# A compositional letter code in high-level visual cortex

# explains how we read jumbled words

Aakash Agrawal[1], K.V.S. Hari[2] & S. P. Arun[3*]

[1]Centre for BioSystems Science & Engineering, [2]Department of Electrical

Communication Engineering & [3]Centre for Neuroscience

Indian Institute of Science, Bangalore, 560012, India


*Correspondence to: S. P. Arun (sparun@iisc.ac.in)

## ABSTRACT

We read words and even jubmled wrods effortlessly, but the neural representations underlying this remarkable ability remain unknown. We hypothesized that word processing is driven by a visual representation that is compositional i.e. with string responses systematically related to letters. To test this hypothesis, we devised a model in which neurons tuned to letter shape respond to longer strings by linearly summing letter responses. This letter model explained human performance in both visual search as well as word reading tasks. Brain imaging revealed that viewing a string activates this compositional letter code in the lateral occipital (LO) region, and that subsequent comparisons to known words are computed by the visual word form area (VWFA). Thus, seeing a word activates a compositional letter code that enables efficient reading.

**INTRODUCTION**

14

15       Reading is a recent cultural invention, yet we are remarkably efficient at reading

16   words and even jmulbed wrods (Figure 1A). What makes a jumbled word easy or hard

17   to read? This question has captured the popular imagination through demonstrations

18   such as the Cambridge University effect (Rawlinson, 1976; Grainger and Whitney,

19   2004), depicted in Figure 1A. Reading a word or a jumbled word can be influenced by

20   a variety of factors (Norris, 2013; Grainger, 2018). Word reading is easy when similar

21   shapes are substituted (Perea et al., 2008; Perea and Panadero, 2014), when the first

22   and last letters are preserved (Rayner et al., 2006), when there are fewer

23   transpositions (Gomez et al., 2008) and when word shape is preserved (Norris, 2013;

24   Grainger, 2018). Word reading is also easier for words with frequent bigrams or

25   trigrams, for frequent words and for shuffled words that preserve intermediate units

26   such as consonant clusters or morphemes (Norris, 2013; Grainger, 2018). Despite

27   these insights, it is not clear how these factors combine, what their distinct

28   contributions are, and more generally, how word representations relate to letter

29   representations.

30       Here, we hypothesized that word reading is enabled by a purely visual

31   representation. To probe purely visual processing, we devised a visual search task in

32   which subjects had to find an oddball target among distractors. This task does not

33   require any explicit reading and is driven by shape representations in visual cortex

34   (Sripati and Olson, 2010a; Zhivago and Arun, 2014). An example visual search array

35   containing two oddball targets is shown in Figure 1B. It can be seen that finding

36   OFRGET is easy among FORGET whereas finding FOGRET is hard (Figure 1B). This

37   difference in visual similarity (Figure 1C) explains why transposing the middle letters

38   renders a word easier to read than transposing its edge letters. This example suggests

39  that word reading could be explained by purely visual processing as indexed by visual

40  search. However, subjects may have been reading during visual search, thereby

41  activating non-visual lexical or linguistic factors.

42      To overcome this confound, we asked whether visual search involving letter

43  strings can be explained using a neurally plausible model containing only visual

44  factors. We drew upon two well-established principles of object representations in

45  high-level visual cortex. First, images that are perceptually similar elicit similar activity

46  in single neurons (Op de Beeck et al., 2001; Sripati and Olson, 2010a; Zhivago and

47  Arun, 2014). Accordingly we used visual search for single letters to create artificial

48  neurons tuned for single letters. Second, the neural response to multiple objects is an

49  average of the response to the individual objects, a phenomenon known as divisive

50  normalization (Zoccolan et al., 2005; Ghose and Maunsell, 2008; Zhivago and Arun,

51  2014). Accordingly, we created neural responses to letter strings as a linear sum of

52  single letter responses. In contrast to an influential proposal that requires neurons

53  tuned to letter combinations (Dehaene et al., 2005, 2010), our model only assumes

54  neurons tuned for letter shape and retinal position, as observed in high-level visual

55  cortex (Lehky and Tanaka, 2016). It does not capture any information about

56  specialized detectors for longer strings, or about other lexical or linguistic factors. We

57  used this model to explain human performance on visual search as well as word

58  recognition tasks. Finally, using brain imaging, we identified the neural substrates for

59  both the letter code as well as subsequent lexical decisions.
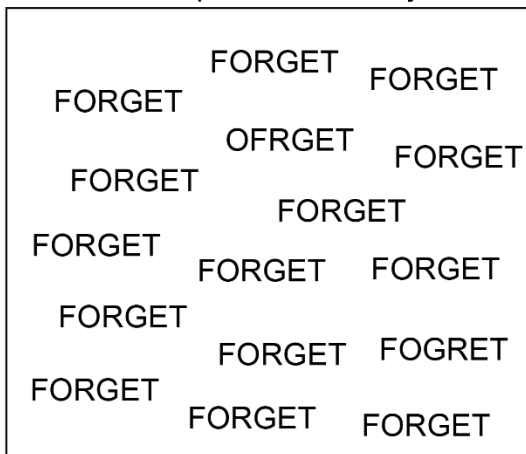
60

61

62

A

AOCCDRNIG TO A RSEEARCH AT CMABRIGDE UINERVTISY, IT DEOSN'T MTTAER IN WAHT OREDR THE LTTEERS IN A WROD ARE, THE OLNY IPRMOETNT TIHNG IS TAHT THE FRIST AND LSAT LTTEER BE AT THE RGHIT PCLAE. THE RSET CAN BE A TOATL MSES AND YOU CAN SITLL RAED IT WOUTHIT A PORBELM. TIHS IS BCUSEAE THE HUAMN MNID DEOS NOT RAED ERVEY LTETER BY ISTLEF, BUT THE WROD AS A WLOHE.
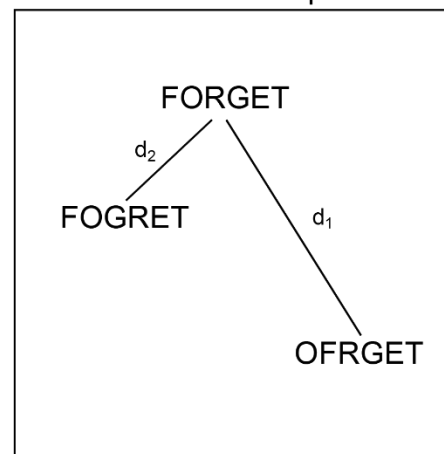


**Figure 1. Reading jumbled words**

(A) We are extremely good at reading jumbled words, as illustrated by the popular Cambridge University effect.

(B) Visual search array showing two oddball targets (OFRGET & FOGRET) among many instances of FORGET. OFRGET is easy to find but not FOGRET.

(C) Schematic representation of these strings in visual search space, arranged such that similar items (corresponding to harder searches) are nearby. Thus, FOGRET is closer to FORGET compared to OFRGET (i.e. $d_1 > d_2$).

