

Task Optimization Leads to Human-like Top-down and Bottom-up Attention during Reading Comprehension

Jiajie Zou

College of Biomedical Engineering and Instrument Sciences, Zhejiang University

Nai Ding (✉ ding_nai@zju.edu.cn)

College of Biomedical Engineering and Instrument Sciences, Zhejiang University

Article

Keywords: reading comprehension, information processing, human attention, deep neural network

Posted Date: November 1st, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1034025/v1>

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

[Read Full License](#)

1 **Task Optimization Leads to Human-like Top-down and Bottom-up Attention**
2 **during Reading Comprehension**

3

4

5

6

Jiajie Zou¹, Nai Ding^{1,2*},

7

8

¹ Key Laboratory for Biomedical Engineering of Ministry of Education,

9

College of Biomedical Engineering and Instrument Sciences,

10

Zhejiang University, Hangzhou 310027, China

11

² Zhejiang lab, Hangzhou 311121, China

12

13

14

15 ***Corresponding author:**

16

Nai Ding,

17

Email: ding_nai@zju.edu.cn

18

Zhejiang lab, Hangzhou 311121, China

19 **Abstract**

20 Attention is a complex system involving multiple interactive components that jointly
21 regulate information processing in the brain. It has been hypothesized that the
22 computational goal of attention is to optimally integrate information under task
23 demands, and evidence has been provided in relatively simple learning and decision
24 making tasks. It remains unclear, however, whether this hypothesis can explain
25 attention distribution in more complex real-world tasks that engage multiple attention
26 systems. Here, taking advantage of the development of attention mechanisms in deep
27 neural network (DNN) models, we investigate whether human attention during real-
28 world reading comprehension tasks can be explained as a consequence of task
29 optimization. In a goal-directed reading task, participants read a passage to answer a
30 question. Eye tracking results show that the attention on each word, quantified by the
31 fixation time, is modulated by both the top-down reading goal and lower-level visual
32 layout and textual features. When trained to perform the same goal-directed reading
33 task, DNN models yield human-level performance and naturally evolve human-like
34 attention distribution, with deep layers tuned to the reading goal and shallow layers
35 tuned to textual features. Further experiments suggest that different training tasks
36 separately contribute to goal-directed and text-based attention. In summary, the results
37 strongly suggest that human attention can be interpreted as a consequence of task
38 optimization during real-world reading tasks.

39

40 **Introduction**

41 Attention is a key mechanism to regulate information processing in the brain and
42 exerts profound influences on perception and cognition^{1,2}. Instead of a single module,
43 the attention system can be functionally divided into multiple components³⁻⁵. For
44 example, goal-directed top-down attention is mainly guided by the task, while
45 stimulus-driven bottom-up attention is mainly determined by the physical features of
46 stimuli⁶. In terms of the neural implementation, it has been suggested that the control
47 of top-down and bottom-up attention engages separable cortical pathways^{7,8} and these
48 attention control pathways are widely connected to other cortical and subcortical areas
49 to effectively control neural processing across sensory modalities at multiple levels
50 along the processing hierarchy⁹⁻¹¹.

51

52 Since attention has multiple function components and engages broad cortical areas, it
53 remains unresolved whether attention plays a unified computational role in
54 information processing. Traditionally, attention is viewed as a mechanism to allocate
55 limited central processing resources^{12,13}. More recent studies, however, propose that
56 attention is a mechanism to optimize learning and task performance, even in
57 conditions where the processing resource is not clearly constrained¹⁴⁻¹⁷, and this
58 optimization view has been used to explain the attention distribution in a range of
59 learning and decision making tasks¹⁸⁻²⁰. Little empirical evidence has been provided,
60 however, whether task optimization can explain attention distribution in more

61 ecologically valid complex tasks that engage multiple forms of attention, since it is
62 generally challenging to solve optimization problems for complex tasks.

63

64 Here, taking advantage of recent development of attention mechanisms in DNN^{21,22},
65 we investigate whether task optimization is sufficient to explain human attention
66 during a real-world challenging reading task, i.e., answering high-school level reading
67 comprehension questions²³. DNN are powerful tools to solve complex optimization
68 problems and, using DNN previous studies have shown that task optimization can
69 indeed explain, for example, the response properties of individual neurons²⁴⁻²⁶. We
70 choose the reading comprehension task for 3 reasons. First, the task strongly engages
71 attention, since finding the answer to a question in a long passage poses a challenge
72 for information selection. Second, the task is of high ecological validity: Information
73 searching is one of the most common reading activities in the modern time²⁷ and
74 people often receive practice for this task in school. Third and importantly, DNNs
75 equipped with the self-attention mechanism have reached human-level performance
76 on this task.

77

78 The purpose of our study is to investigate whether DNNs optimized for task
79 performance can generate attention distribution that is quantitatively similar to the
80 attention distribution of humans who perform the same task. We quantified human
81 attention through eye tracking and quantified DNN attention using the attention
82 weight on each word. In the following, we first report how human attention is

83 modulated by stimulus features and the top-down reading goal, and then analyze
84 whether DNN evolves bottom-up and top-down attention like humans. Lastly, we
85 analyze how the training tasks contribute to the distribution of DNN attention and and
86 its similarity to human attention.

87

88

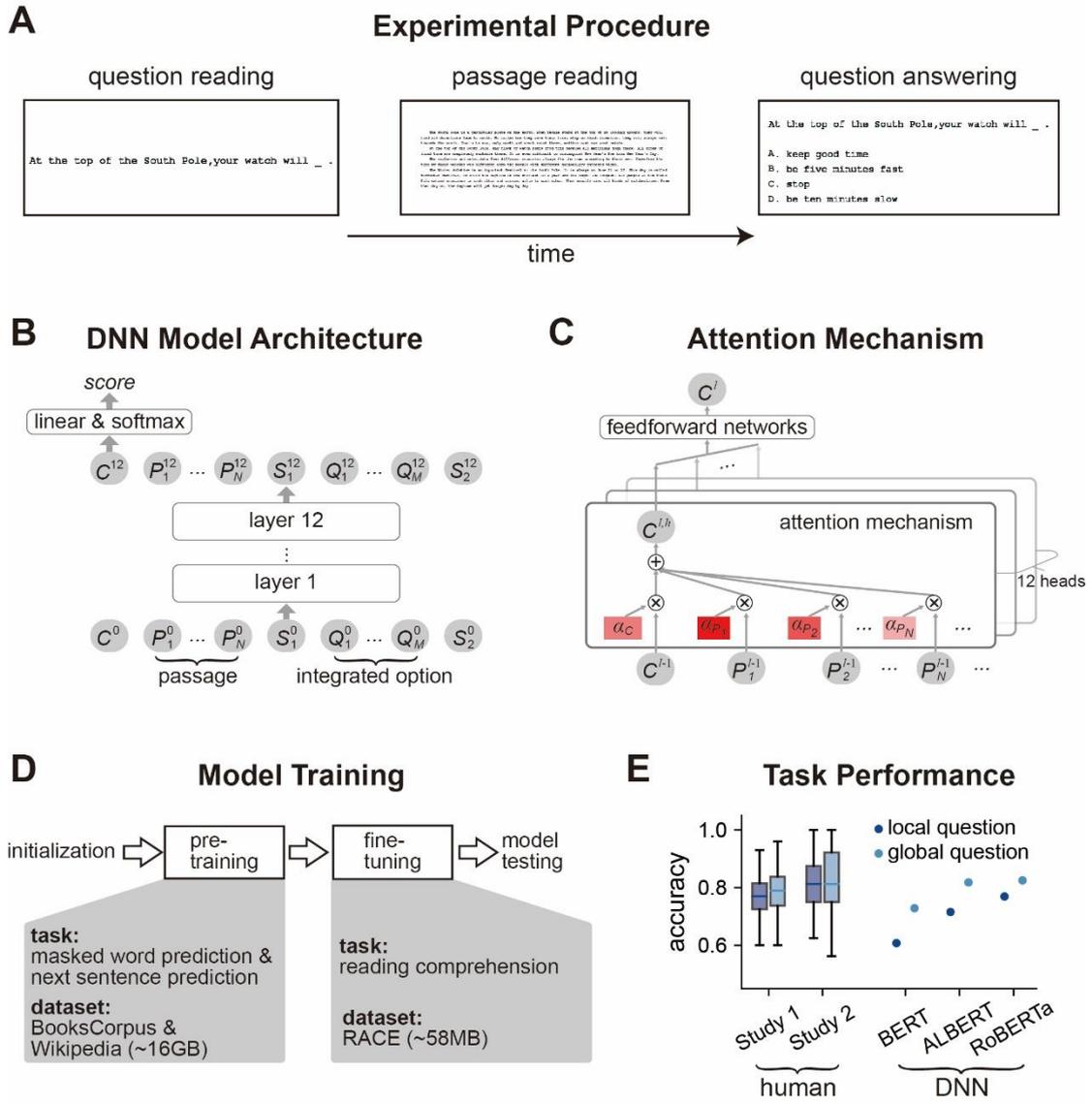
89 **Results**

90 **Human Attention Distribution and Influence Factors**

91 In Study 1, the participants ($N = 25$ for each question) first read a question and then
92 read a passage based on which the question should be answered (Fig. 1A). After
93 reading the passage, the participants read 4 options related to the question and had to
94 choose which option was the most suitable answer. Eight hundred question/passage
95 pairs were presented, and the questions fell into two broad categories, i.e., local and
96 global questions (see *Materials and Methods* for details). Local questions require
97 attention to details while global questions concern the general understanding of a
98 passage. The participants correctly answered 77.94% questions on average (Fig. 1E,
99 77.49% and 78.77% for local and global questions, respectively).

