 tons every year. Some people say the figure could be ten times as high.

It is not only our coasts that suffer from oil pollution. Many shell fish, for example, now have high amounts of poisonous substances. Next time you have shell fish to eat, how can you be sure that they are free from oil pollution? You cannot see the effects and you cannot taste them, either. It is really quite a problem.

prediction based on layout features

prediction based on task relevance

All the oceans of the world will be dead in the future unless action is taken at once. How can this happen?

We have already seen that people allow all sorts of waste products to flow into the sea. It is almost impossible to measure how much waste water and industrial waste in our oceans, but we can find out how much oil is poured into them legally and illegally. It is illegal to pour oil into the sea close to the shore, but when a ship is many miles out to sea, there are no such rules. Scientists have different ideas as to the amount of oil putting into the sea, but the lowest figure for oil poured in European waters alone is nearly 200 000 tons every year. Some people say the figure could be ten times as high.

It is not only our coasts that suffer from oil pollution. Many shell fish, for example, now have high amounts of poisonous substances. Next time you have shell fish to eat, how can you be sure that they are free from oil pollution? You cannot see the effects and you cannot taste them, either. It is really quite a problem.

All the oceans of the world will be dead in the future unless action is taken at once. How can this happen?

We have already seen that people allow all sorts of waste products to flow into the sea. It is almost impossible to measure how much waste water and industrial waste in our oceans, but we can find out how much oil is poured into them legally and illegally. It is illegal to pour oil into the sea close to the shore, but when a ship is many miles out to sea, there are no such rules. Scientists have different ideas as to the amount of oil putting into the sea, but the lowest figure for oil poured in European waters alone is nearly 200 000 tons every year. Some people say the figure could be ten times as high.

It is not only our coasts that suffer from oil pollution. Many shell fish, for example, now have high amounts of poisonous substances. Next time you have shell fish to eat, how can you be sure that they are free from oil pollution? You cannot see the effects and you cannot taste them, either. It is really quite a problem.

prediction based on DNN attention

DNN attention in the last layer

All the oceans of the world will be dead in the future unless action is taken at once. How can this happen?

We have already seen that people allow all sorts of waste products to flow into the sea. It is almost impossible to measure how much waste water and industrial waste in our oceans, but we can find out how much oil is poured into them legally and illegally. It is illegal to pour oil into the sea close to the shore, but when a ship is many miles out to sea, there are no such rules. Scientists have different ideas as to the amount of oil putting into the sea, but the lowest figure for oil poured in European waters alone is nearly 200 000 tons every year. Some people say the figure could be ten times as high.

It is not only our coasts that suffer from oil pollution. Many shell fish, for example, now have high amounts of poisonous substances. Next time you have shell fish to eat, how can you be sure that they are free from oil pollution? You cannot see the effects and you cannot taste them, either. It is really quite a problem.

All the oceans of the world will be dead in the future unless action is taken at once. How can this happen?

We have already seen that people allow all sorts of waste products to flow into the sea. It is almost impossible to measure how much waste water and industrial waste in our oceans, but we can find out how much oil is poured into them legally and illegally. It is illegal to pour oil into the sea close to the shore, but when a ship is many miles out to sea, there are no such rules. Scientists have different ideas as to the amount of oil putting into the sea, but the lowest figure for oil poured in European waters alone is nearly 200 000 tons every year. Some people say the figure could be ten times as high.

It is not only our coasts that suffer from oil pollution. Many shell fish, for example, now have high amounts of poisonous substances. Next time you have shell fish to eat, how can you be sure that they are free from oil pollution? You cannot see the effects and you cannot taste them, either. It is really quite a problem.

128

129 **Fig. 2. Examples of the human attention distribution and the human attention**
130 **distribution predicted by different features.** Panels A and B separately show the
131 attention distribution for two passages and the corresponding questions are shown in
132 the parenthesis. Human attention is quantified by the total fixation time per unit area.
133 Textual features include word properties, e.g., word frequency and the position of a
134 word in the passage. Layout features include visual features that can be processed
135 without recognizing individual words. Task relevance contains human annotation about
136 the contribution of each word to question answering. DNN attention includes all the
137 layers and attention heads, and the DNN attention in the last layer is shown separately
138 (averaged over attention heads).

139 **Results**

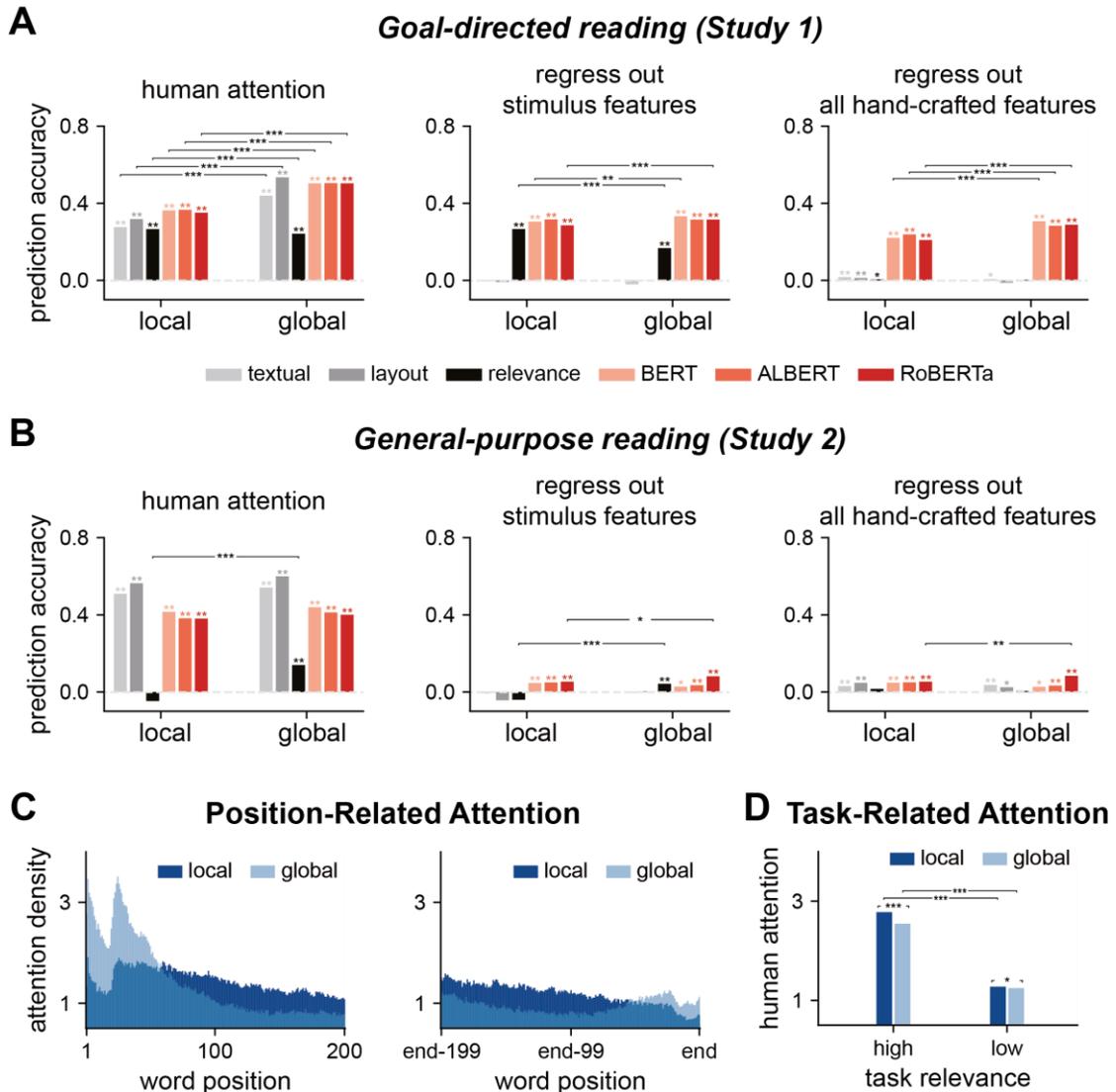
140 **Human Attention Distribution and Influence Factors**

141 In Study 1, the participants ($N = 25$ for each question) first read a question and then
142 read a passage based on which the question should be answered. After reading the
143 passage, the participants read 4 options related to the question and had to choose which
144 option was the most suitable answer. Eight hundred question/passage pairs were
145 presented, and the questions fell into two broad categories, i.e., local and global
146 questions (see *Materials and Methods* for details). Local questions require attention to
147 details while global questions concern the general understanding of a passage. The
148 participants correctly answered 77.94% questions on average (Fig. 1B, 77.49% and
149 78.77% for local and global questions, respectively).