**RESULTS**

72

73    We performed five key experiments. In Experiment 1, subjects performed visual

74    search involving single letters, and we used this to construct artificial neurons tuned

75    for letter shape. In Experiments 2-4, we show that search for longer strings can be

76    predicted using these artificial neurons with a simple compositional rule. In Experiment

77    5, we show that this model also explains human performance on a commonly studied

78    word recognition task. Finally, in Experiment 6, we measured brain activations during

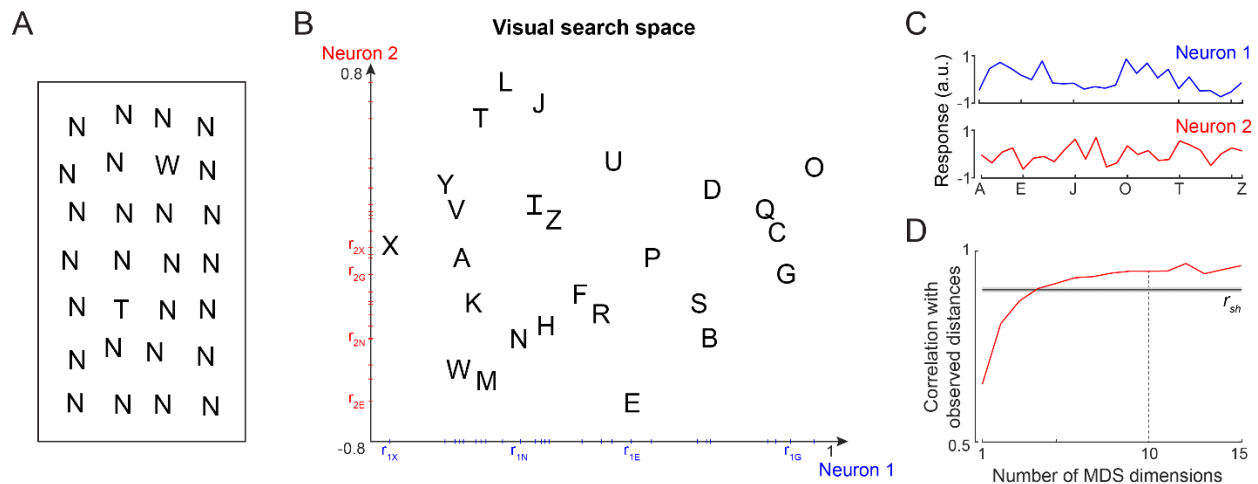79    word recognition to elucidate the underlying neural representations.

80

81    **Experiment 1: Single letter searches**

82    We recruited 16 subjects to perform an oddball visual search task involving

83    uppercase letters (n = 26), lowercase letters (n = 26) and digits (n = 10). An example

84    search is shown in Figure 2A. Subjects were highly consistent in their responses (split-

85    half correlation between average search times of odd- and even-numbered subjects:

86    r = 0.87, p < 0.00005). We calculated the reciprocal of search times for each letter pair

87    which is a measure of distance between them (Arun, 2012). These letter dissimilarities

88    were significantly correlated with previously reported subjective dissimilarity ratings

89    (Section S1).

90    Since shape dissimilarity in visual search matches closely with neural

91    dissimilarity in visual cortex (Sripati and Olson, 2010a; Zhivago and Arun, 2014), we

92    asked whether these letter distances can be used to reconstruct the underlying neural

93    responses to single letters. To do so, we performed a multidimensional scaling (MDS)

94    analysis, which finds the n-dimensional coordinates of all letters such that their

95    distances match the observed visual search distances. In the resulting plot for 2

96    dimensions for uppercase letters (Figure 2B), nearby letters correspond to small

97   distances i.e. long search times. The coordinates of letters along a particular

98   dimension can then be taken as the putative response of a single neuron. For example,

99   the first dimension represents the activity of a neuron that responds strongest to the

100  letter O and weakest to X (Figure 2C). Likewise the second dimension corresponds to

101  a neuron that responds strongest to L and weakest to E (Figure 2C). We note that the

102  same set of distances can be obtained from a different set of neural responses: a

103  simple coordinate axis rotation would result in another set of neural responses with an

104  equivalent match to the observed distances. Thus, the estimated activity from MDS

105  represents one possible solution to how neurons should respond to individual letters

106  so as to collectively produce behaviour.

107       As expected, increasing the number of MDS dimensions led to increased match

108  to the observed letter dissimilarities (Figure 2D). Taking 10 MDS dimensions, which

109  explain nearly 95% of the variance, we obtained the single letter responses of 10 such

110  artificial neurons. We used these single letter responses to predict their response to

111  longer letter strings in all the experiments. Varying this choice yielded qualitatively

112  similar results. Analogous results for all letters and numbers are shown in Section S1.

113

**Figure 2. Single letter discrimination (Experiment 1)**

(A) Visual search array showing two oddball targets (W & T) among many Ns. It can be seen that finding W is harder compared to finding T. The actual experiment comprised search arrays with only one oddball target among 15 distractors.

(B) Visual search space for uppercase letters obtained by multidimensional scaling of observed dissimilarities. Nearby letters represent hard searches. Distances in this 2D plot are highly correlated with the observed distances (r = 0.82, p < 0.00005). Letter activations along the x-axis are taken as responses of Neuron 1 (*blue*), and along the y-axis are taken as Neuron 2 (*red*), etc. The tick marks indicate the response of each letter along that neuron.

(C) Responses of Neuron 1 and Neuron 2 shown separately for each letter. Neuron 1 responds best to O, whereas Neuron 2 responds best to L.

(D) Correlation between observed distances and MDS embedding as a function of number of MDS dimensions. The *black* line represents the split-half correlation with error bars representing s.d calculated across 100 random splits.

**131  Experiment 2: Bigram searches**

132      Next we proceeded to ask whether searches for longer strings can be explained

133  using single letter responses. In Experiment 2, we asked subjects to perform oddball

134  searches involving bigrams. An example search is depicted in Figure 3A. It can be

135  seen that, finding TA among AT is harder than finding UT among AT. Thus, letter

136  transpositions are more similar compared to letter substitutions, consistent with the

137  classic results on reading (Norris, 2013; Grainger, 2018). To characterize the effect of

138  bigram frequency, we included both frequent bigrams (e.g. IN, TH) and infrequent

139  bigrams (e.g. MH, HH). As before, subjects were highly consistent in their performance

140  (split-half correlation between odd and even numbered subjects across all bigrams: r

141  = 0.82, p < 0.00005).

142      Next we asked whether bigram search performance can be explained using

143  neurons tuned to single letters estimated from Experiment 1. The essential principle

144  for constructing bigram responses is depicted in Figure 3B. In monkey visual cortex,

145  the response of single neurons to two simultaneously presented objects is an average

146  of the single object responses (Zoccolan et al., 2005; Zhivago and Arun, 2014; Pramod

147  and Arun, 2018). This averaging can easily be biased through changes in divisive

148  normalization (Ghose and Maunsell, 2008). Therefore we took the response of each

149  neuron to a bigram to be a weighted sum of its responses to the constituent letters

150  (Figure 3B). Specifically, the response of a neuron to the bigram AB is given by $r_{AB} =$

151  $w_1 r_A + w_2 r_B$, where $r_{AB}$ is the response to AB, $r_A$ and $r_B$ are its responses to the

152  constituent letters A & B, and $w_1$, $w_2$ are the summation weights reflecting the

153  importance of letters A & B in the summation. Note that if $w_1 = w_2$, the bigram response

154  to AB and BA will be identical. Thus, discriminating letter transpositions necessarily

155  requires asymmetric summation in at least one of the neurons.

156      To summarize, the letter model for bigrams has two unknown spatial weighting

157      parameters for each of the 10 neurons, resulting in 2 x 10 = 20 free parameters. To

158      calculate dissimilarities between a pair of bigrams, we calculated the Euclidean

159      distance between the 10-dimensional response vectors corresponding to the two

160      bigrams. The data collected in the experiment comprised dissimilarities (1/RT) from

161      1,176 searches involving all possible pairs of 49 bigrams. To estimate the model

162      parameters, we optimized them to match the observed bigram dissimilarities using

163      standard nonlinear fitting algorithms (see Methods).