100

101



102

103

104 **Fig. 1. Experimental procedure and DNN model.** (A) The experimental procedure
 105 in Study 1. In each trial, participants read a question first, and then read the
 106 corresponding passage, and finally proceeded to read and answer the question which
 107 is coupled with 4 options. Here, the question and options are shown in a larger font
 108 size for illustration purposes. (B) Architecture of the DNN models used for the
 109 reading comprehension task. The input to the models consists of all words in the
 110 passage and an integrated option, and also 3 special tokens, i.e., CLS, SEP₁, and SEP₂
 111 (denoted as C, S₁, and S₂), separating the passage and the integrated option. The CLS
 112 token is the decision variable and its final representation is used to calculate a score
 113 that reflects how likely the option is the correct answer. The DNN model has 12
 114 layers, and each layer has 12 parallel self-attention modules, i.e., attention heads. (C)
 115 Illustration of the DNN attention mechanism in a layer. In the models, each
 116 word/token is represented by a vector, and information is integrated across

117 words/tokens only in the self-attention module. The vectorial representation of the
118 CLS token is a weighted sum of the vectorial representations of all words and tokens,
119 and the attention weight for each word in the passage, i.e., α_{P_n} , is the DNN attention
120 analyzed in this study. Output of the self-attention model, i.e., $C^{l,h}$, is further
121 processed by feedforward networks and other operations that do not engage
122 information integration across words. **(D)** DNN models are trained in two steps. The
123 pre-training process aims to learn general statistical regularities in a language based
124 on large corpora, while the fine-tuning process trains models to perform the reading
125 comprehension task. **(E)** Performance of humans and DNN models on the reading
126 comprehension task.

127

128

129 While the participants read the passage, their eye gaze was monitored using an eye
130 tracker, and their attention to each word was quantified by the total fixation time on
131 the word. The results showed that longer words were fixated for longer time (Fig. S1),
132 consistent with previous studies²⁸. Nevertheless, the fixation time on a word was
133 expected to be positively related with the area the word occupied even when attention
134 was uniformly distributed across the visual field. Therefore, here we further extracted
135 the *attention density* by dividing the total fixation time on a word by the area the word
136 occupied, and used this measure in subsequent analyses. The attention density clearly
137 deviated from a uniform distribution (see Fig. 2 for examples). To probe into the
138 factors modulating human attention distribution, we quantified how the human
139 attention distribution was influenced by multiple sets of features in the following.

140

141 We first analyzed whether textual features, e.g., word length, word frequency, and a
142 word's position in a sentence, could predict human attention distribution using linear

143 regression. The prediction accuracy, i.e., the correlation coefficient between the
144 predicted and actual attention density, was significantly above chance ($P = 0.002$,
145 permutation test, FDR corrected). Furthermore, the prediction accuracy was
146 significantly higher for global questions than for local questions ($P = 1.4 \times 10^{-4}$,
147 bootstrap, FDR corrected) (Fig. 3A, the left plot). We then used the same regression
148 analysis to analyze whether the visual layout of a passage could also affect attention
149 distribution. Here, layout features referred to features induced by line changes (see
150 *Materials and Methods* for details), which could be processed without word
151 recognition. The prediction accuracy for layout features was also statistically
152 significant ($P = 0.002$, permutation test, FDR corrected).

153

154 Textual features and layout features characterized properties of the stimulus that were
155 invariant across tasks. In the following, we investigated whether the task, i.e., to
156 answer a specific question, also modulated human attention distribution. To
157 characterize the top-down influence of task, we acquired annotations indicating each
158 word's contribution to question answering, i.e., task relevance (see *Materials and*
159 *Methods*). As shown in the left plot of Fig. 3A, we found that task relevance could
160 indeed significantly predict human attention distribution ($P = 0.002$, permutation test,
161 FDR corrected). Since task relevance was not a well-established modulator of reading
162 attention, we further analyzed whether the task relevance effect could be explained by
163 the well-established textual and layout effects. In this analysis, we first regressed out
164 the influence of textual and layout features from the human attention distribution, and

165 found that the residual attention distribution could still be predicted by task relevance
166 ($P = 0.003$, permutation test, FDR corrected) (Fig. 3A, middle plot). These results
167 showed that the top-down reading goal, quantified by task relevance, could modulate
168 human attention, on top of lower-level stimulus features, i.e., textual and layout
169 features.

170

171 The linear regression analyses revealed that textual features, layout features, and task
172 relevance all modulated human attention (see Fig. 2 for examples). The prediction
173 accuracy for different features ranged between 0.2 and 0.6, comparable to the
174 prediction accuracy of visual saliency models when predicting human attention to
175 images^{29,30}. Further analyses also revealed how these features modulated human
176 attention. For example, we found that participants generally attended more to the
177 beginning of a passage (Fig. 3C). Furthermore, this effect was stronger for global
178 questions, which potentially explained why stimulus features could better predict the
179 attention distribution for global questions. Additionally, it was also found that
180 participants attended more to words that are more relevant to the question answering
181 task (Fig. 3D).

182

183

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.

prediction based on layout 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.

prediction based on DNN attention

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

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 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.

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.

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.

prediction based on layout 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.

prediction based on DNN attention

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 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.

DNN attention in the last layer

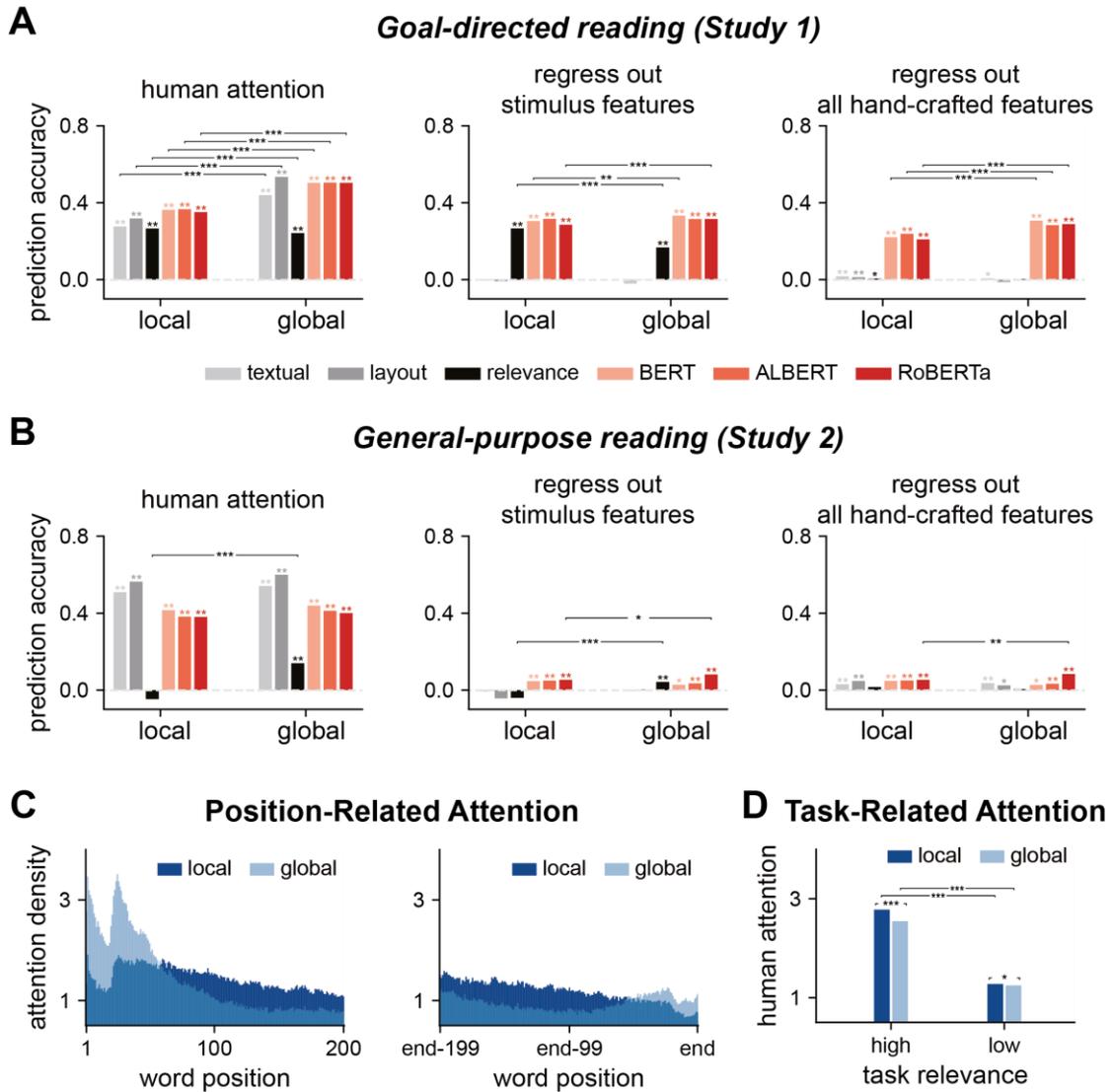
All the oceans of the world will be dead in the future unless action is taken at once. How can this happen? We have already seen that people allow all sorts of waste products to flow into the sea. It is almost impossible to measure how much waste water and industrial waste in our oceans, but we can find 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.