150

151 While the participants read the passage, their eye gaze was monitored using an eye
152 tracker, and their attention to each word was quantified by the total fixation time on the
153 word. The results showed that longer words were fixated for longer time (Fig. S1),
154 consistent with previous studies⁴¹. Nevertheless, the fixation time on a word was
155 expected to be positively related with the area the word occupied even when attention
156 was uniformly distributed across the visual field. Therefore, here we further extracted
157 the *attention density* by dividing the total fixation time on a word by the area the word
158 occupied, and used this measure in subsequent analyses. The attention density clearly
159 deviated from a uniform distribution (see Fig. 2 for examples). To probe into the factors
160 modulating human attention distribution, we quantified how the human attention
161 distribution was influenced by multiple sets of features in the following.

Predicting Human Attention Based on Different Features



162

163

164 **Fig. 3. Predicting human attention using different features.** (A and B) Panels A and

165 B show the results of Study 1 and Study 2, respectively. The left plots show how well

166 different sets of features can predict human attention. In the middle and right plots,

167 some features are regressed out from human attention, and the residual human attention

168 is predicted by other features. Prediction accuracy that is significantly higher than

169 chance is denoted by stars of the same color as the bar. (C) The influence of word

170 position on human attention. Humans generally attend more to the beginning of a

171 passage, especially for global questions. (D) The influence of task relevance on human

172 attention. Humans allocate more attention to words that are more relevant to question

173 answering. * $P < 0.05$; ** $P < 0.01$; *** $P < 0.001$.

174

175 We first analyzed whether textual features, e.g., word length, word frequency, and a
176 word's position in a sentence, could predict human attention distribution using linear
177 regression. The prediction accuracy, i.e., the correlation coefficient between the
178 predicted and actual attention density, was significantly above chance ($P = 0.002$,
179 permutation test, FDR corrected). Furthermore, the prediction accuracy was
180 significantly higher for global questions than for local questions ($P = 1.4 \times 10^{-4}$,
181 bootstrap, FDR corrected) (Fig. 3A, the left plot). We then used the same regression
182 analysis to analyze whether the visual layout of a passage could also influence attention
183 distribution. Here, layout features referred to features induced by line changes (see
184 *Materials and Methods* for details), which could be processed without word recognition.
185 The prediction accuracy for layout features was also statistically significant ($P = 0.002$,
186 permutation test, FDR corrected).

187

188 Textual features and layout features characterized properties of the stimulus that were
189 invariant across tasks. In the following, we investigated whether the task, i.e., to answer
190 a specific question, also modulated human attention distribution. To characterize the
191 top-down influence of task, we acquired annotations indicating each word's
192 contribution to question answering, i.e., task relevance (see *Materials and Methods*).
193 As shown in the left plot of Fig. 3A, we found that task relevance could indeed
194 significantly predict human attention distribution ($P = 0.002$, permutation test, FDR
195 corrected). Since task relevance was not a well-established modulator of reading
196 attention, we further analyzed whether the task relevance effect could be explained by

197 the well-established textual and layout effects. In this analysis, we first regressed out
198 the influence of textual and layout features from the human attention distribution, and
199 found that the residual attention distribution could still be predicted by task relevance
200 ($P = 0.003$, permutation test, FDR corrected) (Fig. 3A, middle plot). These results
201 showed that the top-down reading goal, quantified by task relevance, could modulate
202 human attention, on top of lower-level stimulus features, i.e., textual and layout features.
203

204 The linear regression analyses revealed that textual features, layout features, and task
205 relevance all modulated human attention (see Fig. 2 for examples). The prediction
206 accuracy for different features ranged between 0.2 and 0.6, comparable to the prediction
207 accuracy of visual saliency models when predicting human attention to images^{31,32}.
208 Further analyses also revealed how these features modulated human attention. For
209 example, we found that participants generally attended more to the beginning of a
210 passage (Fig. 3C). Furthermore, this effect was stronger for global questions, which
211 potentially explained why stimulus features could better predict the attention
212 distribution for global questions. Additionally, it was also found that participants
213 attended more to words that are more relevant to the question answering task (Fig. 3D).

214

215 **Attention Distributions in Humans and DNN**

216 We then investigated whether DNN attention was comparable to human attention. The
217 general architecture of the models was illustrated in Fig. 1C. The input to the models
218 included all the words in the passage, integrated option, and 3 special tokens. One of

219 the special token, i.e., CLS, was the decision variable, based on the final representation
220 of which the DNN models decided whether an option was the correct answer or not. In
221 the following, we analyzed the attention weight between the CLS token and each word
222 in the passage (see *Materials and Methods* for details). In each layer of the DNN models,
223 the vectorial representation of the CLS token was updated by a weighted sum of the
224 vectorial representations of all input words and tokens. Therefore, the attention weight
225 on a word could reflect how heavily the word contributed to the decision variable, i.e.,
226 the CLS token.

227

228 We analyzed 3 DNN models, i.e., BERT¹⁷, ALBERT¹⁸, and RoBERTa¹⁹, and the
229 question answering performance of the 3 DNN models was within the range of human
230 performance (Fig. 1B). Each of the 3 DNN models had 12 layers and each layer had 12
231 heads, each of which had a separate set of attention weights (Fig. 1CD). Consequently,
232 each word had 144 attention weights (12 layers \times 12 heads). In the following, we first
233 tested whether the DNNs learned human-like attention distributions by attempting to
234 decode human attention distribution from the 144 DNN attention weights using linear
235 regression. Then, we analyzed whether the attention weights in different layers showed
236 different properties.

237

238 Although the DNN models were only trained to perform the reading comprehension
239 task and were blind to the human fixation data, it was found that the DNN attention
240 weights could significantly predict human attention distribution ($P = 0.002$, permutation

241 test, FDR corrected), and the prediction accuracy was higher for global questions than
242 for local questions ($P = 1.4 \times 10^{-4}$, bootstrap, FDR corrected) (Fig. 3A, left plot). The
243 prediction accuracy of DNN attention weights was higher than that of textual features
244 and task relevance. When compared with the predictions based on layout features, the
245 predictions based on DNN attention weights were higher for local questions and lower
246 for global questions. It should be mentioned, however, that layout features, which were
247 induced by line changes, were not available in the input to DNN models.

248

249 DNN attention weights could model the human attention distribution, but did they
250 capture information beyond the hand-crafted features, i.e., textual features, layout
251 features, and task relevance features? We found that when the influences of textual and
252 layout features were regressed out, the residual human attention distribution could still
253 be explained by the DNN attention weights (Fig. 3A, the middle plot). This result
254 suggested that the DNN attention weights contained information beyond basic stimulus
255 features. Additionally, when the stimulus features and task relevance features were both
256 regressed out, the residual human attention distribution remained significantly
257 predicted by the DNN attention weights (Fig. 3A, the right plot). Therefore, DNN
258 attention weights could model human attention and capture information beyond basic
259 hand-crafted features.

260

261 **Task Modulation in Humans**

262 To further confirm that human attention received top-down modulation from the task,
263 we conducted Study 2 as a control study. In Study 2, participants first read a passage

264 without prior knowledge about the specific question to answer. After the first-pass
265 passage reading, the participants read the question and were then allowed to read the
266 passage again before answering the question. We analyzed the attention density during
267 the first-pass reading of the passage, which was referred to as general-purpose reading.