164      This letter model yielded excellent fits to the observed data ($r = 0.85$, $p <$

165      $0.00005$; Figure 3C). To assess whether the model explains all the systematic

166      variance in the data, we calculated an upper bound estimated from the inter-subject

167      consistency (see Methods). This consistency measure ($r_{data} = 0.90$) was close to the

168      model fit, suggesting that the model captured nearly all the systematic variance in the

169      data. As predicted in the schematic figure (Figure 3B), the estimated spatial

170      summation weights were unequal (absolute difference between $w_1$ and $w_2$, mean ± sd:

171      $0.07 ± 0.04$). To assess whether this difference is statistically significant, we randomly

172      shuffled the observed dissimilarities and estimated these weights. The absolute

173      difference between shuffled weights was significantly smaller than for the original

174      weights (average absolute difference: $0.03 ± 0.02$; $p < 0.005$, sign-rank test across 10

175      neurons).

176      According to an influential account of word reading, specialized detectors are

177      formed for frequently occurring combinations of letters (Dehaene et al., 2005). If this

178      were the case, searches involving frequent bigrams (e.g. TH, ND) or two letter words

179      (e.g. AN, AM) should produce larger model errors compared to infrequent bigrams,

180      since our model does not incorporate any bigram-selective units. Alternatively, if

181  bigram discrimination was driven entirely by single letters, we should find no difference

182  in errors.   In keeping with this latter prediction, we observed no visually obvious

183  difference in model fits for frequent bigram pairs or word-word pairs compared to other

184  bigram pairs (Figure 3C). To quantify this observation, we asked whether the model

185  error for each bigram pair, calculated as the absolute difference between observed

186  and predicted dissimilarity, covaried with the average bigram frequency of the two

187  bigrams (for both frequent bigrams and words). This revealed a weak negative

188  correlation whereby frequent bigram pairs showed smaller errors ($r = -0.06$, $p = 0.04$

189  across 1176 bigram pairs). This is the opposite of what would be expected if there

190  were specialized detectors. To further investigate possible bigram frequency effects,

191  we compared the model error for the 20 bigram pairs with the largest mean bigram

192  frequency with the 20 pairs with the lowest mean bigram frequency. This too revealed

193  no systematic difference (mean ± sd of residual error: $0.10 \pm 0.08$ for the 20 most

194  frequent bigrams and words; $0.11 \pm 0.09$ for 20 least frequent bigrams; $p = 0.80$, rank-

195  sum test). Thus, model errors are not systematically different for frequent compared

196  to infrequent bigram pairs. We conclude that bigram search can be explained entirely

197  using single neurons tuned to single letters.

198

199  **Experiment 3: Upright versus inverted bigrams**

200       In the letter model described above, the response to bigrams is a weighted sum

201  of the single letter responses. As detailed earlier, a critical prediction of this model is

202  that the response to transposed bigrams such as AB & BA will be different only if the

203  summation weights are unequal. By contrast, repeated letter bigrams such as AA &

204  BB will remain discriminable regardless of the nature of summation, since their

205  response will be proportional to the respective single letter responses. Since reading

206     expertise can modulate sensitivity to letter transpositions, we reasoned that familiarity

207     might modulate the summation to make it more asymmetric. We therefore predicted

208     that this would make transposed letter searches (with AB as target and BA as

209     distractor, or vice-versa) easier to discriminate in a familiar upright orientation

210     compared to the (unfamiliar) inverted orientation. By contrast, searches involving

211     repeated letter bigrams (with AA as target and BB as distractor), which also have a

212     change in two letters, will remain equally easy in both upright and inverted orientations.

213       We tested this prediction in Experiment 3 by asking subjects to perform

214     searches involving upright and inverted bigrams. The essential findings are

215     summarized in Figure 3D. As predicted, subjects discriminated repeated letter bigrams

216     (AA-BB searches) equally well at both upright and inverted orientations, but were

217     substantially faster at discriminating transposed letter pairs (AB-BA searches) in the

218     upright orientation (Figure 3D; for detailed analyses see Section S2). We obtained

219     similar results on comparing upright and inverted trigrams as well (Section S2).
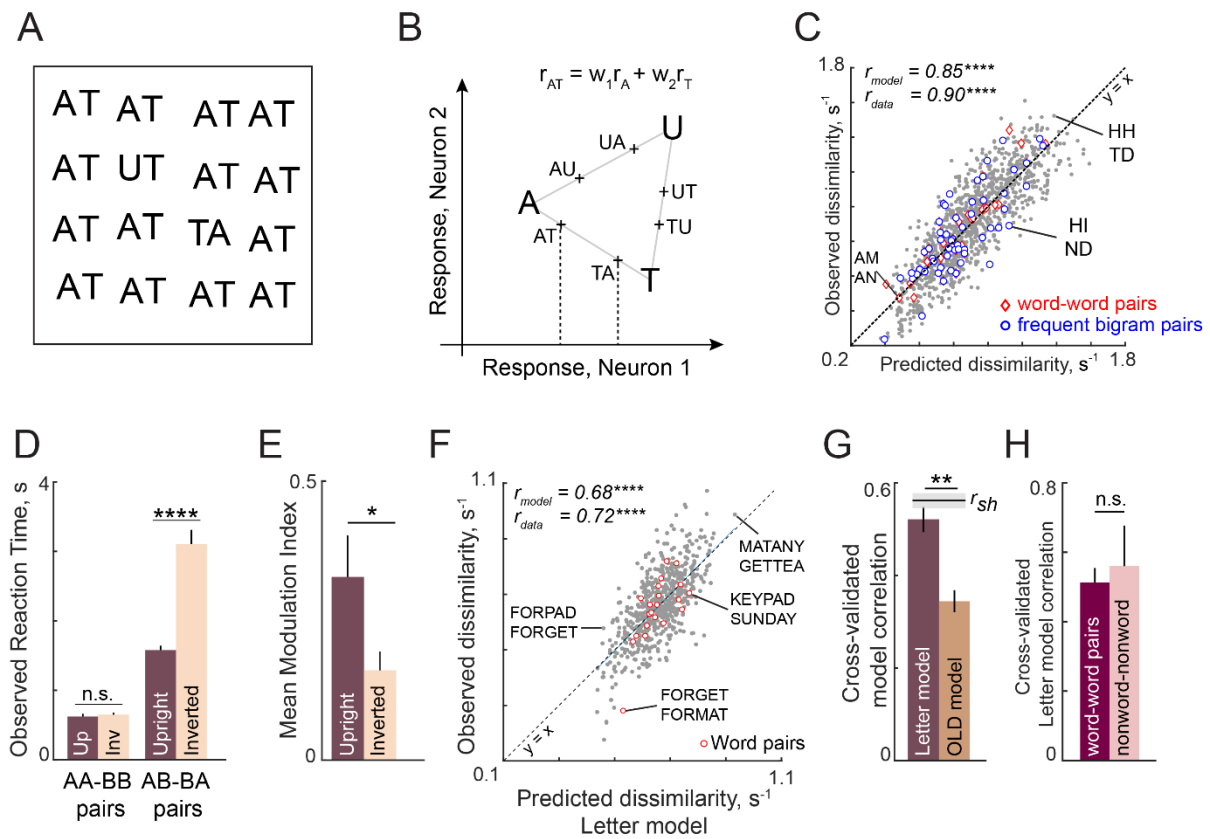
220       We conclude that familiarity leads to asymmetric spatial summation.

221

222

223

224

**Figure 3. Discrimination of strings is explained using single letters (Expts 2-4)**

(A) Example search array with two oddball targets (UT & TA) among the bigram AT. It can be seen that UT is easier to find than TA, showing that letter substitution causes a bigger visual change compared to transposition.