184

185 **Fig. 2. Examples of human attention distribution and the attention distribution**
186 **predicted by different features.** Panels A and B separately show the attention
187 distribution for two passages and the corresponding questions are shown in the
188 parenthesis. Human attention is quantified by the total fixation time per unit area.
189 Textual features include word properties, e.g., word frequency and the position of a
190 word in the passage. Layout features include visual features that can be processed
191 without recognizing individual words, e.g., line changes. Task relevance refers to
192 human annotation about the contribution of each word to question answering. DNN
193 attention includes all the layers and attention heads, and the DNN attention in the last
194 layer is shown separately (average over attention heads).

Predicting Human Attention Based on Different Features



195

196

197 **Fig. 3. Predicting human attention using different features.** (A and B) Panels A
 198 and B show the results of Study 1 and Study 2, respectively. The left plots show how
 199 well different sets of features can predict human attention. In the middle and right
 200 plots, some features are regressed out from human attention, and the residual human
 201 attention is predicted by other features. Prediction accuracy that is significantly higher
 202 than chance is denoted by stars of the same color as the bar. (C) The influence of
 203 word position on human attention. Humans generally attend more to the beginning of
 204 a passage, especially for global questions. (D) The influence of task relevance on
 205 human attention. Humans allocate more attention to words that are more relevant to
 206 question answering. * $P < 0.05$; ** $P < 0.01$; *** $P < 0.001$.

207

208

209 **Attention Distributions in Humans and DNN**

210 We then investigated whether DNN can generate attention that was comparable to
211 human attention, when optimized to perform the reading task. The general
212 architecture of the models was illustrated in Fig. 1B. The input to the models included
213 all the words in the passage, integrated option, and 3 special tokens. The integrated
214 option was the question concatenated with an option. One of the special token, i.e.,
215 CLS, was the decision variable, based on the final representation of which the DNN
216 models decided whether an option was the correct answer or not. In the following, we
217 analyzed the attention weight between the CLS token and each word in the passage
218 (see *Materials and Methods* for details). In each layer of the DNN models, the
219 vectorial representation of the CLS token was updated by a weighted sum of the
220 vectorial representations of all input words and tokens. Therefore, the attention weight
221 on a word could reflect how heavily the word contributed to the decision variable, i.e.,
222 the CLS token.

223

224 We analyzed 3 DNN models, i.e., BERT³¹, ALBERT³², and RoBERTa³³, and the
225 question answering performance of the 3 DNN models was within the range of human
226 performance (Fig. 1E). Each of the 3 DNN models had 12 layers and each layer had
227 12 parallel modules that were referred to as attention heads (Fig. 1BC). Each attention
228 head had a separate set of attention weights, and the 12 attention heads were
229 integrated within each layer. Consequently, each word had 144 attention weights (12

230 layers \times 12 heads). In the following, we first tested whether the DNNs could generate
231 human-like attention distributions by attempting to decode human attention
232 distribution from the 144 DNN attention weights using linear regression. Then, we
233 analyzed whether the attention weights in different layers showed different attention
234 properties.

235

236 Although the DNN models were only trained to perform the reading comprehension
237 task and were blind to the human fixation data, it was found that the DNN attention
238 weights could significantly predict human attention distribution ($P = 0.002$,
239 permutation test, FDR corrected), and the prediction accuracy was higher for global
240 questions than for local questions ($P = 1.4 \times 10^{-4}$, bootstrap, FDR corrected) (Fig. 3A,
241 left). The prediction accuracy of DNN attention weights was higher than that of
242 textual features and task relevance. Compared with the predictions based on layout
243 features, the predictions based on DNN attention weights were higher for local
244 questions and lower for global questions. It is worth noting that layout features, which
245 were induced by line changes, were not available in the input to DNN models.

246

247 DNN attention weights could model the human attention distribution, but did they
248 capture information beyond the hand-crafted features, i.e., textual features, layout
249 features, and task relevance features? We found that when the influences of textual
250 and layout features were regressed out, the residual human attention distribution could
251 still be explained by the DNN attention weights (Fig. 3A, the middle plot). This result

252 suggested that the DNN attention weights contained information beyond basic
253 stimulus features. Additionally, when the stimulus features and task relevance were
254 both regressed out, the residual human attention distribution remained significantly
255 predicted by the DNN attention weights (Fig. 3A, the right plot). Therefore, DNN
256 attention weights can model human attention and capture information beyond basic
257 hand-crafted features.

258

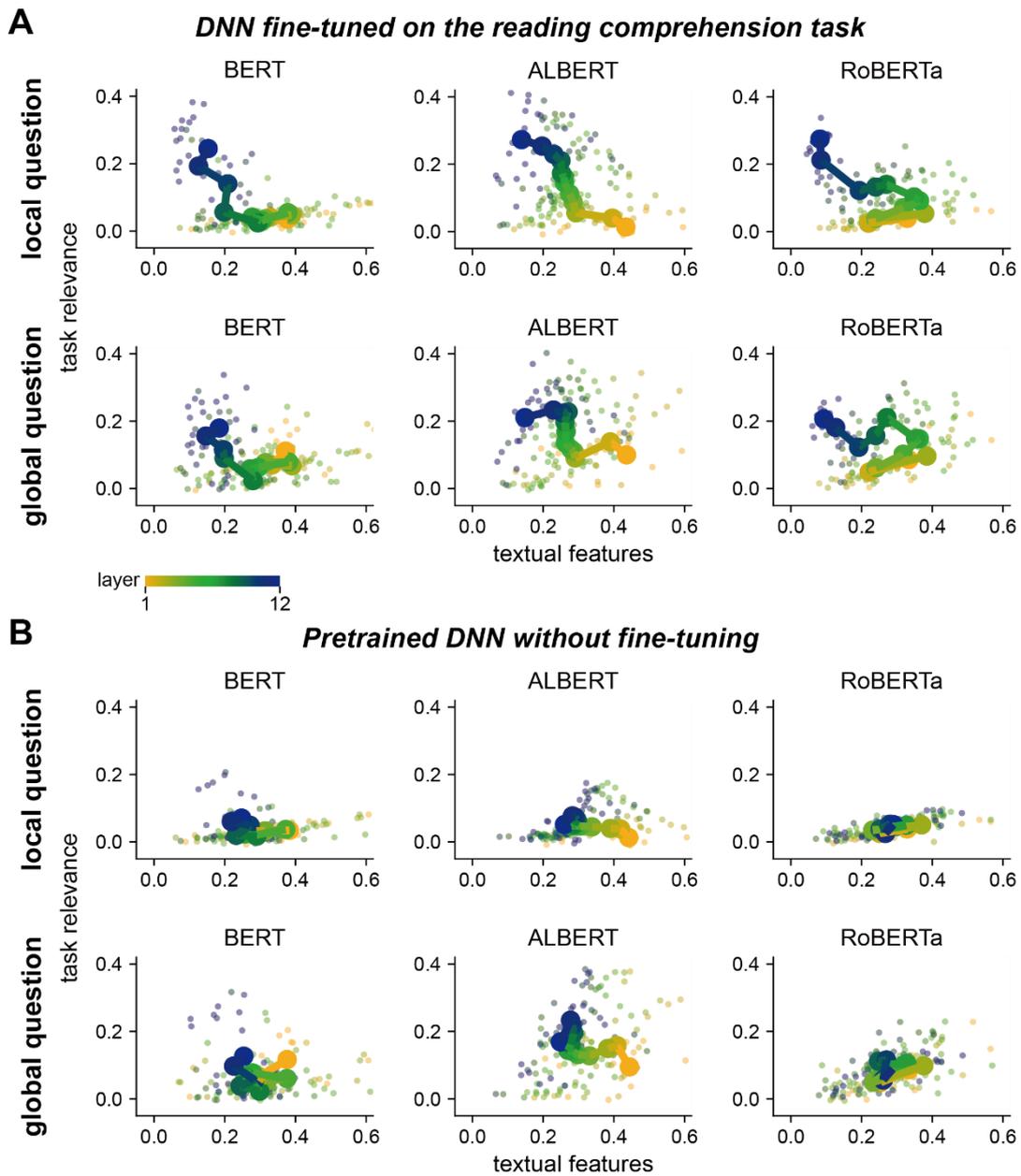
259 **Task Modulation in Humans**

260 To further confirm that human attention received top-down modulation from the task,
261 we conducted Study 2 as a control study. In Study 2, participants first read a passage
262 without prior knowledge about the question to answer. After the first-pass passage
263 reading, the participants read the question and were then allowed to read the passage
264 again before answering the question. We analyzed the attention density during the
265 first-pass reading of the passage, which was referred to as general-purpose reading.

266

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

Properties of DNN Attention in Different Layers



274

275

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

285 Comparing the results obtained from Study 1 and Study 2, it was evident that human
286 attention could be modulated by the specific reading goal, i.e., the question to answer,
287 on top of textual and layout features. Goal-directed top-down attention, characterized
288 in Study 1, could be modeled by either human-annotated task relevance or the DNN
289 attention weights. In the absence of a specific reading goal, human attention in Study
290 2 was mainly influenced by stimulus features, e.g., textual and layout features, which
291 were also captured by the DNN attention weights.

292

293 **DNN Attention in Different Layers**

294 In the human attention system, bottom-up and top-down attention are implemented in
295 separate neural pathways^{7,8}, and in the following, we analyzed whether attention to
296 stimulus features and attention to task relevant information are also implemented in
297 separate parts of DNN. Since previous studies have shown that different layers in
298 DNNs encoded different types of information³⁴, we analyzed whether the properties
299 of DNN attention weights differed across layers. Since the layout features were not
300 available to the DNNs, we only considered textual features as stimulus features in this
301 analysis. As shown in Fig. 4A, the attention weights in different layers were sensitive
302 to different features. In general, shallow layers were more strongly influenced by
303 textual features while deeper layers were more strongly influenced by the task
304 relevance. This trend was observed in all 3 DNN models and was especially obvious
305 for local questions. The transitional trajectory across layers, however, was model-
306 dependent in the 2-dimensional feature space. In Fig. 2, examples were shown for the

307 attention weights averaged over all 12 heads in the last layer of BERT, which
308 resembled the human-annotated task relevance.