268

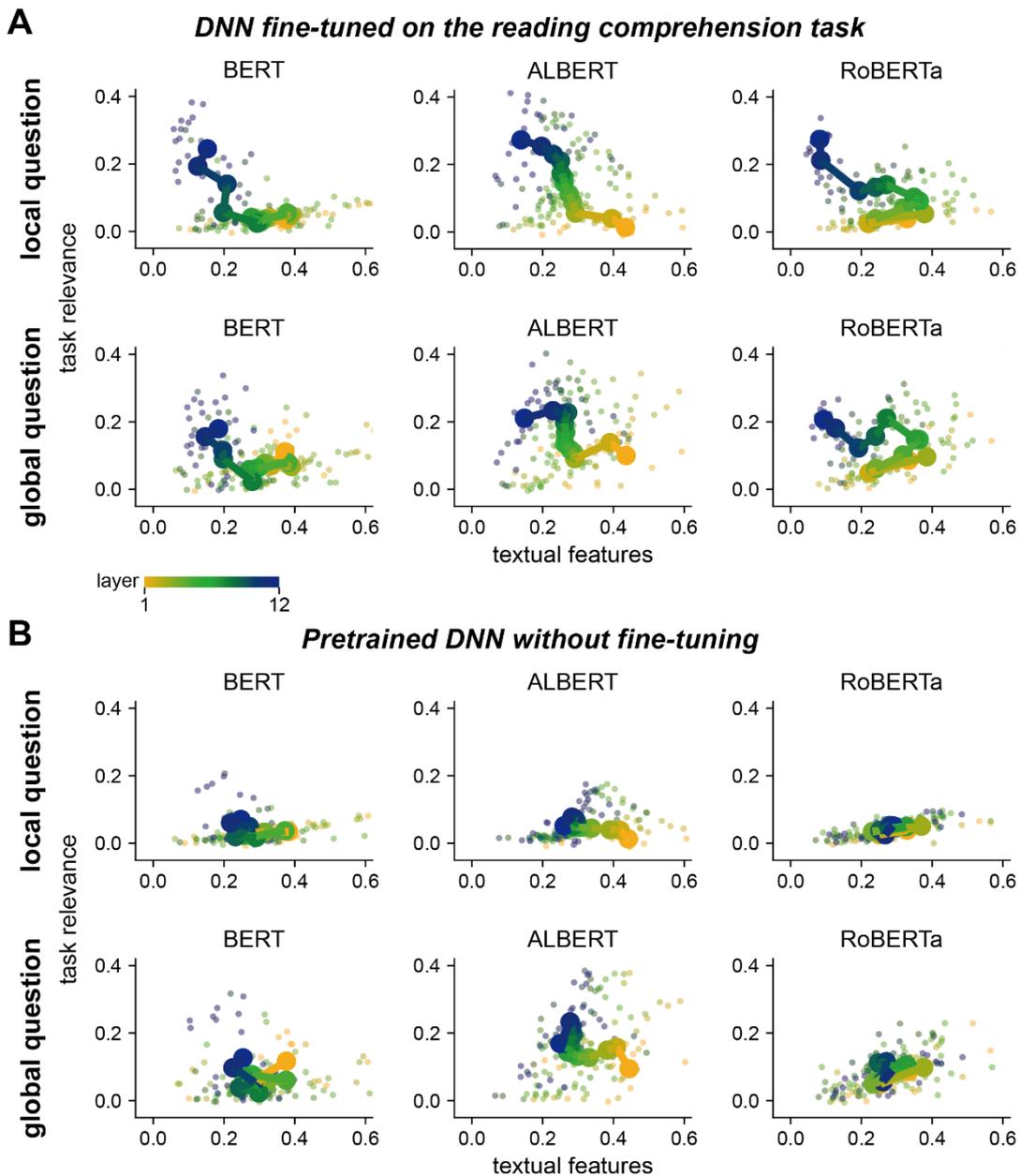
269 For local questions, textual and layout features, but not task relevance, could predict
270 human attention distribution during general-purpose reading ($P = 0.003, 0.003, \text{ and } 1$
271 for textual features, layout features, and task relevance, permutation test, FDR
272 corrected). For global questions, all three features could predict human attention
273 distribution ($P = 0.003, 0.003, \text{ and } 0.003$ for all 3 features, permutation test, FDR
274 corrected). DNNs could also predict human attention distribution during general-
275 purpose reading, but most of the effect was explained by textual and layout features
276 (Fig. 3B, the middle plot).

277

278 Comparing the results obtained from Study 1 and Study 2, it was evident that human
279 attention could be modulated by the specific reading goal, i.e., the question to answer,
280 on top of textual and layout features. Goal-directed top-down attention, characterized
281 in Study 1, could be modeled by either human-annotated task relevance or the DNN
282 attention weights. In the absence of a specific reading goal, human attention in Study 2
283 was mainly influenced by stimulus features, e.g., textual and layout features, which
284 were also captured by the DNN attention weights.

285

Properties of DNN Attention in Different Layers



286

287 **Fig. 4. Influence of stimulus features and top-down task on each DNN layer.** The
288 same regression analyses in Fig. 3 are employed to analyze how the DNN attention is
289 affected by lower-level stimulus features and top-down task relevance. Panels A and B
290 show the results for the DNNs fine-tuned based on the reading comprehension task and
291 the pre-trained DNNs that receive no fine-tuning. Each small dot shows the result from
292 an attention head, and each large dot shows the average over heads of the same layer.
293 Color indicates layer number. Shallow layers of both fine-tuned and pre-trained DNN
294 are more sensitive to stimulus features. Deep layers of fine-tuned DNN, but not pre-
295 trained DNN, are sensitive to task relevance.

296 **DNN Attention in Different Layers**

297 Previous studies have shown that different layers in DNN encoded different types of
298 information⁴²⁻⁴⁴. In the following, we analyzed whether the properties of DNN attention
299 weights differed across layers. Since human attention was influenced by both bottom-
300 up stimulus features and top-down task goal, in the following we also analyzed how
301 these features influenced the attention weights in each DNN layer. Since the layout
302 features were not available to the DNNs, we only considered textual features as
303 stimulus features in this analysis. As shown in Fig. 4A, the attention weights in different
304 layers were sensitive to different features. In general, shallow layers were more strongly
305 influenced by textual features while deeper layers were more strongly influenced by the
306 task relevance. This trend was observed in all 3 DNN models and was especially
307 obvious for local questions. The transitional trajectory across layers, however, was
308 model-dependent in the 2-dimensional feature space. In Fig. 2, examples were shown
309 for the attention weights averaged over all 12 heads in the last layer of BERT, which
310 resembled the human-annotated task relevance.