(B) Schematic diagram of how the bigram response is obtained from letter responses. Consider two neurons selective to single letters A, T & U. These letters can be represented in a 2D space in which the response to each neuron lies along one axis. For each neuron, we take the response to a bigram to be a weighted sum of the single letter responses. Thus, the bigram response lies along the line joining the two stimuli. Note that the bigrams AT and TA can be distinguished only if there is unequal summation. In the schematic, the first position is taken to have higher magnitude, as a result of which the response to AT is closer to A than to T.

(C) Observed dissimilarities between bigram pairs plotted against predictions of the letter model for word-word pairs (*red diamonds*), frequent bigram pairs (*blue circles*) and all other bigram pairs (*gray dots*), for Experiment 2. Model correlation is shown at the top left, along with the data consistency for comparison. Asterisks indicate the statistical significance of the correlations (**** is $p < 0.00005$).

(D) Average observed search reaction time for upright (dark) and inverted (pale) bigram searches for repeated letter pairs (AA-BB pairs) and transposed letter pairs (AB-BA pairs) in Experiment 3. Asterisks indicate statistical significance of the main effect of orientation in an ANOVA (see text for details; **** is $p < 0.00005$).

(E) Mean modulation index of the summation weights, calculated as $|w1-w2|/|w1+w2|$, where w1 and w2 are the bigram summation weights, averaged across the 10 neurons in the letter model for upright (dark) and inverted (pale) bigrams. The asterisk indicates statistical significance calculated on a sign-rank test comparing the modulation index across 10 neurons (* is $p < 0.05$).

253  (F) Observed dissimilarities between 6-letter strings in visual search (Experiment 4)
254      plotted against predicted dissimilarities from the single letter model for word-word
255      pairs (*red dots*) and all other pairs (*gray dots*). Model correlation is shown at the
256      top left with data consistency for comparison. Asterisks indicate statistical
257      significance of the correlations (**** is p < 0.00005).
258  (G) Cross-validated model correlation for the letter model (*dark*) and the Orthographic
259      Levenshtein distance (OLD) model (*light*). For each model, the cross-validated
260      correlation is the correlation between model predictions trained on one half of the
261      data and the observed response times from the other half. The upper bound on
262      model fits is the split-half correlation ($r_{sh}$) shown in black with shaded error bars
263      representing standard deviation across 1000 random splits. The asterisk indicates
264      statistical significance of the comparison obtained by estimating the fraction of
265      bootstrap samples in which the observed difference was violated (** is p < 0.005).
266  (H) Cross-validated letter model correlation for word-word pairs and nonword-nonword
267      pairs.
268

## Generalization to longer strings

270      To investigate whether these results would generalize to longer strings which

271  can contain frequent words, we performed several additional visual search

272  experiments using 3, 4, 5 and 6-letter uppercase strings (Section S4). In Experiment

273  4, subjects performed visual search involving six-letter strings that were either valid

274  compound words (e.g. FORGET, TEAPOT) or pseudowords (FORPOT, TEAGET).

275  The single letter model yielded excellent fits to the data (Figure 3F). These fits were

276  superior to a widely used measure of string similarity, the Orthographic Levenshtein

277  Distance (OLD) model (Figure 3G). Importantly, the letter model fits were equivalent

278  for both word-word pairs and nonword-nonword pairs (Figure 3H). These and other

279  analyses are described in Section S3.

280      The letter model also yielded excellent fits across all string lengths tested. We

281  also tested lowercase and mixed-case strings because word shape is thought to play

282  a role when letters vary in size or have upward and downward deflections (Pelli and

283  Tillman, 2007). Even here, the letter model, without any explicit representation of

284  overall word shape, was able to accurately predict most of the search performance.

285  These results are detailed in Section S4.

286    The letter model described is neurally plausible and compositional, but is based

287    on dissimilarities between letters presented in isolation. It could be that the

288    representation of a letter within a bigram, although compositional, differs from its

289    representation when seen in isolation. To explore these possibilities we developed an

290    alternate model in which bigram dissimilarities can be predicted using a sum of

291    (unknown) part dissimilarities at different locations. The resulting model, which we

292    denote as the part sum model, yielded comparable fits to the data. It is completely

293    equivalent to the letter model under certain conditions. Unlike the letter model which

294    is nonlinear and could suffer from multiple local minima, the part sum model is linear

295    and its parameters can be estimated uniquely using standard linear regression. Its

296    complexity can be drastically reduced using simplifying assumptions without affecting

297    model fits. These results are detailed in Section S5.

298

299    **Experiment 5: Lexical decision task**

300    The above experiments show that discrimination of strings in visual search can

301    be explained by neurons tuned for single letter shape with letter responses that

302    combine linearly. Could the same shape representation drive reading behaviour? We

303    evaluated this possibility through two separate word recognition experiments.

304    In Experiment 5, we used a widely used paradigm for word recognition, a lexical

305    decision task (Norris, 2013; Grainger, 2018). Subjects had to indicate whether a string

306    of letters is a word or not using a keypress. The words comprised 4, 5 or 6-letter words

307    and the nonwords consisted of random strings and jumbled words. Subjects were

308    highly accurate in responding to both words and nonwords (mean ± sd: 96 ± 2% for

309    words, 95 ± 3% for nonwords). Importantly, their response times across words and

310    nonwords were consistent as evidenced by a significant split-half correlation

311   (correlation between odd- and even-numbered subjects: r = 0.59 for words, r = 0.73

312   for nonwords, p < 0.00005). Since responses in lexical decision tasks are thought to

313   depend on accumulation of evidence towards or against word status (Ratcliff et al.,

314   2004; Ratcliff and McKoon, 2008), we hypothesized that looking at a string of letters

315   will activate the compositional neural code for the string which is then compared to

316   stored patterns corresponding to known words.

317       We started by characterizing response times for words. To depict the

318   systematic variation in word response times, we plotted them in descending order

319   (Figure 4A). Subjects took longer to respond to infrequent words like MALICE

320   compared to frequent words like MUSIC. If the string is a word, the response time will

321   depend on the strength of the stored pattern, which in turn would depend on lexical

322   factors such as word frequency (Ratcliff et al., 2004; Ratcliff and McKoon, 2008).

323   Indeed, response times for words showed a negative correlation with log word

324   frequency (r = -0.5, p < 0.00005 across 450 words). We also estimated other lexical

325   factors such as the logarithm of the letter frequency (averaged across letters of the

326   string), logarithm of the bigram frequency (averaged across all bigrams in the string),

327   and the number of orthographic neighbours (i.e. number of nearby words), which are

328   all standard measures in linguistic corpora (see Methods).

329       To avoid overfitting, we trained a model based on each factor on one half of the

330   subjects and tested it on the other half. This cross-validated performance is shown for

331   all lexical factors in Figure 4B. It can be seen that the word frequency is the best

332   predictor of word response times (Figure 4B). To assess whether all lexical factors

333   together predict word response times any better, we fit a combined model in which the

334   word response times are modelled as a linear sum of the four factors. The combined

335   model performance was comparable to the performance of the word frequency model

336   alone (Figure 4B). To assess the statistical significance of these results, we performed

337   a bootstrap analysis. On each trial, we trained all models on the dissimilarity obtained

338   from considering only one randomly chosen half of subjects. We calculated the

339   correlation between each model's predictions on the other half of the data, and

340   repeated this procedure 1000 times. Across these samples, the word frequency model

341   performance rarely fell below all other individual models ($p < 0.005$). We conclude that

342   word response times are determined primarily by word frequency.