309

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

311 All the 3 DNN models were pre-trained based on large-scale corpora and fine-tuned
312 based on the reading comprehension task (Fig. 1D). How was the DNN attention
313 distribution separately influenced by the pre-training and fine-tuning processes? We
314 addressed this question by analyzing the attention weights in pre-trained DNN models
315 that did not receive fine-tuning (Fig. 4B). It was found that the attention weights of
316 pre-trained DNNs were sensitive to textual features in shallow layers but not sensitive
317 to task relevance in deeper layers. This result suggests that top-down attention in
318 DNNs emerged during fine-tuning using the reading comprehension task, while text-
319 based attention existed after pre-training.

320

321 We then asked how the attention weights of DNNs changed during fine-tuning and
322 whether such changes were related to the performance of question answering. In the
323 following, we analyzed the properties of models that received different steps of fine-
324 tuning. Furthermore, since fine-tuning process was stochastic, we fine-tuned 10 times
325 (see *Materials and Methods*). We found that, in deep layers, the properties of attention
326 weights significantly changed during fine-tuning (Fig. S3 and Fig. S4). In the last
327 layer, for example, it was clear that the DNN attention weights became more sensitive
328 to task relevance during fine-tuning, coinciding with the improvement in task

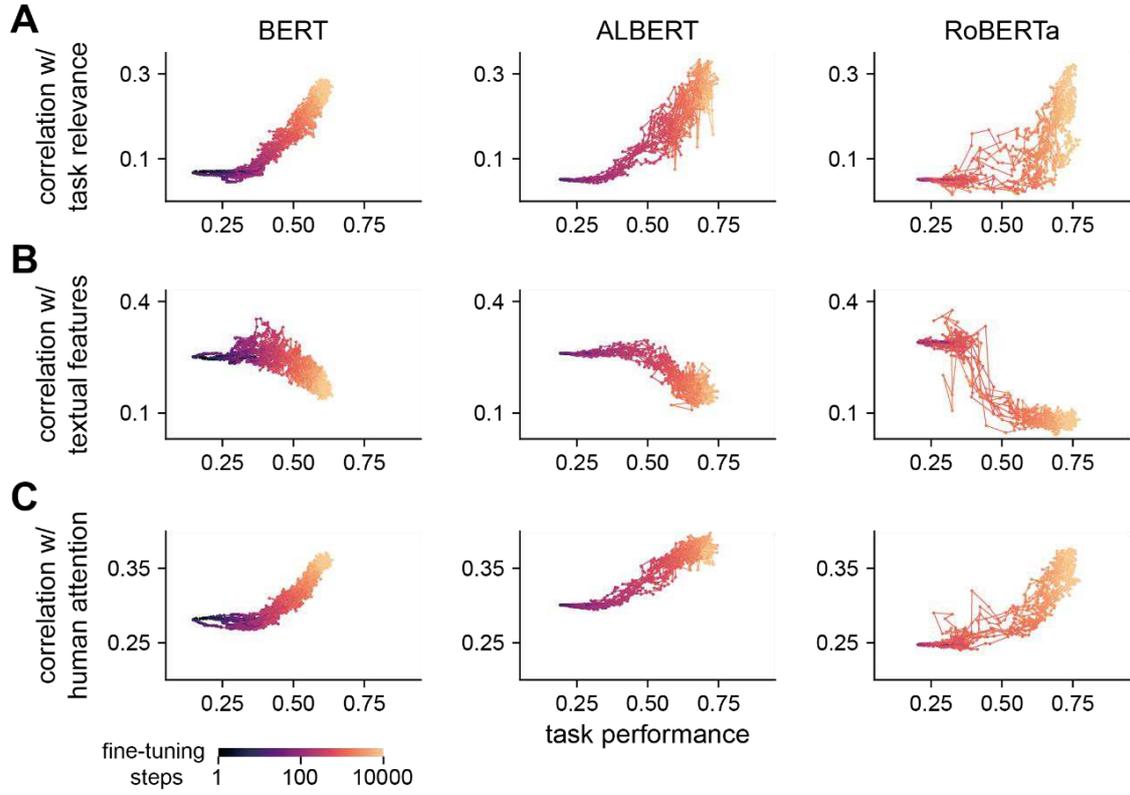
329 performance (Fig. 5AD), especially for local questions (Fig. 5A). The trend is less
330 clear for global questions and a potential explanation is that global questions concern
331 the main topic of the passage and can be answered by paying attention to different
332 sets of words. Deep layer's sensitivity to textual features, however, dropped during
333 fine-tuning (Fig. 5BE). Therefore, fine-tuning directed deep layers' attention towards
334 task-relevant information, compromising the sensitivity to textual features.

335 Additionally, we found that the similarity between DNN attention weights and human
336 attention was also boosted by fine-tuning for local but not global questions (Fig.
337 5CF). These results show that fine-tuning boosted task performance by generating
338 goal-directed human-like attention.

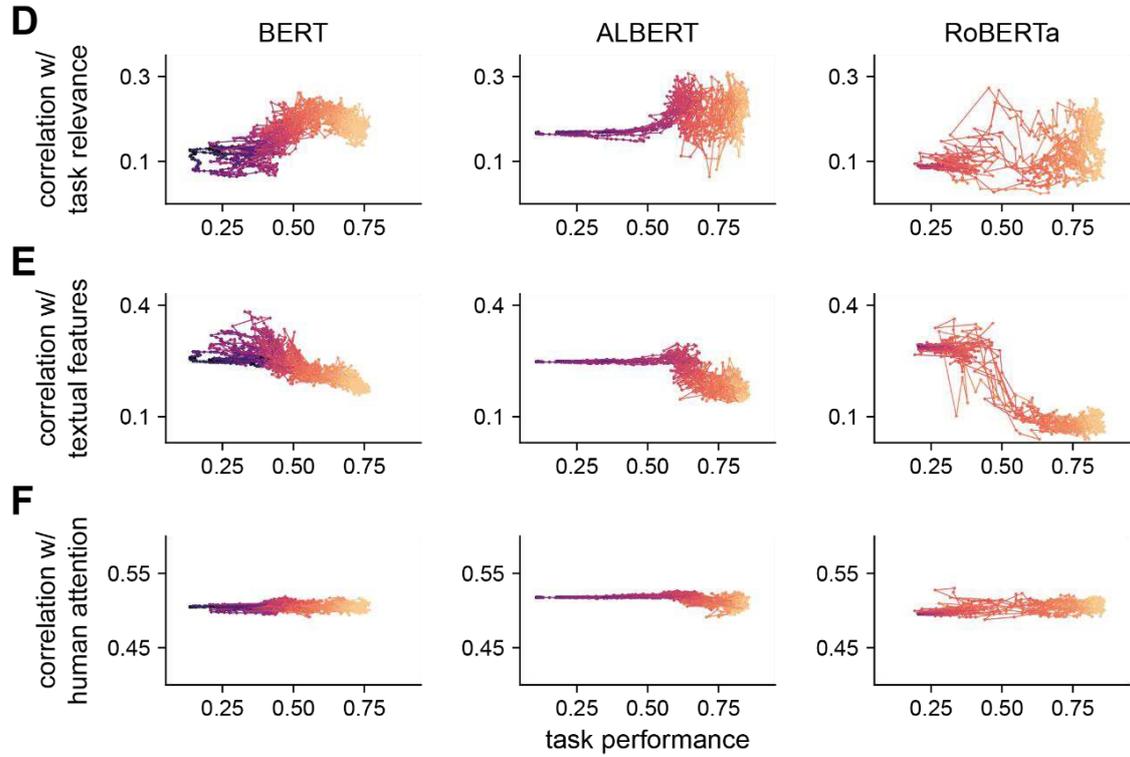
339

Influence of Fine-Tuning on DNN Attention and Task Performance

Local question



Global question



340

341

342

343 **Fig. 5. Influence of fine-tuning on DNN attention and task performance.** Each
344 model is fine-tuned 10 times. Each data point denotes the result at a fine-tuning step
345 (color coded), and steps from each run of fine-tuning are connected by a line. (A, B,
346 D, and E) The effect of fine-tuning on the attention weights in the last layer of DNNs
347 for local questions (A and B) and global questions (D and E). Fine-tuning enhances
348 the sensitivity to top-down task relevance while reducing the sensitivity to lower-level
349 textual features, which correlates with the increase in task performance. (C and F)
350 Influence of fine-tuning on the similarity between DNN and human attention. For
351 local questions (C), fine-tuning clearly increases the similarity between DNN and
352 human attention, coinciding with the increase in task performance. For global
353 questions (F), the similarity between DNN and human attention is high even without
354 fine-tuning and is not further boosted by fine-tuning. For ALBERT, 2 out of the 10
355 runs of fine-tuning are unstable, showing sharp drops in task performance during fine-
356 tuning. Results of these 2 runs are not shown here but separately shown in Fig. S2.