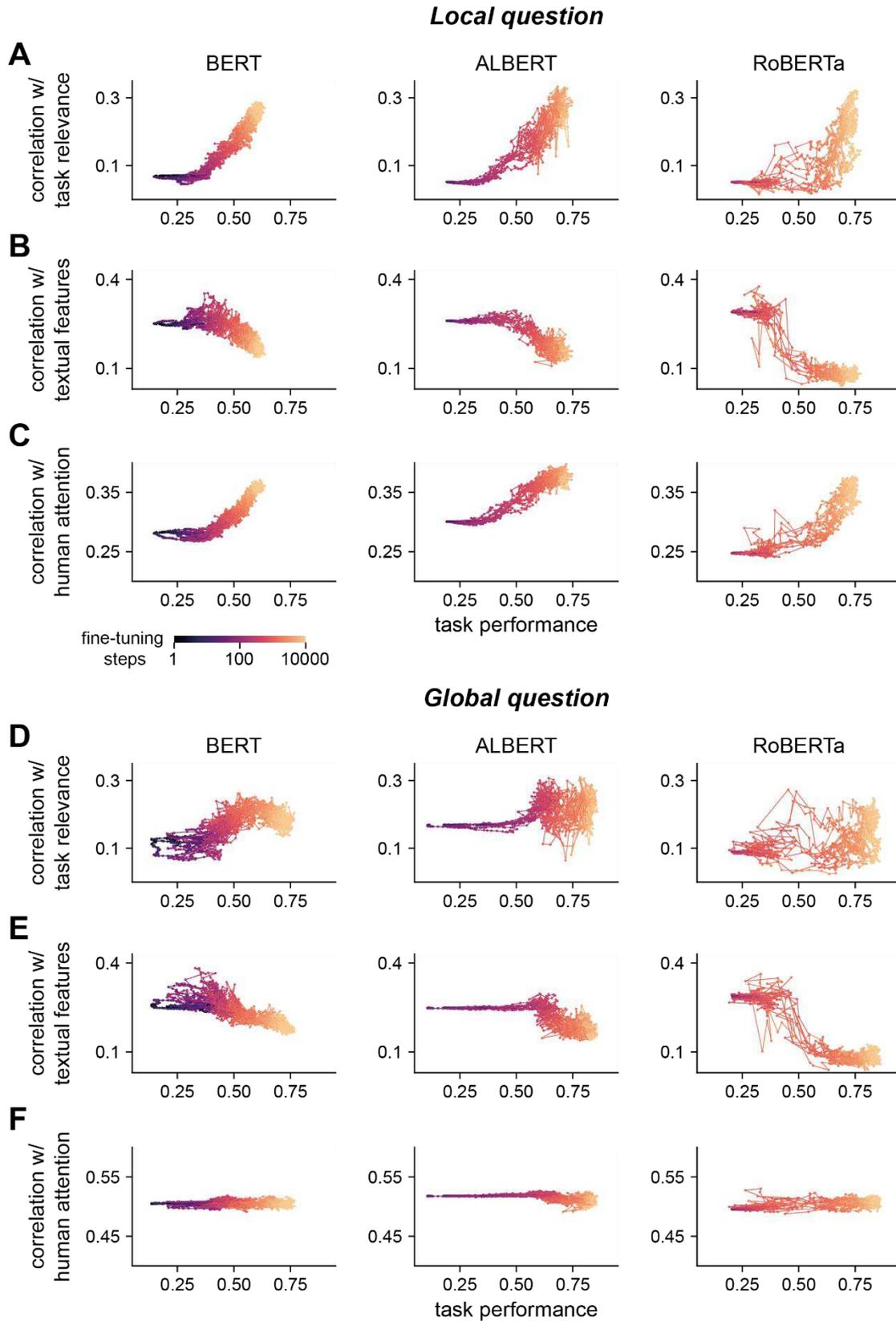
311

312 **Evolution of DNN Attention during Fine-Tuning**

313 All the 3 DNN models were pre-trained based on large-scale corpora and fine-tuned
314 based on the reading comprehension task (see *Materials and Methods*). Was the DNN
315 attention mechanism mainly shaped by the pre-training process or the fine-tuning
316 process? We addressed this question by analyzing the attention weights in pre-trained
317 DNN models that did not receive fine-tuning (Fig. 4B). It was found that the attention
318 weights of pre-trained DNN were sensitive to textual features in shallow layers but not
319 sensitive to task relevance in deeper layers, suggesting that top-down attention in DNNs

320 emerged during fine-tuning using the reading comprehension task.
321
322 We then asked how the attention weights of DNN changed during fine-tuning and
323 whether such changes were related to the performance of question answering. During
324 fine-tuning, the structure of the DNN model remained but the parameters were adjusted.
325 In the following, we analyzed the properties of models that received different steps of
326 fine-tuning. Furthermore, since fine-tuning process was stochastic, we fine-tuned 10
327 times (see *Materials and Methods*). We found that, in deep layers, the properties of
328 attention weights significantly changed during fine-tuning (Fig. S3 and Fig. S4). In the
329 last layer, for example, it was clear that the DNN attention weights became more
330 sensitive to task relevance during fine-tuning, coinciding with the improvement in task
331 performance (Fig. 5AD), especially for local questions (Fig. 5A). The trend is less clear
332 for global questions and a potential explanation is that global questions concern the
333 main topic of the passage and can be answered by paying attention to different sets of
334 words. Deep layer's sensitivity to textual features, however, dropped during fine-tuning
335 (Fig. 5BE). Therefore, fine-tuning directed deep layers' attention towards task relevant
336 information, sacrificing the sensitivity to textual features. Additionally, we found that
337 the similarity between DNN attention weights and human attention was also boosted
338 by fine-tuning for local questions (Fig. 5C). This result further demonstrated that
339 human-like attention in DNNs was the consequence of optimization of the reading
340 comprehension task, instead of the consequence of more general pre-training language
341 tasks.

Influence of Fine-Tuning on DNN Attention and Task Performance



342

343

344

345

Fig. 5. Influence of fine-tuning on DNN attention and task performance. Each model is fine-tuned 10 times. Each data point denotes the result during a fine-tuning step (color coded), and steps from each run of fine-tuning was connected by a line.

346 (A, B, D, and E) The effect of fine-tuning on the attention weights in the last layer of
347 DNN for local questions (A and B) and global questions (D and E). Fine-tuning
348 enhances the sensitivity to top-down task relevance while reducing the sensitivity to
349 lower-level textual features, which correlates with the increase in task performance.
350 (C and F) Influence of fine-tuning on the similarity between DNN and human
351 attention. For local questions (C), fine-tuning clearly increases the similarity between
352 DNN and human attention, coinciding with the increase in task performance. For
353 global questions (F), the similarity between DNN and human attention is high even
354 without fine-tuning and is not further boosted by fine-tuning. For ALBERT, 2 out of
355 the 10 runs of fine-tuning are unstable, showing sharp drops in task performance
356 during fine-tuning. Results of these 2 runs are not shown here but separately shown in
357 Fig. S2.

358

359

360 **Discussion**

361 Since attention is a key mechanism for both the biological brain and artificial neural
362 networks, it provides a common ground to quantitatively compare biological neural
363 computations and artificial neural computations. Such comparisons, however, are
364 challenging since biological and artificial neural networks are investigated in different
365 fields using very different approaches. The current study attempts to bridge this gap by
366 building a large eye tracking dataset for a real-world reading task that is of interest to
367 both the psychology and AI community. Based on these data, it is shown that, when
368 optimized to perform a reading comprehension task, DNNs naturally evolve human-
369 like attention distribution. On the one hand, the results indicate that the attention
370 mechanism in DNNs could indeed be of biological relevance. On the other hand, it
371 provides a plausible computational explanation for human attention distribution.

372

373 **Computational models of biological attention**

374 Deep neural network models of biological attention are best studied in vision. A large
375 number of models are proposed to predict bottom-up visual saliency^{31,32}, and recently
376 DNN models are also employed to model top-down visual attention. It is shown that,
377 through either implicit^{34,35} or explicit training^{9,36}, DNN can predict which parts of a
378 picture relates to a verbal phrase, a task similar to goal-directed visual search³⁰. The
379 current study distinguishes from these studies in that the DNN model is not trained to
380 predict human attention. Instead, the DNN models naturally generate human-like
381 distribution when trained to perform the same task that humans perform. Therefore, the
382 current study suggests that DNN models can potentially serve as a mechanistic, instead
383 of a descriptive model of human attention during reading comprehension. Previous
384 studies have also proposed mechanistic models of biological attention. It has been
385 proposed that attention can be interpreted a mechanism to implement optimal decision
386 making. For example, when faced with multiple conflicting cues, the brain can use
387 attention to modulate, i.e., weight, the neural representation of each cue. It has been
388 proposed that the brain attends to more informative and reliable cues to make an optimal
389 decision⁴⁵⁻⁴⁷. The current results are generally consistent with this idea since both
390 human and DNN attend to words that are relevant to task solving.

391

392 **Attention during human reading**

393 How human readers allocate attention during reading is an extensively studied topic.

394 Eye tracking studies have shown that the readers fixate longer at, e.g., longer words,

395 words of lower-frequency, words that are less predictable based on the context, and
396 words at the beginning of a line^{48,49}. A number of models, e.g., the E-Z reader^{37,50} and
397 SWIFT⁵¹, have been proposed to predict the eye movements during reading, either
398 based on basic oculomotor properties or lexical processing³⁷. These models can
399 generate fine-grained predictions, e.g., which letter in a word will be fixated first. A
400 limitation of these models, however, is that they are generally developed to explain the
401 reading of simple sentences, instead of complex sentences or multi-line text. In contrast
402 to these studies, the current study focuses on macroscopic distribution of attention, i.e.,
403 the distribution of total fixation time in the units of words, when readers read
404 challenging multi-line text. Future studies can potentially integrate classic eye
405 movement models with DNNs to explain the dynamic eye movement trajectory,
406 possibly with a letter-based spatial resolution.

407

408 A more important difference between the current and previous studies on reading
409 attention is that the current study investigates how attention is affected by a specific
410 top-down reading goal. In previous studies, readers are generally instructed to read a
411 sentence in a normal manner, not aimed to extract a specific kind of information. In the
412 current study, however, readers know in advance what question they have to answer
413 and this kind of reading can be viewed as a kind of information seeking behavior⁵², and
414 is also referred to as the reading-to-do task³⁸. Previous studies have shown the reader's
415 task may have heterogeneous influences on attention, depending on the task difficulty
416 and skill level of readers^{53,54}. Here, the task is demanding and the readers are highly

417 skilled to perform the task: The reading comprehension questions are selected from
418 exams and the time to answer each question is limited, leading to about 80% question
419 answering accuracy (Fig. 1). The participants are skilled since all Chinese students have
420 extensive practice in such reading comprehension questions in high school. Future work
421 is needed to quantify how the task and reading skills modulate human attention and
422 whether these effects can also be modeled by DNN models.