343       Next we investigated the factors determining the nonword response times. The

344   nonword responses are plotted in descending order in Figure 4C. Subjects took longer

345   to respond to jumbled words like PENICL (original word: PENCIL) with fewer

346   transpositions compared to VTAOCE (original word: OCTAVE) with more

347   transpositions. We hypothesized that, if a string is a nonword, the response will be

348   slow if there is a nearby stored pattern corresponding to a word, and fast otherwise

349   (Dufau et al., 2012; Yap et al., 2015). Likewise the response is likely to be faster if the

350   nearest word is highly familiar (i.e. frequent in the lexicon). Specifically, nonword

351   response times will be inversely proportional to the dissimilarity of the nonword to the

352   nearest word (Figure 4D), and also inversely proportional to the frequency of the

353   nearest word (Figure 4D).

354       To test this prediction, we took the letter model with 10 neurons with single letter

355   tuning and optimized the spatial summation weights to match the reciprocal of the

356   nonword responses for each word length. The model yielded excellent fits to the data

357   ($r = 0.70$, $p < 0.00005$; Figure 4E). This model fit was comparable to the data

358   consistency ($r_{data} = 0.84$). Importantly, this model was able to explain classic

359   phenomena in orthographic processing. Specifically, subjects took longer to respond

360   to nonwords obtained by transposing a letter of a word, compared to nonwords

361  obtained through letter substitution – these trends were present in the model

362  predictions as well (Figure 4F). Likewise, subjects took longer when the middle letters

363  were transposed compared to when the edge letters were transposed – as did the

364  model predictions (Figure 4F). These effects replicate the classic orthographic

365  processing effects reported across many studies (Grainger et al., 2012; Norris, 2013;

366  Ziegler et al., 2013; Grainger, 2018).

367       Next we asked whether a widely used measure of orthographic distance could

368  explain the same data. We selected the Orthographic Levenshtein Distance (OLD), in

369  which the net distance between two strings is calculated as the minimum number of

370  letter additions, transpositions and deletions required to transform one string into

371  another. The OLD model yielded relatively poorer predictions of the data (r = 0.36, p
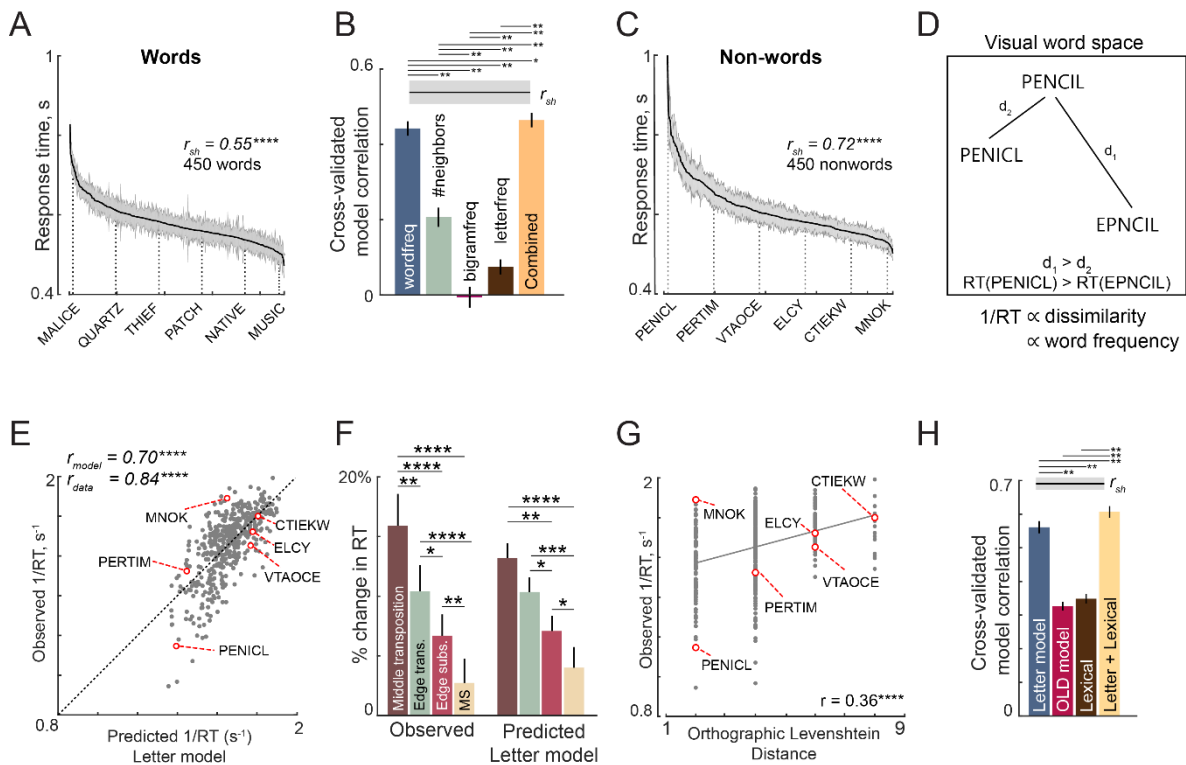
372  < 0.00005; Figure 4G).

373       We compared the letter model with two alternate models: the OLD model and

374  a model based on lexical factors. The OLD model is as described above. In the lexical

375  model, the nonword response time is modelled as a linear sum of log word frequency,

376  log mean bigram frequency of words, log mean bigram frequency of nonwords, #

377  orthographic neighbours, log letter frequency. Since all three models have different

378  numbers of free parameters, we compared their performance using cross-validation:

379  we trained each model on one-half of the subjects and evaluated it on the other half

380  of the subjects. The resulting cross-validated model fits are shown in Figure 4H. The

381  letter model outperformed both the OLD model and the lexical model (model

382  correlations: r = 0.56 ± 0.02, 0.33 ± 0.01 and 0.35 ± 0.01 for the neural, OLD and

383  lexical models; fraction of bootstrap samples with neural <  other models: p < 0.005;

384  Figure 4H). To be absolutely certain that the superior fit of the letter model was not

385  simply due to having more free parameters, we compared the lexical model with a

386   reduced version of the letter model with only 5 free parameters (SID model; Section

387   S5). Even this reduced model yielded fits were better than the lexical model (SID

388   model correlation: r = 0.48 ± .02). Finally, a combined model – in which the neural and

389   lexical model predictions were linearly combined – proved to explain more variance

390   than either model (Figure 4H).

391   In sum, we conclude that word response times are explained by word frequency

392   and nonword response times are explained by the distance between the nonword and

393   the nearest word calculated using the compositional neural code. As a further test of

394   the ability of this compositional code to explain word reading, we performed an

395   additional experiment in which subjects had to recognize the identity of a jumbled

396   word. Here too, response times were explained best by the letter model compared to

397   lexical and OLD models (Section S6).

398

399

**Figure 4. Lexical decision task behaviour (Experiment 5)**

(A) Response times for words in the lexical decision task, sorted in descending order. The solid line represents the mean categorization time for words and the shaded bars represent s.e.m. Some example words are indicated using dotted lines. The split-half correlation between subjects ($r_{sh}$) is indicated on the top.

(B) Cross-validated model correlation between observed and predicted word response times across all words for various models: log word frequency (*blue*), number of orthographic neighbours (*orange*), log mean bigram frequency (*purple*), log mean letter frequency (*cyan*) and a combined model containing all these factors (*red*). Shaded error bars indicate mean ± sd of the correlation across 1000 random splits of the observed data. The asterisk indicates statistical significance of the comparison obtained by estimating the fraction of bootstrap samples in which the observed difference was violated (* is $p < 0.05$, ** is $p < 0.005$).

(C) Response times for nonwords in the lexical decision task, sorted in descending order. Conventions as in (A).

(D) Schematic of visual word space, with one stored word (PENCIL) and two nonwords (PENICL & EPNCIL). We hypothesize that subjects would take longer to categorize a nonword when it is similar to a word, i.e. RT for PENICL would be larger than for EPNCIL. Thus, 1/RT would be proportional to this dissimilarity, and also to word frequency.