357

358

359

360

361 **Discussion**

362 The current study quantifies human attention during a real-world reading
363 comprehension task, and demonstrates that DNN models optimized to perform the
364 same task generate human-like attention distribution. Similar to the fact that the
365 human attention system is divided into dual streams separately controlling bottom-up
366 and top-down attention, the shallow and deep layers of DNNs separately implement
367 text-based and goal-directed attention. These results strongly suggest that task
368 optimization provides a computational-level explanation for human attention
369 distribution even in real-world complex reading tasks.

370

371 Although rarely investigated, goal-directed reading comprehension provides a suitable
372 paradigm to study attention, and has close relationship to visual search and other
373 information seeking tasks. In the current study, for example, each passage contains a
374 few hundred of words but usually only a few words directly contribute to question
375 answering, highlighting the need for information selection. Furthermore, attributable
376 to the crowding effect³⁵, the participant could only recognize a couple of words in one
377 fixation and therefore question-relevant key words cannot pop out. Consequently, the
378 participants have to sequentially sample a passage until they could construct their
379 answer. The sampling process, however, is highly active, as evidenced by the uneven
380 distribution of attention. Therefore, goal-directed reading provides a real-world
381 scenario to study attention during complex information seeking behavior³⁶.

382

383 **Computational models of attention**

384 Since attention is a key mechanism in the brain, a large number of computational
385 models of attention have been proposed. According to Marr's 3 levels of analysis³⁷,
386 some models investigate the computational goal of attention. It has been hypothesized
387 that attention can be interpreted as a mechanism to implement optimal learning and
388 decision making. For example, the brain generally attends to information sources that
389 can reduce the uncertainty of a decision¹⁴⁻¹⁶, and a higher degree of attention can
390 enhance the fidelity of neural coding and influence decision^{20,38-40}. The current study
391 supports the view that attention is a mechanism to optimize task performance and
392 shows that task optimization can lead to both stimulus-driven bottom-up and goal-
393 directed top-down attention, and can predict human attention in a challenging real-
394 world reading task.

395

396 The DNN models analyzed here are trained in two steps (Fig. 1D). The results show
397 that the pre-training process leads to text-based attention while the fine-tuning process
398 leads to goal-directed attention (Fig. 4 and Fig. 5). The pre-training process is also
399 implemented through task optimization. One pre-training task that is used in all 3
400 models analyzed here is to predict a word based on the context. Another task used to
401 pre-train BERT and ALBERT is to predict the order of two sentences selected from a
402 passage. The purpose of pre-training tasks is to let models learn the general statistical
403 regularity in a language based on large corpora⁴¹, and this process has greatly
404 promoted model performance³¹⁻³³. The results here show that the predictive pre-

405 training tasks can lead to text-based attention. Previous studies have also suggested
406 that predictivity can explain reading time^{42,43} and attention to images⁴⁴.

407

408 A separate class of models investigate what are the factors contributing to human
409 attention distribution. A large number of models are proposed to predict bottom-up
410 visual saliency^{29,30,45}, and recently DNN models are also employed to model top-down
411 visual attention. It is shown that, through either implicit^{46,47} or explicit training^{48,49},
412 DNNs can predict which parts of a picture relates to a verbal phrase, a task similar to
413 goal-directed visual search⁵⁰. The current study distinguishes from these studies in
414 that the DNN model is not trained to predict human attention. Instead, the DNN
415 models naturally generate human-like distribution when trained to perform the same
416 task that humans perform, suggesting that task optimization is a potential cause for the
417 attention distribution during reading.

418

419 **Attention during human reading**

420 How human readers allocate attention during reading is an extensively studied topic.
421 Eye tracking studies have shown that the readers fixate longer at longer words, words
422 of lower-frequency, words that are less predictable based on the context, and words at
423 the beginning of a line, etc⁵¹. A number of models, e.g., the E-Z reader^{52,53} and
424 SWIFT⁵⁴, have been proposed to predict the eye movements during reading, either
425 based on basic oculomotor properties or lexical processing⁵³. Some models are ideal
426 observer models that minimize the time or the number of saccades required to read a

427 sentence^{55,56}, but these models do not model how the reading time is affected by the
428 reading goal.

429

430 Traditional eye movement models can generate fine-grained predictions, e.g., which
431 letter in a word will be fixated first, for the reading of simple sentences, but have only
432 been occasionally tested for complex sentences or multi-line text⁵⁷. In contrast to
433 these studies, the current study focuses on goal-directed reading of complex multi-line
434 text and models attention in the units of words. Future studies can potentially integrate
435 classic eye movement models with DNNs to explain the dynamic eye movement
436 trajectory, possibly with a letter-based spatial resolution.

437

438 In most previous studies, readers are instructed to read a sentence in a normal manner,
439 not aimed to extract a specific kind of information⁵⁸. In the current study, however,
440 readers know in advance what question they have to answer and this kind of reading
441 is also referred to as the reading-to-do task⁵⁹. Previous studies have shown the
442 reader's task may have heterogeneous influences on attention, depending on the task
443 difficulty and skill level of readers^{60,61}. Here, the task is demanding and the readers
444 are highly skilled to perform the task. Future work is needed to quantify how the task
445 and reading skills modulate human attention and whether these effects can also be
446 modeled by DNN models. Furthermore, when comparing goal-directed and general-
447 purpose reading, it can be observed that attention is more strongly influenced by
448 textual and visual features during general-purpose reading (Fig. 3AB). This result is

449 consistent with previous results that when the task is to search for a target word, the
450 reading time is much less influenced by, e.g., word frequency compared with normal
451 reading⁶².

452

453 **Attention mechanisms in DNN**

454 In transformer-based DNN models, the roles self-attention plays are highly diverse.
455 Since self-attention assigns a weight between every pair of inputs (including words
456 and special tokens such as CLS), it can capture a number of relationships between
457 words, e.g., co-reference and syntactic dependency^{63,64}. In the current study, however,
458 we only analyze a small portion of the self-attention weights that are directly relevant
459 to the task, i.e., the attention weights between CLS and words in the passage. The
460 attention weights analyzed here can be interpreted as a selective information
461 integration mechanism, describing how different words in a passage contribute to the
462 decision variable, i.e., CLS. A couple of recent studies have also compared human
463 attention with several kinds of attention mechanisms in DNNs, and all found some
464 degree of similarity between DNN attention and human attention⁶⁵⁻⁶⁷. The current
465 study, however, demonstrate how the task and textual features separately modulate
466 human and DNN attention, and also analyze how the DNN attention properties vary
467 across layers and gradually evolve during training.

468

469 Additionally, whether attention can increase the interpretability of DNN models is a
470 heavily debated topic. A number of studies have shown that the DNN attention

471 weights are higher for words that are more important for the task^{68,69}. Most of these
472 studies, however, are based on visual inspection of a couple of examples. Other
473 studies, however, reveal low correlation between attention weights and other
474 measures of the importance of words, and therefore raise concern about the
475 interpretability of attention weights^{70,71}. The importance of a word, however, can be
476 measured in many different ways, and no correlation with some importance measures
477 does not indicate no contribution to task solving in other ways. The current study
478 demonstrated that 3 transformer-based models generate human-like attention through
479 task optimization. It remains unclear whether other DNN models show similar
480 properties. Nevertheless, the dataset and methods developed here can be easily
481 applied to test whether other models also evolve human-like attention distribution,
482 serving as a probe to test the biological plausibility of NLP models.

483

484 In sum, the current study demonstrates that, when DNNs and humans perform the
485 same reading comprehension task with comparable accuracy, the DNN attention
486 weights optimized for task performance resemble human attention indexed by eye
487 movements. The results suggest that human attention distribution is shaped by the
488 demand to optimally perform the task. Furthermore, since explainability has become a
489 main issue for DNN models, the results also suggest that, in some models, DNN
490 attention can be interpreted as an approximation of human attention. The large set of
491 eye tracking data in the current study can also motivate future computational

- 492 modeling of human attention during natural reading tasks and to test whether NLP
- 493 models exhibit human-like attention distribution.

494 **Materials and Methods**

495 **Participants**

496 Study 1 enrolled 102 participants (19-30 years old, mean age, 22.9 years; 54 female).
497 Study 2 enrolled a separate group of 18 participants (21-26 years old, mean age, 23.4
498 years; 10 female). Additional 11 participants were recruited but failed to pass the
499 calibration process for eye tracking and therefore did not participate in the reading
500 experiments. All participants were native Chinese speakers with normal or corrected-
501 to-normal vision and were college students or graduate students from Zhejiang
502 University, and were thus above the level required to answer high-school-level
503 reading comprehension questions. English proficiency levels were further guaranteed
504 by the following criterion for screening participants: a minimum score of 6 on IELTS,
505 80 on TOEFL, or 425 on CET6¹. The experimental procedures were approved by the
506 Research Ethics Committee of the College of Medicine, Zhejiang University (2019–
507 047). The participants provided written consent and were paid.

508

509 **Experimental materials**

510 The reading materials were selected and adapted from the large-scale RACE dataset, a
511 collection of reading comprehension questions in English exams for middle and high
512 schools in China²³. We selected eight hundred high-school level questions from the
513 test set of RACE and each question was associated with a distinct passage (117 to 456

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.

514 words per passage). All questions were multiple-choice questions with 4 alternatives
515 including only one correct option among them. The questions fell into 6 types, i.e.,
516 Cause ($N = 200$), Fact ($N = 200$), Inference ($N = 120$), Theme ($N = 100$), Title ($N =$
517 100), and Purpose ($N = 80$). The Cause, Fact, and Inference questions were concerned
518 with the location, extraction, and comprehension of specific information from a
519 passage, and were referred to as local questions. Questions of Theme, Title, and
520 Purpose tested the understanding of a passage as a whole, and were referred to as
521 global questions.