423

424 **Attention mechanisms in DNN**

425 In DNN NLP models, attention is a mechanism to selectively integrate information
426 across words (or other units of representations), and is typically implemented by
427 assigning different weights to different words. How the attention mechanism is
428 integrated with the neural network model and the role it plays, however, differ across
429 models^{36,55}. In some models, attention is an explicit information integration mechanism.
430 For example, in many models^{22,39,56-58}, the representations of all words in a passage are
431 integrated with different weights to compute a representation of the passage. In these
432 models, ideally, words that contribute more to task solving should receive higher
433 weights.

434

435 In transformer-based models, the roles self-attention plays are highly diverse. Since
436 self-attention assigns a weight between every pair of inputs (including words and
437 special tokens such as CLS), it can capture a number of relationship between words,
438 e.g., co-reference and syntactic dependency²⁴⁻²⁶. In the current study, however, we only

439 analyze a small portion of the self-attention weights that are directly relevant to the task,
440 i.e., the attention weights between CLS and words in the passage. The attention weights
441 analyzed here can be interpreted as a selective information integration mechanism,
442 describing how different words in a passage contribute to the decision variable, i.e.,
443 CLS. The current study demonstrated that 3 transformer-based models generate human-
444 like attention through task optimization. It remains unclear whether other DNN models
445 show similar properties. Nevertheless, the dataset and methods developed here can be
446 easily applied to test whether other models also evolve human-like attention
447 distribution, serving as a probe to test the biological plausibility of NLP models.

448

449 **Interpretation of the attention mechanisms in DNN**

450 Whether attention can increase the interpretability of DNN models is a topic that
451 receives a considerable amount of debates. A number of studies have shown that the
452 DNN attention weights are higher for words that are more important for the task^{22,23,58,59}.
453 Most of these studies, however, are based on visual inspection of a couple of examples.
454 Other studies, however, find low correlation between attention weights and other
455 measures of the importance of words, and therefore raise concern about whether the
456 attention weights are interpretable^{25,26,60,61}. The importance of a word, however, can be
457 measured in many different ways, and no correlation with some importance measures
458 does not indicate no contribution to task solving in other ways. Here, we show that, for
459 3 transformer-based models, the DNN attention weights correlate with human attention
460 and are modulated by both textual features and the task. Furthermore, by analyzing the

461 fine tuning process (Fig. 5C), we found that DNN with human-like attention perform
462 better at the task.

463

464 Here, it is revealed that different layers in DNN show different attention properties,
465 with the deep layers being more sensitive to task relevance. This result is consistent
466 with previous findings that artificial neurons in deep layers of DNN encode more
467 abstract information, e.g., object information in convolutional networks⁴² and syntactic
468 information in BERT^{43,44}. A recent study has also compared human eye movements and
469 attention weights in the last layer of BERT, when participants evaluate whether a
470 passage is an appropriate answer to a question. It is demonstrated that the human
471 fixation time is more similar to the attention weights in BERT than the simple TF-IDF
472 weights²⁷.

473

474 In sum, the current study demonstrates that, when DNN and humans perform the same
475 reading comprehension task with comparable accuracy, the DNN attention weights
476 resemble human attention measured by eye tracking. The results suggest that human
477 attention distribution is shaped by the demand to optimally perform the task and the
478 DNN attention can be interpreted as an approximation of human attention. The large
479 set of eye tracking data in the current study can also motivate future computational
480 modeling of human attention during natural reading tasks and be applied to test whether
481 other NLP models exhibit human-like attention distribution.

482

483

484 **Materials and Methods**

485 **Participants**

486 Study 1 enrolled 102 participants (19-30 years old, mean age, 22.9 years; 54 female).

487 Study 2 enrolled a separate group of 18 participants (21-26 years old, mean age, 23.4

488 years; 10 female). All participants were native Chinese speakers and were college

489 students or graduate students at Zhejiang University, and were thus above the level

490 required to answer high-school-level reading comprehension questions. English

491 proficiency levels were further guaranteed by the following criterion for screening

492 participants: a minimum score of 6 on IELTS, 80 on TOEFL, or 425 on CET6¹. The

493 experimental procedures were approved by the Research Ethics Committee of the

494 College of Medicine, Zhejiang University (2019–047). The participants provided

495 written consent and were paid.

496

497 **Experimental materials**

498 The reading materials were selected and adapted from the large-scale RACE dataset, a

499 collection of reading comprehension questions in English exams for middle and high

500 schools in China³⁹. We selected eight hundreds of high-school level questions from the

501 test set of RACE and each question was associated with a distinct passage (117 to 456

502 words per passage). All questions were multiple-choice questions with 4 alternatives

503 including only one correct option among them. The questions fell into 6 types, i.e.,

¹ The National College English Test (CET) is a national English test system developed to examine the English proficiency of undergraduate students in China. CET includes tests of two levels: a lower level test CET4 and a higher level test CET6.

504 Cause ($N = 200$), Fact ($N = 200$), Inference ($N = 120$), Theme ($N = 100$), Title ($N =$
505 100), and Purpose ($N = 80$). The Cause, Fact, and Inference questions were concerned
506 with the location, extraction, and comprehension of specific information from a passage,
507 and were referred to as local questions. Questions of Theme, Title, and Purpose tested
508 the understanding of a passage as a whole, and were referred to as global questions. We
509 further acquired annotations about the relevance of each word to the question answering
510 task. Details about the question types and the annotation procedures could be found in
511 reference⁶².

512

513 **Experimental procedures**

514 **Study 1:** Study 1 included all 800 passages, and different question types were
515 separately tested in different experiments, hence six experiments in total. Each
516 experiment included 25 participants and one participant could participate in multiple
517 experiments. Before each experiment, participants were given a familiarization session
518 with 5 questions that were not used in the formal experiment. During the formal
519 experiment, questions were presented in a randomized order. Considering the quantities
520 of questions, for Cause and Fact questions, the experiment was carried out in 3 separate
521 days (one third questions on each day), and for other question types the experiment was
522 carried out in 2 days (fifty percent of questions on each day).

523

524 The experiment procedure in Study 1 was illustrated in Fig. 1A. In each trial,
525 participants first read a question, pressed the space bar to read the corresponding

526 passage, and then pressed it again to read the question coupled with 4 options and
527 answer the question. The time limit for passage reading was 120 s. To encourage the
528 participants to read as quickly as possible, the bonus they received for a specific
529 question would decrease linearly over time. They did not receive any bonus for the
530 question, however, if they gave a wrong answer. Furthermore, before answering the
531 comprehension question, the participants reported whether they were confident that
532 they could correctly answer the question. After answering the question, they also rated
533 their confidence about their answer on the scale of 1-4 (low to high). The confidence
534 ratings were not analyzed.

535

536 **Study 2:** Study 2 included 96 reading passages and questions, with 16 questions for
537 each question type that were randomly selected from the questions used in Study 1. The
538 study was carried out in 2 days, and none of the participants participated in Study 1.
539 The familiarization procedure was identical to that in Study 1.

540

541 The procedure of Study 2 was similar to that of Study 1, and the main difference was
542 that a 90-s first-pass passage reading stage was introduced at the beginning of each trial.
543 During the first-pass passage reading, participants had no prior information of the
544 relevant question. The participants could press the space bar to terminate the first-pass
545 reading stage and to read a question. Then, participants read the passage for the second
546 time with a time limit of 30 s, before proceeding to answer the question. In Study 2, the
547 correctness of the answer was also the prerequisite for bonus, and the amount of bonus
548 decreased linearly with the duration of second-pass passage reading.

549

550 **Stimulus presentation and eye tracking**

551 The text was presented using the bold Courier New font, and each letter occupied 14 ×
552 27 pixels. We set the maximum number of letters on each line to 120 and used double
553 space. We separated paragraphs by indenting the first line of each new paragraph.
554 Participants sat about 880 mm from a monitor, at which each letter horizontally
555 subtended approximately 0.25 degrees of visual angle.