(E) Observed reciprocal response times for nonwords in the lexical decision task plotted against letter model predictions fit to the full dataset (450 nonwords). Some example nonwords are depicted.

(F) Percent change in response time (nonword-RT – word-RT)/word-RT for middle & edge letter transpositions and for middle & edge substitutions for observed data (*left*) and for letter model predictions (*right*). MS: middle substitution. In both cases, asterisks represent statistical significance comparing the means of the corresponding groups using a rank-sum test (* is $p < 0.05$, ** is $p < 0.005$, etc.).

429  (G) Observed reciprocal response times plotted against the Orthographic Levenshtein
430      Distance (OLD), a popular model for edit distance between strings.
431  (H) Cross-validated model correlation between observed and predicted nonword RTs
432      for the letter model, OLD model, lexical model and the combined neural+lexical
433      model. Conventions are as in (B).
434

435  **Brain activations during lexical decisions (Experiments 6-7)**

436      The above results show that visual discrimination of strings can be explained

437  using a letter-based compositional neural code, and that dissimilarities calculated

438  using this code can explain human performance on reading tasks. Here, we sought to

439  uncover the brain regions that represent this code and guide eventual lexical

440  decisions.

441      In Experiment 6, we recorded neural activations using fMRI while subjects

442  performed a lexical decision task. Since lexical decision times for nonwords can be

443  predicted using visual dissimilarity, we performed a separate experiment to directly

444  estimate visual dissimilarities using visual search (Experiment 7; see Methods).

445  Additionally, we estimated the semantic dissimilarity between words in order to

446  compare visual and semantic representations in different ROIs (see Methods).

447  Importantly, the visual search and semantic dissimilarities were uncorrelated (r = 0.03,

448  p = 0.55), thereby allowing us to identify regions with distinct or overlapping

449  visual/semantic representations. The visual and semantic representations are

450  visualized in Section S7.

451      We identified several possible regions of interest (ROIs) using a combination of

452  functional localizers and anatomical considerations (see Methods). These included the

453  early and mid-level visual areas (V1-V3 & V4), the object-selective lateral occipital

454  region (LO), and two language areas: the visual word form area (VWFA) which

455  selectively responds to words and a broad region in the temporal gyrus reading

456    network (TG). The flattened brain map of a representative subject with these ROIs is

457    shown in Figure 5A.

458         For the main event-related block, subjects had to make a response on each

459    trial to indicate whether a string displayed on the screen was a word or not. The stimuli

460    consisted of 5-letter words, nonwords. Subjects also viewed single letters, to which

461    they had to make no response. Subjects were highly accurate (mean ± std of accuracy:

462    94 ± 4%) and showed consistent response time variations (split-half correlation

463    between odd and even subjects: $r_{sh}$ = 0.54 & 0.79 for words and nonwords, p <

464    0.00005). As before, the lexical decision time for words was negatively correlated with

465    word frequency (r = -0.42, p < 0.05). Likewise, the lexical decision times for nonwords

466    were strongly correlated with the word-nonword dissimilarity measured in visual

467    search in Experiment 7 (r = -0.68, p < 0.00005). These results reconfirm the findings

468    of the previous experiment performed outside the scanner.

469         We first compared the overall brain activation levels for words, nonwords and

470    letters in each ROI. While V4 showed greater activation for words compared to

471    nonwords, VWFA and TG regions showed greater activation to nonwords compared

472    to words, presumably reflecting greater engagement to discriminate nonwords that are

473    highly similar to words (Section S7). Although the visual regions did not show

474    differential overall activations, there could still be differential activation at the

475    population level for words and nonwords. To assess this possibility, we built linear

476    classifiers to discriminate words from nonwords using the voxel population activity in

477    each ROI. This revealed above-chance classification in all ROIs. Further,

478    discriminating words from nonwords was significantly easier for nonwords that were

479    obtained by substituting letters compared to those obtained by transposing letters

480    (Section S7). Correspondingly, in behaviour, subjects were faster at responding to

481    substituted nonwords compared to transposed nonwords (response times, mean ± sd:

482    1.03 ± 0.08 s for 16 substituted nonwords, 1.20 ± 0.15 s for 16 transposed nonwords,

483    p < 0.005, rank-sum test comparing average response times). A detailed analysis of

484    these results is presented in Section S7.

485

486    **Neural correlates of the compositional letter code**

487        Next we sought to compare the neural representations in each ROI with visual

488    search and semantic representations. The visual search and semantic representations

489    can be quite distinct, as depicted in Figure 5B: in visual search space, TRAIL and

490    TRIAL can be quite similar since one is obtained from the other by transposing letters,

491    but the word PATH is quite distinct. By contrast, in semantic space, TRAIL and PATH

492    have similar meanings and usage whereas TRIAL is quite distinct. Indeed, visual

493    search and semantic dissimilarities across words were uncorrelated for the words in

494    experiment (r = 0.03, p = 0.55).

495        To investigate these issues, we calculated the neural dissimilarity between

496    each stimulus pair in a given ROI as the cross-validated Mahalanobis distance

497    between the voxel-wise activations evoked by the two stimuli. We then averaged this

498    dissimilarity across subjects to get an average neural dissimilarity for that ROI. We

499    then compared this neural dissimilarity in each ROI with visual dissimilarities estimated

500    from visual search. This match to visual search dissimilarity is shown in Figure 5C.

501    Among the ROIs tested, only the LO dissimilarities showed a significant correlation

502    (correlation between 1024 pairwise dissimilarities involving $^{32}C_2$ words, $^{32}C_2$

503    nonwords, and 32 word-nonword pairs: r = 0.16, p < 0.00005; Figure 5C). A searchlight

504    analysis confirmed that the match to visual search dissimilarities was strongest in a

505    region centred around the bilateral LO region (Section S7). Thus, neural dissimilarity

506 in the LO region match best with the visual dissimilarities observed in visual search.

507 We therefore conclude that LO is the likely neural substrate for the compositional letter

508 code.

509 To further investigate the link between the compositional letter code and the LO

510 representation, we performed several additional analyses. First, we asked whether the

511 neural activation of each voxel in LO could be explained using a linear sum of the

512 single letter activations. Importantly, these model fits were equally good for words and

513 nonwords. This parallels our finding that dissimilarity in visual search was predicted

514 equally well for word-word and nonword-nonword pairs (Figure 3H). Both these results

515 suggest that there are no specialized detectors for letter combinations (Section S7).

516 Second, we confirmed that both the neural tuning for single letters, and the summation

517 weights estimated from the behavioural data in the letter model were qualitatively

518 similar to the observed tuning for single letters and summation weights observed in

519 the voxel activations for the LO region (Section S7).

520 In sum, we conclude that the LO region is the likely neural substrate for the

521 compositional letter code predicted from behaviour.

522

523 **Neural basis of semantic space**

524 Next we compared neural representations in each ROI to semantic space. The

525 match to semantic space was significant only in the LO and TG regions (correlation

526 between 496 pairwise dissimilarities between words: $r = 0.18 \pm 0.05$ for LO, $0.22 \pm$

527 0.04 for TG; Figure 5D).