522

523 In a separate online experiment, we acquired annotations about the relevance of each
524 word to the question answering task. For each passage, a participant were allowed to
525 annotate up to 5 key words that were considered relevant to answering the
526 corresponding question. Each passage was annotated by N participants ($N = 26$),
527 producing N versions of annotated key words. Each version of annotation was then
528 validated by a separate participant. In the validation procedure, the participant was
529 required to answer the question solely based on the key words in a specific annotation
530 version; if the person could not derive the correct answer, this version of annotation
531 was discarded. The percent of questions correctly answered in the validation
532 procedure was 75.9% and 67.6%, for local and global questions respectively. The
533 correct rate for local questions was comparable to the correct rate in the eye tracking
534 experiments, in which the readers had access to the whole passage (Fig. 1), suggesting
535 that as few as 5 key words could well support question answering. The correct rate for

536 global questions was lower than the eye tracking experiments but remained much
537 higher than chance, i.e., 25%. For the M versions of annotation that passed the
538 validation procedure, if a word was annotated in K versions, the relevance of the word
539 to question answering, i.e., task relevance, was M/K . More details about the question
540 types and the annotation procedures could be found in reference⁷².

541

542 **Experimental procedures**

543 **Study 1:** Study 1 included all 800 passages, and different question types were
544 separately tested in different experiments, hence six experiments in total. Each
545 experiment included 25 participants and one participant could participate in multiple
546 experiments. Before each experiment, to familiarize the participants with the
547 experimental procedure and the type of questions to answer, participants were given a
548 familiarization session with 5 questions that were not used in the formal experiment
549 or in the analysis. During the formal experiment, questions were presented in a
550 randomized order. Considering the quantities of questions, for Cause and Fact
551 questions, the experiment was carried out in 3 separate days (one third questions on
552 each day), and for other question types the experiment was carried out in 2 days (fifty
553 percent of questions on each day).

554

555 The experiment procedure in Study 1 was illustrated in Fig. 1A. In each trial,
556 participants first read a question, pressed the space bar to read the corresponding
557 passage, and then pressed it again to read the question coupled with 4 options and

558 answer the question. The time limit for passage reading was 120 s. To encourage the
559 participants to read as quickly as possible, the bonus they received for a specific
560 question would decrease linearly from 1.5 to 0.5 RMB over time. They did not
561 receive any bonus for the question, however, if they gave a wrong answer.
562 Furthermore, before answering the comprehension question, the participants reported
563 whether they were confident that they could correctly answer the question (yes or no).
564 Participants selected yes for 90.47% of questions (89.62% and 92.04% for local and
565 global questions, respectively). After answering the question, they also rated their
566 confidence about their answer on the scale of 1-4 (low to high). The mean confidence
567 rating was 3.25 (3.28 and 3.18 for local and global question, respectively), suggesting
568 that the participants were confident about their answers.

569

570 **Study 2:** Study 2 included 96 reading passages and questions, with 16 questions for
571 each question type that were randomly selected from the questions used in Study 1.

572 The study was carried out in 2 days, and none of the participants participated in Study
573 1. The familiarization procedure was identical to that in Study 1.

574

575 The procedure of Study 2 was similar to that of Study 1, except a difference that a 90-
576 s first-pass passage reading stage was introduced at the beginning of each trial. During
577 the first-pass passage reading, participants had no prior information of the relevant
578 question. The participants could press the space bar to terminate the first-pass reading
579 stage and to read a question. Then, participants read the passage for the second time

580 with a time limit of 30 s, before proceeding to answer the question. The payment
581 method was similar to Study 2, and the bonus was calculated based on the duration of
582 second-pass passage reading.

583

584 **Stimulus presentation and eye tracking**

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

590

591 Eye tracking data were recorded from the left eye with 500-Hz sampling rate (Eyelink
592 Portable Duo, SR Research). The experiment stimuli were presented on a 24-inch
593 monitor (1920 \times 1080 resolution; 60 Hz refresh rate) and administered using
594 MATLAB Psychtoolbox⁷³. Each experiment started with a 13-point calibration and
595 validation of eye tracker, and the validation error was required to be below 0.5° of
596 visual angle. Furthermore, before each trial, a 1-point validation was applied, and if
597 the calibration error was higher than 0.5°, a recalibration was carried out. Head
598 movements were minimized using a chin and forehead rest.

599

600 **DNN models**

601 We tested 3 popular transformer-based DNN models, i.e., BERT³¹, ALBERT³², and

602 RoBERTa³³. ALBERT and RoBERTa were both adapted from BERT, and had the
603 same basic structure. RoBERTa differed from BERT in its pre-training procedure³³
604 while ALBERT applied factorized embedding parameterization and cross-layer
605 parameter sharing to reduce memory consumption³². Following previous works^{32,33},
606 each option was independently processed. For the i^{th} option ($i = 1, 2, 3, \text{ or } 4$), the
607 question and the option were concatenated to form an integrated option. As shown in
608 Fig. 1B, for the i^{th} option, the input to DNNs was the following sequence:

609

$$610 \quad C_i, P_1, P_2, \dots, P_N, S_{i,1}, O_{i,1}, O_{i,2}, \dots, O_{i,M}, S_{i,2},$$

611

612 where $C_i, S_{i,1}$, and $S_{i,2}$ denoted special tokens, i.e., the CLS, SEP₁, and SEP₂ tokens,
613 separating different components of the input. P_1, P_2, \dots, P_N denoted all the N words of
614 a passage, while $O_{i,1}, O_{i,2}, \dots, O_{i,M}$ denoted all the M words of the i^{th} integrated option.

615 Each of the token was represented by a vector. The vectorial representation was
616 updated in each layer, and in the following the output of the l^{th} layer was denoted as a
617 superscript, e.g., C_i^l . Following previous works^{32,33}, we calculated a score for each
618 option, which indicated the possibility that the option was the correct answer. The
619 score was calculated by first applying a linear transform to the final representation of
620 the CLS token, i.e.,

621

$$622 \quad s_i = \Phi C_i^{12},$$

623 where C_i^{12} was the final output representation of CLS and Φ was a vector learned

624 from data. The score was independently calculated for each option and then
625 normalized using the following equation:

626

$$627 \quad score_i = \frac{exp(s_i)}{\sum_{i=1}^4 exp(s_i)}.$$

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

630

631 We fine-tuned the DNNs based on the training set of RACE. In the analyses shown in
632 Fig. 5, the fine-tuning process was independently run 10 times. Each time, the training
633 samples were fed in with a randomized order and nodes in the dropout layer were
634 randomly eliminated. Results from the first run of fine-tuning was used for the main
635 analysis reported in Figs. 2-4. All models were implemented based on HuggingFace⁷⁴
636 and all hyperparameters for fine-tuning were adopted from previous studies (Table
637 S1).

638

639 To isolate how the fine-tuning process modulated DNN attention, we also tested the
640 pre-trained DNN that was not fine-tuned on RACE dataset, and compared it with the
641 fine-tuned model (Fig. 4). Furthermore, we quantified how the properties of DNN
642 attention changed throughout the fine-tuning process by analyzing models that
643 received different steps of fine tuning. The steps we sampled were exponentially
644 spaced between 1 and the maximum fine-tuning steps.

645

646 **DNN attention**

647 In each attention head, the attention mechanism calculated an attention weight
 648 between any pair of inputs, including words and special tokens. The vectorial
 649 representation of each input was then updated by the weighted sum of the vectorial
 650 representations of all inputs²². In other words, the models we considered were all
 651 context-dependent models, in which the representation of each word was modeled by
 652 integrating the representations of all inputs. Since only the CLS token was directly
 653 related to question answering, here we analyzed the attention weights that were used
 654 to calculate the vectorial representation of CLS (illustrated in Fig. 1C). For each layer,
 655 the output of an attention head, i.e., C^h , was computed using the following equations.
 656 For the sake of clarity, we denote the input words and tokens generally as X_i .

657

658
$$C^h = \sum_{i=1}^{N+M+2} \alpha_i V_i = \alpha_C V_C + \sum_{n=1}^N \alpha_{Pn} V_{Pn} + \alpha_{S1} V_{S1} + \sum_{m=1}^M \alpha_{Om} V_{Om} + \alpha_{S2} V_{S2},$$

659
$$\alpha_i = \frac{\exp(Q_C K_i^T)}{\sum_{i=1}^{N+M+2} \exp(Q_C K_i^T)},$$

660
$$V_i = X_i W^V + b^V, K_i = X_i W^K + b^K, Q_C = X_C W^Q + b^Q,$$

661

662 where W^v , W^q , W^k , b^v , b^q , and b^k were parameters to learn from the data, and α_i was the
 663 attention weight between CLS and X_i . The attention weight between CLS and the n^{th}
 664 word in the passage, i.e., α_{Pn} , was compared to human attention. Here, we only
 665 considered the attention weights associated with the correct option.

666

667 Output of the attention module, i.e., C^h , was concatenated over all the 12 heads in
668 each layer, and further processed by position-wise operations to generate the final
669 representation of CLS in the layer²². Additionally, DNNs used byte-pair tokenization
670 which split some words into multiple tokens. We converted the token-level attention
671 weights to word-level attention weights by summing the attention weights over tokens
672 within a word^{63,66}.

673

674 **Human attention analysis and prediction**

675 We analyzed eye fixations during passage reading in Study 1 and the first-pass
676 passage reading in Study 2. The total fixation time of each word was extracted using
677 the SR Research Data Viewer software. We averaged the total fixation time across all
678 participants who correctly answered the question, and measured human attention
679 using the attention density, i.e., the total fixation time divided by the area a word
680 occupied.