556

557 Eye tracking data were recorded from the left eye with 500-Hz sampling rate (Eyelink
558 Portable Duo, SR Research). The experiment stimuli were presented on a 24-inch
559 monitor (1920x1080 resolution; 60 Hz refresh rate) and administered using MATLAB
560 Psychtoolbox⁶³. Each experiment started with a 13-point calibration and validation of
561 eye tracker, and the validation error was required to be below 0.5° of visual angle.
562 Furthermore, before each trial, a 1-point validation was applied, and if the calibration
563 error was higher than 0.5°, a recalibration was carried out. Head movements were
564 minimized using a chin and forehead rest.

565

566 **DNN models**

567 We tested 3 popular transformer-based DNN models, i.e., BERT¹⁷, ALBERT¹⁸, and
568 RoBERTa¹⁹. ALBERT and RoBERTa were both adapted from BERT, and had the same
569 basic structure. RoBERTa differed from BERT in its pre-training procedure¹⁹ while
570 ALBERT applied factorized embedding parameterization and cross-layer parameter

571 sharing to reduce memory consumption¹⁸. Following previous works^{18,19}, each option
 572 was independently processed. For the i^{th} option ($i = 1, 2, 3, \text{ or } 4$), the question and the
 573 option were concatenated to form an integrated option. As shown in Fig. 1C, for the i^{th}
 574 option, the input to DNN was the following sequence:

$$575$$

$$576 \quad C_i, P_1, P_2, \dots, P_N, S_{i,1}, O_{i,1}, O_{i,2}, \dots, O_{i,M}, S_{i,2},$$

577 where C_i , $S_{i,1}$, and $S_{i,2}$ denoted special tokens, i.e., the CLS, SEP₁, and SEP₂ tokens,
 578 separating different components of the input. P_1, P_2, \dots, P_N denoted all the N words of
 579 a passage, while $O_{i,1}, O_{i,2}, \dots, O_{i,M}$ denoted all the M words of the i^{th} integrated option.
 580 Each of the token was represented by a vector. The vectorial representation was updated
 581 in each layer, and in the following the output of the l^{th} layer was denoted as a superscript,
 582 e.g., C_i^l . Following previous works^{18,19}, we calculated a score for each option, which
 583 indicated the possibility that the option was the correct answer. The score was
 584 calculated by first applying a linear transform to the final representation of the CLS
 585 token, i.e.,

$$586$$

$$587 \quad s_i = \Phi C_i^{12},$$

588 where C_i^{12} was the final output representation of CLS and Φ was a vector learned from
 589 data. The score was independently calculated for each option and then normalized using
 590 the following equation:

$$591$$

$$592 \quad score_i = \frac{\exp(s_i)}{\sum_{i=1}^4 \exp(s_i)}.$$

593 The answer to a question was determined as the option with highest score, and all the
594 models were trained to maximize the logarithmic score of the correct option.

595

596 We fine-tuned DNN based on the training set of RACE. In the analyses shown in Figs.
597 5 and 6, the fine-tuning process was independently run 10 times. Each time, the training
598 samples were fed in with a randomized order and nodes in the dropout layer were
599 randomly eliminated. Results from the first run of fine-tuning was used for the main
600 analysis reported in Figs. 2-4. All models were implemented based on HuggingFace⁶⁴
601 and all hyperparameters for fine-tuning were adopted from previous studies (Table S1).
602 To isolate how the fine-tuning process modulated DNN attention, we also tested the
603 pre-trained DNN that was not fine-tuned on RACE dataset, and compared it with the
604 fine-tuned model (Fig. 4). Furthermore, we quantified how the properties of DNN
605 attention changed throughout the fine-tuning process by analyzing models that received
606 different steps of fine tuning. The steps we sampled were exponentially spaced between
607 1 and the maximum fine-tuning steps.

608

609 **DNN attention**

610 In each attention head, the attention mechanism calculated an attention weight between
611 any pair of inputs, including words and special tokens. The vectorial representation of
612 each input was then updated by the weighted sum of the vectorial representations of all
613 inputs²⁰. In other words, the models we considered were all context-dependent models,
614 in which the representation of each word was modeled by integrating the

615 representations of all inputs. Since only the CLS token was directly related to question
616 answering, here we analyzed the attention weights that were used to calculate the
617 vectorial representation of CLS (illustrated in Fig. 1D). For each layer, the output of an
618 attention head was computed using the following equations. For the sake of clarity, we
619 denote the input words and tokens generally as X_i .

620

$$621 \quad C^h = \sum_{i=1}^{N+M+2} \alpha_i V_i = \alpha_C V_C + \sum_{n=1}^N \alpha_{Pn} V_{Pn} + \alpha_{S1} V_{S1} + \sum_{m=1}^M \alpha_{Om} V_{Om} + \alpha_{S2} V_{S2},$$

$$622 \quad \alpha_i = \frac{\exp(Q_C K_i^T)}{\sum_{i=1}^{N+M+2} \exp(Q_C K_i^T)},$$

$$623 \quad V_i = X_i W^V + b^V, K_i = X_i W^K + b^K, Q_C = X_C W^Q + b^Q,$$

624 where W^V , W^Q , W^K , b^V , b^Q , and b^K were parameters to learn from the data. The attention
625 weight between CLS and the n^{th} word in the passage, i.e., α_{Pn} , was compared to human
626 attention. Here, we only considered the attention weight associated with the correct
627 option.

628

629 Output of the attention module, i.e., C^h , was concatenated over all the 12 heads in each
630 layer, and further processed by position-wise operations to generate the final
631 representation of CLS in the layer²⁰. Additionally, DNN used byte-pair tokenization
632 which split some words into multiple tokens. We converted the token-level attention
633 weights to word-level attention weights by summing the attention weights over tokens
634 within a word^{24,27}.

635

636 **Human attention analysis and prediction**

637 We analyzed eye fixations during passage reading in Study 1 and the first-pass passage
638 reading in Study 2. For each word, the total fixation time was the sum of the duration
639 across all fixations that fell into the square area the word occupied. We averaged the
640 total fixation time across all participants who correctly answered the question, and
641 measured human attention using the attention density, i.e., the total fixation time
642 divided by the area a word occupied.

643

644 We employed linear regression to test whether a set of features could explain human
645 attention distribution. Four sets of features were analyzed, i.e., textual features, layout
646 features, task relevance, and DNN attention weights. The textual features included word
647 length, logarithmic word frequency estimated based on the British National Corpus ⁶⁵,
648 ordinal position of a word in a sentence, ordinal position of a word in a passage, and
649 ordinal sentence number of a word. The layout features referred to the visual layout of
650 text, i.e., features induced by line changes, including the coordinate of the left most
651 pixel of a word, ordinal position of a word in a paragraph, ordinal row number of a
652 word in a paragraph, ordinal row number of a word in a passage. Task relevance was
653 annotated by humans, and the DNN attention weights included the 144 attention
654 weights from all layers and attention heads. In the regression analysis, human attention
655 density on word w was modeled using the following equation.

656

657
$$attention_density_w = \sum_{j=1}^J \beta_j F_{w,j} + b + \varepsilon_w,$$

658 where F and ε denoted the features being considered and the residual error, respectively.
659 The parameters β and b were fitted to minimize the mean square error. Each feature and
660 the human attention distribution were normalized within a passage by taking the z-score.
661 The prediction accuracy, i.e., the correlation between predicted attention and actual
662 human attention, was calculated based on five-fold cross-validation. Each question type
663 was separately modeled.

664

665 **Statistical tests**

666 We employed a one-sided permutation test to test whether the attention distribution
667 predicted by a set of features significantly correlated with human attention. Five
668 hundreds of chance-level prediction accuracy was calculated by predicting shuffled
669 human attention. Specifically, the human attention density was shuffled across words
670 and was predicted by word features which were not shuffled. The procedure was
671 repeated 500 times, creating 500 chance-level prediction accuracy. If the actual
672 correlation was greater than N out of the 500 chance-level correlation, the significance
673 level was $(N + 1)/501$.