528 The above analysis shows that neural activations in LO are correlated with both

529 visual search and semantic dissimilarities, but these correlations cannot be directly

530 compared since they are based on different pairs of stimuli. To investigate whether the

531    neural representation in LO matches better with visual search or with semantic space,

532    we compared the match for word-word pairs alone. This revealed no significant

533    difference between the two correlations (r = 0.16 ± .04 for LO with visual search, r =

534    0.16 ± 0.05 for LO with semantic space; p = 0.49 across 1000 bootstrap samples). We

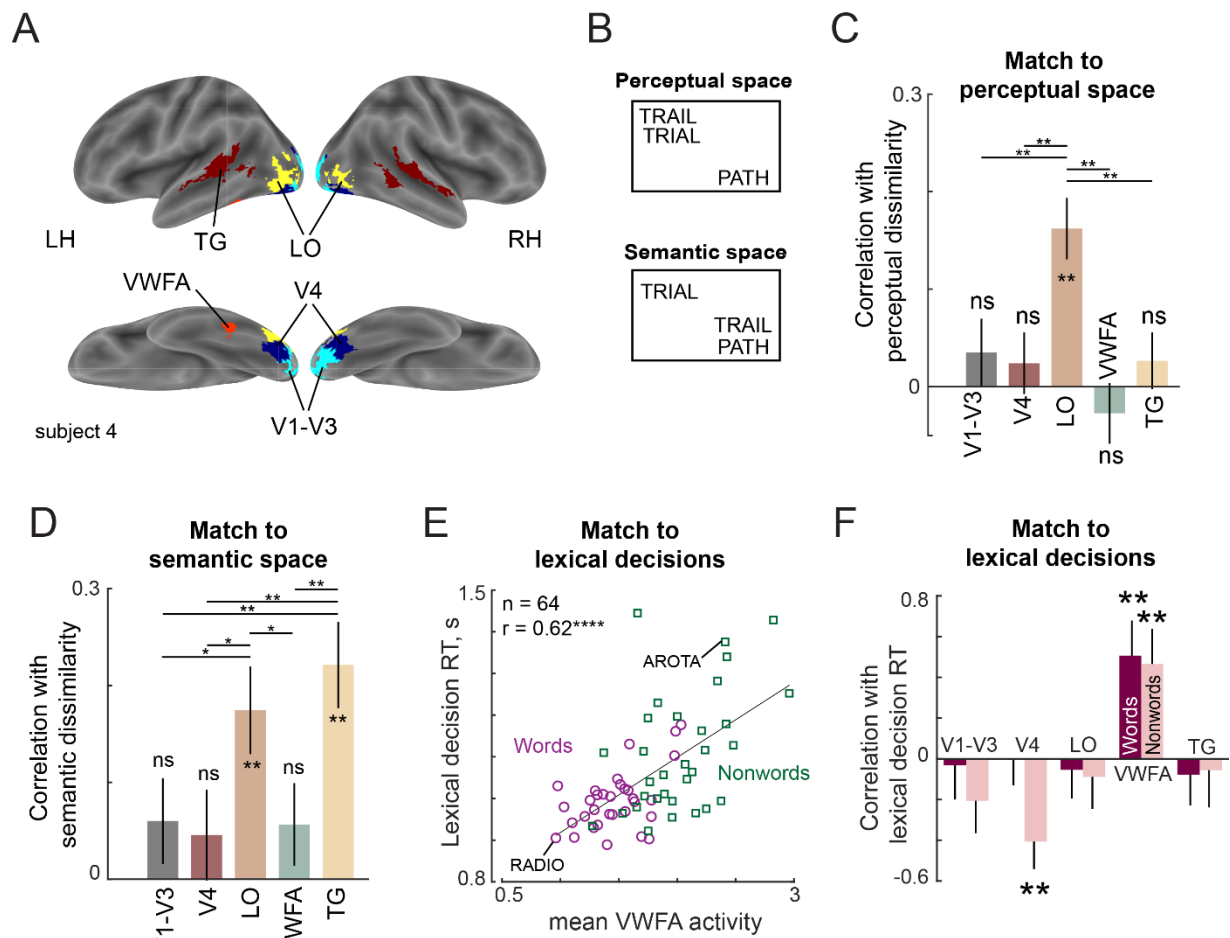535    conclude that both LO and TG regions represent semantic space.

536

537    **Neural basis of lexical decisions**

538      If the LO region represents each string (word or nonword) using a compositional

539    code, then according to the preceding experiments, lexical decisions for words and

540    nonwords must involve some comparison with stored word representations. Recall

541    that lexical decision times for words are correlated with word frequency, and lexical

542    decision times for nonwords are correlated with word-nonword dissimilarity. We

543    therefore asked whether these lexical decision times are correlated with the average

544    activity (across voxels & subjects) in  a given ROI. The resulting correlations are shown

545    in Figure 5F. Across the ROIs, only the VWFA showed a consistently positive

546    correlation with lexical decision times for both words and nonwords (Figure 5E). A

547    searchlight analysis confirmed that there was indeed a peak in the correlation with

548    lexical decision times centred on the VWFA, with additional peaks in the parietal and

549    frontal regions (Section S7). This is consistent with word frequency effects observed

550    in these regions (Kronbichler et al., 2004), but we have observed similar effects for

551    nonwords as well. We conclude that lexical decisions are driven by the VWFA.

552

553

**Figure 5. Lexical task fMRI (Experiment 6)**

(A) ROIs for an example subject, showing V1–V3 (cyan), V4 (blue), LO (yellow), VWFA (red) and TG (maroon).

(B) Example difference between perceptual and semantic spaces. In perceptual space, the representation of TRAIL is closer to its visual similar counterpart TRIAL, whereas in semantic space, its representation is closer to its synonym PATH.

(C) Correlation between neural dissimilarity in each ROI with behavioural dissimilarity for strings (Experiment 7). Error bars indicate standard deviation of the correlation between the group behavioural dissimilarity and ROI dissimilarities calculated repeatedly by resampling of dissimilarity values with replacement across 1000 iterations. Asterisks along the length of each bar indicate statistical significance of the correlation between group behaviour and group ROI dissimilarity (** is p < 0.005 across 1000 bootstrap samples). Horizontal lines indicate the fraction of bootstrap samples in which the observed difference was violated (* is p < 0.05, ** is p < 0.005, etc.). All significant comparisons are indicated.

(D) Correlation between neural dissimilarity in each ROI with semantic dissimilarity for words. Other details are same as in (C).

(E) Correlation between mean VWFA activity (averaged across subjects and voxels) with mean lexical decision time for both words (purple circles) and nonwords (green squares). Each point corresponds to one string and example word and nonword is highlighted. Asterisks indicate statistical significance (**** is p < 0.00005).

(F) Correlation between lexical decision time and mean activity within each ROI separately for words and nonwords. Error bars indicate standard deviation across 1000 bootstrap splits. Asterisks indicate statistical significance (** is p < 0.005).

**DISCUSSION**

579

580      Our main finding is that viewing a string activates a compositional letter code,

581   consisting of neurons tuned to single letters whose response to longer strings is a

582   linear sum of single letter responses. This code accurately explains human

583   performance on visual search as well as word reading tasks. It is encoded by the LO

584   region in the high-level visual cortex, and subsequent comparisons required for lexical

585   decisions are computed in the VWFA. Below we discuss these findings in relation to

586   the existing literature.

587

588   **Relation to models of reading**

589      Our compositional letter code stands in stark contrast to existing models of

590   reading. Existing models of reading assume explicit encoding of letter position and do

591   not account for letter shape (Gomez et al., 2008; Davis, 2010; Norris and Kinoshita,

592   2012; Norris, 2013). By contrast, our model encodes letter shape explicitly and position

593   implicitly through asymmetric spatial summation. The implicit coding of letter position

594   avoids the complication of counting transpositions (Yarkoni et al., 2008; Yap et al.,

595   2015). Our model can thus easily be extended to any language by simply estimating

596   letter dissimilarities and then estimating the unknown summation weights.