681

682 We employed linear regression to test whether a set of features could explain human
683 attention distribution. Four sets of features were analyzed, i.e., textual features, layout
684 features, task relevance, and DNN attention weights. The textual features included
685 word length, logarithmic word frequency estimated based on the British National
686 Corpus⁷⁵, ordinal position of a word in a sentence, ordinal position of a word in a
687 passage, and ordinal sentence number of a word. The layout features referred to the
688 visual layout of text, i.e., features induced by line changes, including the coordinate of

689 the left most pixel of a word, ordinal position of a word in a paragraph, ordinal row
690 number of a word in a paragraph, ordinal row number of a word in a passage. Task
691 relevance was annotated by humans (see *Experimental materials* for details), and the
692 DNN attention weights included the 144 attention weights from all layers and
693 attention heads. In the regression analysis, human attention density on word w was
694 modeled using the following equation.

695

$$696 \quad attention_density_w = \sum_{j=1}^J \beta_j F_{w,j} + b + \varepsilon_w,$$

697

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

704

705 **Statistical tests**

706 We employed a one-sided permutation test to test whether a set of features were
707 statistically significant to predict human attention. Five hundred chance-level
708 prediction accuracy was calculated by predicting shuffled human attention.
709 Specifically, the human attention density was shuffled across words and was predicted
710 by word features which were not shuffled. The procedure was repeated 500 times,

711 creating 500 chance-level prediction accuracy. If the actual correlation was greater
712 than N out of the 500 chance-level correlation, the significance level was $(N + 1)/501$.

713

714 The comparison between global and local questions were based on bias-corrected and
715 accelerated bootstrap⁷⁶. For example, to test whether the prediction accuracy differed
716 between the 2 types of questions, all global questions were resampled with
717 replacement 50000 times and each time the prediction accuracy was calculated based
718 on the resampled questions, resulting in 50000 resampled prediction accuracy. If the
719 prediction accuracy for local questions was greater (or smaller) than N out of the
720 50000 resampled accuracy for global questions, the significance level of their
721 difference was $2(N + 1)/50001$. When multiple comparisons were performed, the p-
722 value was further adjusted using the false discovery rate (FDR) correction.

723

724 **Acknowledgements**

725 We thanks David Poeppel and Erik D. Reichle for valuable comments on earlier
726 versions of this manuscript; Jonathan Simon, Bingjiang Lyu, and members of the
727 Ding lab for thoughtful discussions and feedback; Yuran Zhang, Anqi Dai, Zhonghua
728 Tang, and Yuhan Lu for assistance with experiments. Work supported by National
729 Natural Science Foundation of China 31771248 and Major Scientific Research Project
730 of Zhejiang Lab 2019KB0AC02

731

732

733 **Author contributions**

734 Nai Ding acquired the funding, conceived and coordinated the project, analyzed data,
735 and wrote the manuscript. Jiajie Zou implemented the experiments and models,
736 analyzed data, and wrote the manuscript.

737

738 **Competing interests**

739 The authors declare no competing interests.

740

741

742 **References**

- 743 1 Posner, M. I. & Petersen, S. E. The attention system of the human brain. *Annu.*
744 *Rev. Neurosci.* **13**, 25-42 (1990).
- 745 2 Treisman, A. M. & Gelade, G. A feature-integration theory of attention. *Cogn.*
746 *Psychol.* **12**, 97-136 (1980).
- 747 3 Posner, M. I. & Boies, S. J. Components of attention. *Psychol. Rev.* **78**, 391
748 (1971).
- 749 4 Chun, M. M., Golomb, J. D. & Turk-Browne, N. B. A taxonomy of external
750 and internal attention. *Annu. Rev. Psychol.* **62**, 73-101 (2011).
- 751 5 Awh, E., Belopolsky, A. V. & Theeuwes, J. Top-down versus bottom-up
752 attentional control: A failed theoretical dichotomy. *Trends Cogn. Sci.* **16**, 437-
753 443 (2012).
- 754 6 Egeth, H. E. & Yantis, S. Visual attention: Control, representation, and time
755 course. *Annu. Rev. Psychol.* **48**, 269-297 (1997).
- 756 7 Corbetta, M. & Shulman, G. L. Control of goal-directed and stimulus-driven
757 attention in the brain. *Nat. Rev. Neurosci.* **3**, 201-215 (2002).
- 758 8 Buschman, T. J. & Miller, E. K. Top-down versus bottom-up control of
759 attention in the prefrontal and posterior parietal cortices. *Science* **315**, 1860-
760 1862 (2007).
- 761 9 Knudsen, E. I. Fundamental components of attention. *Annu. Rev. Neurosci.* **30**,
762 57-78 (2007).
- 763 10 Atiani, S. *et al.* Emergent selectivity for task-relevant stimuli in higher-order
764 auditory cortex. *Neuron* **82**, 486-499 (2014).
- 765 11 Desimone, R. & Duncan, J. Neural mechanisms of selective visual attention.
766 *Annu. Rev. Neurosci.* **18**, 193-222 (1995).
- 767 12 Broadbent, D. Perception and communication Pergamon. *Elmsford, New York*
768 (1958).
- 769 13 Kahneman, D. *Attention and effort*. Vol. 1063 (Citeseer, 1973).
- 770 14 Dayan, P., Kakade, S. & Montague, P. R. Learning and selective attention. *Nat.*
771 *Neurosci.* **3**, 1218-1223 (2000).

- 772 15 Gottlieb, J., Hayhoe, M., Hikosaka, O. & Rangel, A. Attention, reward, and
773 information seeking. *J. Neurosci.* **34**, 15497-15504 (2014).
- 774 16 Yu, A. J. & Dayan, P. Inference, attention, and decision in a Bayesian neural
775 architecture. In *Proc. Advances in neural information processing systems*
776 1577-1584 (MIT Press, 2005).
- 777 17 Radulescu, A., Niv, Y. & Ballard, I. Holistic reinforcement learning: the role of
778 structure and attention. *Trends Cogn. Sci.* **23**, 278-292 (2019).
- 779 18 Najemnik, J. & Geisler, W. S. Optimal eye movement strategies in visual
780 search. *Nature* **434**, 387-391 (2005).
- 781 19 Navalpakkam, V., Koch, C., Rangel, A. & Perona, P. Optimal reward
782 harvesting in complex perceptual environments. *Proc. Natl Acad. Sci. USA*
783 **107**, 5232-5237 (2010).
- 784 20 Yang, S. C.-H., Lengyel, M. & Wolpert, D. M. Active sensing in the
785 categorization of visual patterns. *Elife* **5**, e12215 (2016).
- 786 21 Bahdanau, D., Cho, K. & Bengio, Y. Neural machine translation by jointly
787 learning to align and translate. In *Proc. 3rd International Conference on*
788 *Learning Representations (ICLR, 2015)*.
- 789 22 Vaswani, A. *et al.* Attention is all you need. In *Proc. Advances in neural*
790 *information processing systems* 5998-6008 (Curran Associates, 2017).
- 791 23 Lai, G., Xie, Q., Liu, H., Yang, Y. & Hovy, E. Race: Large-scale reading
792 comprehension dataset from examinations. In *Proc. 2017 Conference on*
793 *Empirical Methods in Natural Language Processing* 785-794 (Association for
794 Computational Linguistics, 2017).
- 795 24 Yamins, D. L. *et al.* Performance-optimized hierarchical models predict neural
796 responses in higher visual cortex. *Proc. Natl Acad. Sci. USA* **111**, 8619-8624
797 (2014).
- 798 25 Kell, A. J., Yamins, D. L., Shook, E. N., Norman-Haignere, S. V. &
799 McDermott, J. H. A task-optimized neural network replicates human auditory
800 behavior, predicts brain responses, and reveals a cortical processing hierarchy.
801 *Neuron* **98**, 630-644. e616 (2018).
- 802 26 Saxe, A., Nelli, S. & Summerfield, C. If deep learning is the answer, what is
803 the question? *Nat. Rev. Neurosci.* **22**, 55-67 (2021).

- 804 27 White, S., Chen, J. & Forsyth, B. Reading-related literacy activities of
805 American adults: Time spent, task types, and cognitive skills used. *J. Lit. Res.*
806 **42**, 276-307 (2010).
- 807 28 Rayner, K. & McConkie, G. W. What guides a reader's eye movements? *Vision*
808 *Res.* **16**, 829-837 (1976).
- 809 29 Borji, A., Sihite, D. N. & Itti, L. Quantitative analysis of human-model
810 agreement in visual saliency modeling: a comparative study. *IEEE Trans.*
811 *Image Process.* **22**, 55-69 (2013).
- 812 30 Bylinskii, Z., Judd, T., Oliva, A., Torralba, A. & Durand, F. What do different
813 evaluation metrics tell us about saliency models? *IEEE Trans. Pattern Anal.*
814 *Mach. Intell.* **41**, 740-757 (2018).
- 815 31 Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep
816 bidirectional transformers for language understanding. In *Proc. 2019*
817 *Conference of the North American Chapter of the Association for*
818 *Computational Linguistics: Human Language Technologies* 4171-4186
819 (Association for Computational Linguistics, 2019).
- 820 32 Lan, Z. *et al.* Albert: A lite bert for self-supervised learning of language
821 representations. In *Proc. International Conference on Learning*
822 *Representations* (ICLR, 2020).
- 823 33 Liu, Y. *et al.* Roberta: A robustly optimized bert pretraining approach. Preprint
824 at <https://arxiv.org/abs/1907.11692> (2019).
- 825 34 Zhou, B., Khosla, A., Lapedriza, A., Oliva, A. & Torralba, A. Object detectors
826 emerge in deep scene CNNs. In *Proc. International Conference on Learning*
827 *Representations* (ICLR, 2015).
- 828 35 Pelli, D. G. *et al.* Crowding and eccentricity determine reading rate. *J. Vis.* **7**,
829 20-20 (2007).
- 830 36 Gottlieb, J., Cohanpour, M., Li, Y., Singletary, N. & Zabeh, E. Curiosity,
831 information demand and attentional priority. *Curr. Opin. Behav.* **35**, 83-91
832 (2020).
- 833 37 Marr, D. Vision: A computational investigation into the human representation
834 and processing of visual information, henry holt and co. *Inc.*, *New York, NY* **2**
835 (1982).