674

675 The comparison between global and local questions were based on bias-corrected and
676 accelerated bootstrap⁶⁶. For example, to test whether the prediction accuracy differed
677 between the 2 types of questions, all global questions were resampled with replacement
678 5000 times and each time the prediction accuracy was calculated based on the
679 resampled questions, resulting in 5000 resampled prediction accuracy. If the prediction

680 accuracy for local questions was greater (or smaller) than N out of the 5000 resampled
681 accuracy for global questions, the significance level of their difference was $2(N +$
682 $1)/5001$. When multiple comparisons were performed, the p-value was further adjusted
683 using the false discovery rate (FDR) correction.

684

685 **Acknowledgements**

686 We thanks David Poeppel, Jonathan Simon, and Xunyi Pan for thoughtful comments
687 on earlier versions of this manuscript; Yuran Zhang, Anqi Dai, Zhonghua Tang, and
688 Yuhan Lu for assistance with experiments; Peiqing Jin and Cheng Luo for helpful
689 discussions. Work supported by National Natural Science Foundation of China
690 31771248 and Major Scientific Research Project of Zhejiang Lab 2019KB0AC02

691

692 **Author contributions**

693 Nai Ding acquired the funding, conceived and coordinated the project, analyzed data,
694 and wrote the manuscript. Jiajie Zou implemented the experiments and models,
695 analyzed data, and wrote the manuscript.

696

697 **Competing interests**

698 The authors declare no competing interests.

699

700

701 **References**

- 702 1 Nasr, K., Viswanathan, P. & Nieder, A. Number detectors spontaneously
703 emerge in a deep neural network designed for visual object recognition.
704 *Science Advances* **5**, eaav7903 (2019).
- 705 2 Kim, G., Jang, J., Baek, S., Song, M. & Paik, S.-B. Visual number sense in
706 untrained deep neural networks. *Science Advances* **7**, eabd6127 (2021).
- 707 3 Sheahan, H., Luyckx, F., Nelli, S., Teupe, C. & Summerfield, C. Neural state
708 space alignment for magnitude generalization in humans and recurrent
709 networks. *Neuron* **109**, 1214-1226. e1218 (2021).
- 710 4 Xu, Y. & Vaziri-Pashkam, M. Limits to visual representational correspondence
711 between convolutional neural networks and the human brain. *Nature*
712 *communications* **12**, 1-16 (2021).
- 713 5 Ettinger, A. What BERT is not: Lessons from a new suite of psycholinguistic
714 diagnostics for language models. In *Proc. Transactions of the Association for*
715 *Computational Linguistics*, 34-48 (Association for Computational Linguistics,
716 2020).
- 717 6 Hale, J., Dyer, C., Kuncoro, A. & Brennan, J. Finding syntax in human
718 encephalography with beam search. In *Proc. 56th Annual Meeting of the*
719 *Association for Computational Linguistics*, 2727-2736 (Association for
720 Computational Linguistics, 2018).
- 721 7 Linzen, T., Dupoux, E. & Goldberg, Y. Assessing the ability of LSTMs to
722 learn syntax-sensitive dependencies. In *Proc. Transactions of the Association*
723 *for Computational Linguistics*, 521-535 (Association for Computational
724 Linguistics, 2016).
- 725 8 Judd, T., Durand, F. & Torralba, A. A benchmark of computational models of
726 saliency to predict human fixations. In *Proc. 11th Joint Conference on*
727 *Computer Vision, Imaging and Computer Graphics Theory and Applications*,
728 134-142 (SciTePress, 2016).
- 729 9 Liu, C., Mao, J., Sha, F. & Yuille, A. Attention correctness in neural image

- 730 captioning. In *Proc. AAAI Conference on Artificial Intelligence*, 4176-4182
731 (AAAI, 2017).
- 732 10 Yamins, D. L. *et al.* Performance-optimized hierarchical models predict neural
733 responses in higher visual cortex. *Proc. Natl Acad. Sci. USA* **111**, 8619-8624
734 (2014).
- 735 11 Kell, A. J., Yamins, D. L., Shook, E. N., Norman-Haignere, S. V. &
736 McDermott, J. H. A task-optimized neural network replicates human auditory
737 behavior, predicts brain responses, and reveals a cortical processing hierarchy.
738 *Neuron* **98**, 630-644. e616 (2018).
- 739 12 Donhauser, P. W. & Baillet, S. Two distinct neural timescales for predictive
740 speech processing. *Neuron* **105**, 385-393. e389 (2020).
- 741 13 Seidenberg, M. S. & McClelland, J. L. A distributed, developmental model of
742 word recognition and naming. *Psychological review* **96**, 523-568 (1989).
- 743 14 Pinker, S. & Prince, A. On language and connectionism: Analysis of a parallel
744 distributed processing model of language acquisition. *Cognition* **28**, 73-193
745 (1988).
- 746 15 McClelland, J. L. & Elman, J. L. The TRACE model of speech perception.
747 *Cognitive psychology* **18**, 1-86 (1986).
- 748 16 Zhou, B., Khosla, A., Lapedriza, A., Oliva, A. & Torralba, A. Object detectors
749 emerge in deep scene CNNs. In *Proc. International Conference on Learning*
750 *Representations*, (ICLR, 2015).
- 751 17 Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep
752 bidirectional transformers for language understanding. In *Proc. 2019*
753 *Conference of the North American Chapter of the Association for*
754 *Computational Linguistics: Human Language Technologies*, 4171-4186
755 (Association for Computational Linguistics, 2019).
- 756 18 Lan, Z. *et al.* Albert: A lite bert for self-supervised learning of language
757 representations. In *Proc. International Conference on Learning*
758 *Representations*, (ICLR, 2020).
- 759 19 Liu, Y. *et al.* Roberta: A robustly optimized bert pretraining approach. Preprint

- 760 at <https://arxiv.org/abs/1907.11692> (2019).
- 761 20 Vaswani, A. *et al.* Attention is all you need. In *Proc. Advances in neural*
762 *information processing systems*, 5998-6008 (Curran Associates, 2017).
- 763 21 Brown, T. B. *et al.* Language models are few-shot learners. In *Proc. Advances*
764 *in Neural Information Processing Systems*, 1877-1901 (Curran Associates,
765 2020).
- 766 22 Yang, Z. *et al.* Hierarchical attention networks for document classification. In
767 *Proc. 2016 Conference of the North American Chapter of the Association for*
768 *Computational Linguistics: Human Language Technologies*, 1480-1489
769 (Association for Computational Linguistics, 2016).
- 770 23 Lin, Z. *et al.* A structured self-attentive sentence embedding. In *Proc.*
771 *International Conference on Learning Representations*, (ICLR, 2017).
- 772 24 Clark, K., Khandelwal, U., Levy, O. & Manning, C. D. What Does BERT
773 Look at? An Analysis of BERT's Attention. In *Proc. 2019 ACL Workshop*
774 *BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, 276-286
775 (Association for Computational Linguistics, 2019).
- 776 25 Kovaleva, O., Romanov, A., Rogers, A. & Rumshisky, A. Revealing the Dark
777 Secrets of BERT. In *Proc. 2019 Conference on Empirical Methods in Natural*
778 *Language Processing and the 9th International Joint Conference on Natural*
779 *Language Processing (EMNLP-IJCNLP)*, 4365-4374 (Association for
780 Computational Linguistics, 2019).
- 781 26 Voita, E., Talbot, D., Moiseev, F., Sennrich, R. & Titov, I. Analyzing multi-
782 head self-attention: Specialized heads do the heavy lifting, the rest can be
783 pruned. In *Proc. 57th Annual Meeting of the Association for Computational*
784 *Linguistics*, 5797-5808 (Association for Computational Linguistics, 2019).
- 785 27 Bolotova, V. *et al.* Do People and Neural Nets Pay Attention to the Same
786 Words: Studying Eye-tracking Data for Non-factoid QA Evaluation. In *Proc.*
787 *29th ACM International Conference on Information & Knowledge*
788 *Management*, 85-94 (ACM, 2020).
- 789 28 Corbetta, M. & Shulman, G. L. Control of goal-directed and stimulus-driven

790 attention in the brain. *Nat. Rev. Neurosci.* **3**, 201-215 (2002).

791 29 Petersen, S. E. & Posner, M. I. The attention system of the human brain: 20
792 years after. *Annu. Rev. Neurosci* **35**, 73-89 (2012).

793 30 Wolfe, J. M. & Horowitz, T. S. Five factors that guide attention in visual
794 search. *Nat. Hum. Behav.* **1**, 1-8 (2017).

795 31 Borji, A., Sihite, D. N. & Itti, L. Quantitative analysis of human-model
796 agreement in visual saliency modeling: a comparative study. *IEEE Trans.*
797 *Image Process.* **22**, 55-69 (2013).

798 32 Bylinskii, Z., Judd, T., Oliva, A., Torralba, A. & Durand, F. What do different
799 evaluation metrics tell us about saliency models? *IEEE Trans. Pattern Anal.*
800 *Mach. Intell.* **41**, 740-757 (2018).

801 33 Tatler, B. W., Hayhoe, M. M., Land, M. F. & Ballard, D. H. Eye guidance in
802 natural vision: reinterpreting salience. *J. Vis.* **11**, 5-25, doi:10.1167/11.5.5
803 (2011).

804 34 Anderson, P. *et al.* Bottom-up and top-down attention for image captioning
805 and visual question answering. In *Proc. The IEEE Conference on Computer*
806 *Vision and Pattern Recognition (CVPR)*, 6077-6086 (IEEE, 2018).

807 35 Xu, K. *et al.* Show, attend and tell: Neural image caption generation with
808 visual attention. In *Proc. 32nd International Conference on Machine*
809 *Learning*, 2048-2057 (PMLR, 2015).

810 36 Das, A., Agrawal, H., Zitnick, L., Parikh, D. & Batra, D. Human attention in
811 visual question answering: Do humans and deep networks look at the same
812 regions? *Comput. Vis. Image Underst.* **163**, 90-100 (2017).

813 37 Reichle, E. D., Rayner, K. & Pollatsek, A. The EZ Reader model of eye-
814 movement control in reading: Comparisons to other models. *Behav. Brain Sci.*
815 **26**, 445-476 (2003).

816 38 Duffy, T. M. & Kabance, P. Testing a readable writing approach to text
817 revision. *Journal of educational psychology* **74**, 733-748 (1982).

818 39 Lai, G., Xie, Q., Liu, H., Yang, Y. & Hovy, E. Race: Large-scale reading
819 comprehension dataset from examinations. In *Proc. 2017 Conference on*

- 820 *Empirical Methods in Natural Language Processing*, 785-794 (Association for
821 Computational Linguistics, 2017).
- 822 40 Liu, S. S., Zhang, X., Zhang, S., Wang, H. & Zhang, W. M. Neural Machine
823 Reading Comprehension: Methods and Trends. *Appl. Sci.-Basel* **9**, 3698-3742,
824 doi:10.3390/app9183698 (2019).
- 825 41 Rayner, K. & McConkie, G. W. What guides a reader's eye movements? *Vision*
826 *Res.* **16**, 829-837 (1976).
- 827 42 Zeiler, M. D. & Fergus, R. Visualizing and understanding convolutional
828 networks. In *Proc. European conference on computer vision*, 818-833
829 (Springer, 2014).
- 830 43 Lin, Y., Tan, Y. C. & Frank, R. Open Sesame: Getting Inside BERT's
831 Linguistic Knowledge. In *Proc. 2019 ACL Workshop BlackboxNLP: Analyzing*
832 *and Interpreting Neural Networks for NLP*, 241-253 (Association for
833 Computational Linguistics, 2019).
- 834 44 Jawahar, G., Sagot, B. & Seddah, D. What does BERT learn about the
835 structure of language? In *Proc. 57th Annual Meeting of the Association for*
836 *Computational Linguistics*, 3651-3657 (Association for Computational
837 Linguistics, 2019).
- 838 45 Dayan, P., Kakade, S. & Montague, P. R. Learning and selective attention. *Nat.*
839 *Neurosci.* **3**, 1218-1223 (2000).
- 840 46 Yu, A. J. & Dayan, P. Inference, attention, and decision in a Bayesian neural
841 architecture. In *Proc. Advances in neural information processing systems*,
842 1577-1584 (MIT Press, 2005).
- 843 47 Rao, R. P. Bayesian inference and attentional modulation in the visual cortex.
844 *Neuroreport* **16**, 1843-1848 (2005).
- 845 48 Just, M. A. & Carpenter, P. A. A theory of reading: From eye fixations to
846 comprehension. *Psychol. Rev.* **87**, 329-354 (1980).
- 847 49 Rayner, K. Eye movements in reading and information processing: 20 years of
848 research. *Psychol. Bull.* **124**, 372-422 (1998).
- 849 50 Reichle, E. D., Pollatsek, A., Fisher, D. L. & Rayner, K. Toward a model of

850 eye movement control in reading. *Psychol. Rev.* **105**, 125-157 (1998).

851 51 Engbert, R., Nuthmann, A., Richter, E. M. & Kliegl, R. SWIFT: a dynamical
852 model of saccade generation during reading. *Psychol. Rev.* **112**, 777-813
853 (2005).

854 52 Gottlieb, J., Hayhoe, M., Hikosaka, O. & Rangel, A. Attention, reward, and
855 information seeking. *J. Neurosci.* **34**, 15497-15504 (2014).

856 53 van der Schoot, M., Vasbinder, A. L., Horsley, T. M. & van Lieshout, E. C. D.
857 M. The role of two reading strategies in text comprehension: An eye fixation
858 study in primary school children. *J. Res. Read.* **31**, 203-223 (2008).

859 54 Kaakinen, J. K., Hyönä, J. & Keenan, J. M. How prior knowledge, WMC, and
860 relevance of information affect eye fixations in expository text. *J. Exp.*
861 *Psychol.-Learn. Mem. Cogn.* **29**, 447-457 (2003).

862 55 Bahdanau, D., Cho, K. & Bengio, Y. Neural machine translation by jointly
863 learning to align and translate. In *Proc. 3rd International Conference on*
864 *Learning Representations*, (ICLR, 2015).

865 56 Chen, D., Bolton, J. & Manning, C. D. A thorough examination of the
866 cnn/daily mail reading comprehension task. In *Proc. 54th Annual Meeting of*
867 *the Association for Computational Linguistics (Volume 1: Long Papers)*, 2358-
868 2367 (Association for Computational Linguistics, 2016).

869 57 Zhang, Y., Marshall, I. & Wallace, B. C. Rationale-augmented convolutional
870 neural networks for text classification. In *Proc. 2016 Conference on Empirical*
871 *Methods in Natural Language Processing*, 795-804 (Association for
872 Computational Linguistics, 2016).

873 58 Wang, Y., Huang, M. & Zhao, L. Attention-based LSTM for aspect-level
874 sentiment classification. In *Proc. 2016 conference on empirical methods in*
875 *natural language processing*, 606-615 (Association for Computational
876 Linguistics, 2016).

877 59 Zhu, H., Wei, F., Qin, B. & Liu, T. Hierarchical attention flow for multiple-
878 choice reading comprehension. In *Proc. AAAI Conference on Artificial*
879 *Intelligence*, 6077-6085 (AAAI, 2018).

880 60 Wiegrefe, S. & Pinter, Y. Attention is not not Explanation. In *Proc. 2019*
881 *Conference on Empirical Methods in Natural Language Processing and the*
882 *9th International Joint Conference on Natural Language Processing (EMNLP-*
883 *IJCNLP)*, 11-20 (Association for Computational Linguistics, 2019).

884 61 Serrano, S. & Smith, N. A. Is Attention Interpretable? In *Proc. 57th Annual*
885 *Meeting of the Association for Computational Linguistics*, 2931-2951
886 (Association for Computational Linguistics, 2019).

887 62 Zou, J. *et al.* PALRACE: Reading Comprehension Dataset with Human Data
888 and Labeled Rationales. Preprint at <https://arxiv.org/abs/2106.12373> (2021).

889 63 Brainard, D. H. The psychophysics toolbox. *Spatial vision* **10**, 433-436 (1997).

890 64 Wolf, T. *et al.* HuggingFace's Transformers: State-of-the-art natural language
891 processing. In *Proc. 2020 Conference on Empirical Methods in Natural*
892 *Language Processing: System Demonstrations*, 38-45 (Association for
893 Computational Linguistics, 2020).

894 65 Burnard, L. *The British National Corpus, version 3 (BNC XML Edition)*.
895 <<http://www.natcorp.ox.ac.uk/>> (2007).

896 66 Efron, B. & Tibshirani, R. J. *An introduction to the bootstrap*. (CRC press,
897 1994).

898

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SMNC0814.pdf](#)