- 836 38 Krajbich, I., Armel, C. & Rangel, A. Visual fixations and the computation and
837 comparison of value in simple choice. *Nat. Neurosci.* **13**, 1292-1298 (2010).
- 838 39 Jang, A. I., Sharma, R. & Drugowitsch, J. Optimal policy for attention-
839 modulated decisions explains human fixation behavior. *Elife* **10**, e63436
840 (2021).
- 841 40 Wyart, V., Myers, N. E. & Summerfield, C. Neural mechanisms of human
842 perceptual choice under focused and divided attention. *J. Neurosci.* **35**, 3485-
843 3498 (2015).
- 844 41 Radford, A., Narasimhan, K., Salimans, T. & Sutskever, I. Improving language
845 understanding by generative pre-training. (2018).
- 846 42 Hale, J. Information - theoretical complexity metrics. *Lang. Linguist.*
847 *Compass.* **10**, 397-412 (2016).
- 848 43 Smith, N. J. & Levy, R. The effect of word predictability on reading time is
849 logarithmic. *Cognition* **128**, 302-319, doi:10.1016/j.cognition.2013.02.013
850 (2013).
- 851 44 Itti, L. & Baldi, P. Bayesian surprise attracts human attention. *Vision Res.* **49**,
852 1295-1306 (2009).
- 853 45 Tatler, B. W., Hayhoe, M. M., Land, M. F. & Ballard, D. H. Eye guidance in
854 natural vision: reinterpreting salience. *J. Vis.* **11**, 5-25, doi:10.1167/11.5.5
855 (2011).
- 856 46 Anderson, P. *et al.* Bottom-up and top-down attention for image captioning
857 and visual question answering. In *Proc. The IEEE Conference on Computer*
858 *Vision and Pattern Recognition* 6077-6086 (CVPR, 2018).
- 859 47 Xu, K. *et al.* Show, attend and tell: Neural image caption generation with
860 visual attention. In *Proc. 32nd International Conference on Machine Learning*
861 2048-2057 (PMLR, 2015).
- 862 48 Das, A., Agrawal, H., Zitnick, L., Parikh, D. & Batra, D. Human attention in
863 visual question answering: Do humans and deep networks look at the same
864 regions? *Comput. Vis. Image Underst.* **163**, 90-100 (2017).
- 865 49 Liu, C., Mao, J., Sha, F. & Yuille, A. Attention correctness in neural image
866 captioning. In *Proc. AAAI Conference on Artificial Intelligence* 4176-4182
867 (AAAI, 2017).

- 868 50 Wolfe, J. M. & Horowitz, T. S. Five factors that guide attention in visual
869 search. *Nat. Hum. Behav.* **1**, 1-8 (2017).
- 870 51 Just, M. A. & Carpenter, P. A. A theory of reading: From eye fixations to
871 comprehension. *Psychol. Rev.* **87**, 329-354 (1980).
- 872 52 Reichle, E. D., Pollatsek, A., Fisher, D. L. & Rayner, K. Toward a model of
873 eye movement control in reading. *Psychol. Rev.* **105**, 125-157 (1998).
- 874 53 Reichle, E. D., Rayner, K. & Pollatsek, A. The EZ Reader model of eye-
875 movement control in reading: Comparisons to other models. *Behav. Brain Sci.*
876 **26**, 445-476 (2003).
- 877 54 Engbert, R., Nuthmann, A., Richter, E. M. & Kliegl, R. SWIFT: a dynamical
878 model of saccade generation during reading. *Psychol. Rev.* **112**, 777-813
879 (2005).
- 880 55 Legge, G. E., Hooven, T. A., Klitz, T. S., Mansfield, J. S. & Tjan, B. S. Mr.
881 Chips 2002: New insights from an ideal-observer model of reading. *Vision*
882 *Res.* **42**, 2219-2234 (2002).
- 883 56 Liu, Y. & Reichle, E. The emergence of adaptive eye movements in reading. In
884 *Proc. Annual Meeting of the Cognitive Science Society*.
- 885 57 Mancheva, L. *et al.* An analysis of reading skill development using EZ Reader.
886 *J. Cogn. Psychol.* **27**, 657-676 (2015).
- 887 58 Clifton, C. *et al.* Eye movements in reading and information processing: Keith
888 Rayner's 40 year legacy. *J. Mem. Lang.* **86**, 1-19,
889 doi:10.1016/j.jml.2015.07.004 (2016).
- 890 59 Duffy, T. M. & Kabanec, P. Testing a readable writing approach to text
891 revision. *J. Educ. Psychol.* **74**, 733-748 (1982).
- 892 60 van der Schoot, M., Vasbinder, A. L., Horsley, T. M. & van Lieshout, E. C. D.
893 M. The role of two reading strategies in text comprehension: An eye fixation
894 study in primary school children. *J. Res. Read.* **31**, 203-223 (2008).
- 895 61 Kaakinen, J. K., Hyönä, J. & Keenan, J. M. How prior knowledge, WMC, and
896 relevance of information affect eye fixations in expository text. *J. Exp.*
897 *Psychol.: Learn. Mem. Cogn.* **29**, 447-457 (2003).
- 898 62 Rayner, K. & Fischer, M. H. Mindless reading revisited: Eye movements
899 during reading and scanning are different. *Perception & Psychophysics* **58**,

- 900 734-747 (1996).
- 901 63 Clark, K., Khandelwal, U., Levy, O. & Manning, C. D. What Does BERT
902 Look at? An Analysis of BERT’s Attention. In *Proc. 2019 ACL Workshop*
903 *BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP* 276-286
904 (Association for Computational Linguistics, 2019).
- 905 64 Voita, E., Talbot, D., Moiseev, F., Sennrich, R. & Titov, I. Analyzing multi-
906 head self-attention: Specialized heads do the heavy lifting, the rest can be
907 pruned. In *Proc. 57th Annual Meeting of the Association for Computational*
908 *Linguistics* 5797-5808 (Association for Computational Linguistics, 2019).
- 909 65 Sood, E., Tannert, S., Frassinelli, D., Bulling, A. & Vu, N. T. Interpreting
910 attention models with human visual attention in machine reading
911 comprehension. In *Proc. 24th Conference on Computational Natural*
912 *Language Learning* 12–25 (Association for Computational Linguistics).
- 913 66 Bolotova, V. *et al.* Do People and Neural Nets Pay Attention to the Same
914 Words: Studying Eye-tracking Data for Non-factoid QA Evaluation. In *Proc.*
915 *29th ACM International Conference on Information & Knowledge*
916 *Management* 85-94 (ACM, 2020).
- 917 67 Sen, C., Hartvigsen, T., Yin, B., Kong, X. & Rundensteiner, E. Human
918 Attention Maps for Text Classification: Do Humans and Neural Networks
919 Focus on the Same Words? In *Proc. 58th Annual Meeting of the Association*
920 *for Computational Linguistics* 4596-4608.
- 921 68 Yang, Z. *et al.* Hierarchical attention networks for document classification. In
922 *Proc. 2016 Conference of the North American Chapter of the Association for*
923 *Computational Linguistics: Human Language Technologies* 1480-1489
924 (Association for Computational Linguistics, 2016).
- 925 69 Lin, Z. *et al.* A structured self-attentive sentence embedding. In *Proc.*
926 *International Conference on Learning Representations* (ICLR, 2017).
- 927 70 Serrano, S. & Smith, N. A. Is Attention Interpretable? In *Proc. 57th Annual*
928 *Meeting of the Association for Computational Linguistics* 2931-2951
929 (Association for Computational Linguistics, 2019).
- 930 71 Jain, S. & Wallace, B. C. Attention is not Explanation. In *Proc. 2019*
931 *Conference of the North American Chapter of the Association for*
932 *Computational Linguistics* 3543-3556 (Association for Computational

933 Linguistics).

934 72 Zou, J. *et al.* PALRACE: Reading Comprehension Dataset with Human Data
935 and Labeled Rationales. Preprint at <https://arxiv.org/abs/2106.12373> (2021).

936 73 Brainard, D. H. The psychophysics toolbox. *Spat. Vis.* **10**, 433-436 (1997).

937 74 Wolf, T. *et al.* HuggingFace's Transformers: State-of-the-art natural language
938 processing. In *Proc. 2020 Conference on Empirical Methods in Natural*
939 *Language Processing: System Demonstrations* 38-45 (Association for
940 Computational Linguistics, 2020).

941 75 Burnard, L. *The British National Corpus, version 3 (BNC XML Edition)*.
942 [<http://www.natcorp.ox.ac.uk/>](http://www.natcorp.ox.ac.uk/) (2007).

943 76 Efron, B. & Tibshirani, R. J. *An introduction to the bootstrap*. (CRC press,
944 1994).

945

Supplementary Files

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

- [SM1030.pdf](#)