597      Unlike existing models of reading, our compositional letter code is neurally

598   plausible and grounded in well-known principles of object representations. The first

599   principle is that images that elicit similar activity across neurons in high-level visual

600   cortex will appear perceptually similar (Op de Beeck et al., 2001; Sripati and Olson,

601   2010a; Zhivago and Arun, 2014). This is non-trivial because it is not necessarily true

602   in lower visual areas or in image pixels (Ratan Murty and Arun, 2015). We have turned

603   this principle around to construct artificial neurons whose shape tuning matches visual

604    search. The second principle is that the neural response to multiple objects is typically

605    the average of the individual object responses (Zoccolan et al., 2005; Sripati and

606    Olson, 2010b) that can be biased towards a weighted sum (Ghose and Maunsell,

607    2008; Bao and Tsao, 2018). Finally, we note that our letter code assumes no explicit

608    calculations of letter position in a word, since the neurons in our model only need to

609    be tuned for retinal position. We speculate that these neurons may be tuned not only

610    to retinal position but to the relative size and position of letters, as observed in high-

611    level visual cortex (Sripati and Olson, 2010a; Vighneshvel and Arun, 2015).

612

613    **Relation to theories of word recognition**

614        We have found that lexical decisions for nonwords are driven by the dissimilarity

615    between the viewed string and the nearest word. This idea is consistent with the well-

616    known Interactive Activation model (McClelland and Rumelhart, 1981; Rumelhart and

617    McClelland, 1982), where viewing a string activates the nearest word representation.

618    However, the Interactive Activation model does not explain lexical decisions or

619    scrambled word reading, and also does not integrate letter shape and position into a

620    unified code. Our findings are consistent with previous work showing that nonword

621    responses are influenced by the number of orthographic neighbours (Yap et al., 2015).

622    Likewise, we found word frequency to be a major factor influencing lexical decisions,

623    in keeping with previous work (Ratcliff et al., 2004; Dufau et al., 2012; Yap et al., 2015).

624    We have gone further to demonstrate a unified letter-based code that integrates letter

625    shape and position, and localized the underlying neural substrates of the letter code

626    to the LO region, and the comparison process to the VWFA. We propose that the

627    compositional shape code provides a quick match to unscramble a word, failing which

628    subjects may initiate more detailed symbolic manipulation.

629        The success of our letter code challenges the widely held belief that reading

630        expertise should lead to the formation of specialized bigram detectors (Dehaene et al.,

631        2005, 2015; Grainger, 2018). The presence of these specialized detectors should have

632        caused larger model errors for valid words and frequent n-grams, but we observed no

633        such trend (Figure 3). So what happens to visual letter representations upon expertise

634        with reading? Our comparison of upright and inverted bigrams suggests that reading

635        should increase letter discrimination and increase the asymmetry of spatial summation

636        (Figure 3D,E). This is consistent with differences in letter position effects for symbols

637        and letters (Chanceaux and Grainger, 2012; Scaltritti et al., 2018). We propose that

638        both processes may be driven by visual exposure: repeated viewing of letters makes

639        them more discriminable (Mruczek and Sheinberg, 2005), while viewing letter

640        combinations induces asymmetric spatial weighting. Whether these effects require

641        active discrimination such as letter-sound association training or can be induced even

642        by passive viewing will require comparing letter string discrimination under these

643        paradigms.

644

645        **Neural basis of word recognition**

646        Our brain imaging results further elucidate the neural basis of lexical decisions.

647        We have shown that viewing a string activates this compositional letter code in the LO

648        region, and that subsequent comparisons to stored words is driven by the VWFA. We

649        have found that lexical decision response times for both words and nonwords are

650        strongly predicted by the VWFA activity. Since lexical decision times for words are

651        linked to word frequency, this finding implies that VWFA is sensitive to word frequency.

652        This has been confirmed by previous studies (Kronbichler et al., 2004; Dehaene et al.,

653        2015). We have also found that VWFA activity is strongly predictive of nonword

654 response times, which in turn are modulated by the dissimilarity of the nonword to the

655 nearest word. This finding is somewhat surprising because of VWFA's status as a

656 word processing area, but consistent with previous suggestions that it stores word

657 representations (Vinckier et al., 2007) and is modulated by orthographic similarity

658 (Baeck et al., 2015). Our results suggest that VWFA might be involved in comparing

659 viewed strings with known words. We speculate that differences in this comparison

660 process could explain the contradictory findings about VWFA activation to words and

661 nonwords (Baeck et al., 2015).

662

663 **Does the compositional letter code explain orthographic processing?**

664 Our letter code explains many orthographic processing phenomena reported in

665 the literature. Its integrated representation of both letter shape and position explains

666 both letter transposition and substitution effects and their relative importance (Figure

667 4F). Its asymmetric spatial weighting favouring the first letter (Section S3), explains

668 the first-letter advantage observed previously (Scaltritti et al., 2018). It also explains

669 why increasing letter spacing can benefit reading in poor readers, presumably

670 because it increases asymmetry in spatial summation (Zorzi et al., 2012).

671 To elucidate how various jumbled versions of a word are represented according

672 to this neural code, we calculated responses of the letter model trained on data from

673 Experiment 4, and visualized the distances using multidimensional scaling (Figure 6A).

674 It can be seen transposing the edge letters (OFRGET) results in a bigger change than

675 transposing the middle letters (FOGRET), thus explaining many transposed letter

676 effects (Norris, 2013). Likewise, it can be seen that substituting a dissimilar letter

677 (FORXET) leads to a large change compared to substituting a similar letter (FORCET).

678 Replacing G with C in FORGET leads to a smaller change than replacing with X, thus

679    explaining how priming is stronger when similar letters are substituted (Marcet and

680    Perea, 2017). Finally, the letter subset FRGT is closer to FORGET than the same

681    letters reversed (TGRF), thereby explaining subset priming (Grainger and Whitney,

682    2004; Dehaene et al., 2005).

683        Finally, as a powerful demonstration of this code, we used it to arbitrarily

684    manipulate reading difficulty along a sentence (Figure 6B), or across multiple

685    transpositions and even number substitutions (Figure 6C). We propose that this

686    compositional neural code can serve as a powerful baseline for the purely visual

687    shape-based representation triggered by viewing words, thereby enabling the study of

688    higher order linguistic influences on reading processes.

689

690    **Relation between word recognition and reading sentences**

691        We have shown that word recognition can be explained using a compositional

692    visual code based on single letters. While this is an important first step in

693    understanding how we read single words, reading sentences involves sampling many

694    words with each eye movement (Rayner, 1998). Our ability to sample multiple letters

695    or words at a single glance is limited by two factors. The first is our visual acuity, which

696    reduces with eccentricity. The second is crowding, by which letters become

697    unrecognizable when flanked by other letters – this effect increases with eccentricity

698    (Pelli and Tillman, 2008).

699        The visual search experiments in our study involved searching for an oddball

700    target (consisting of multiple letters) among multiple distractors. This would most

701    certainly have involved detecting and making saccades to peripheral targets. By

702    contrast, the word recognition tasks in our study involved subjects looking at words

703    presented at the fovea. Our finding that visual search dissimilarity explains word

704 recognition then implies that shape representations are qualitatively similar in the

705 fovea and periphery. Furthermore, the structure of the letter model suggests a possible

706 mechanistic explanation for crowding. Neural responses might show greater sensitivity

707 to spatial location at the fovea compared to the periphery, leading to more

708 discriminable representations of multiple letters. Alternatively, neural responses to

709 multiple letters might be more predictable from single letters at the fovea but not in the

710 periphery. Both possibilities would predict reduced recognition with closely spaced

711 flankers. Distinguishing these possibilities will require testing neural responses in

712 higher visual areas to single letters and multi-letter strings of both familiar and

713 unfamiliar scripts. Ultimately understanding reading fully will require not only asking

714 how letters combine to form words, but how words combine to form larger units of

715 meaning (Pallier et al., 2011; Nelson et al., 2017).

716

A

**Predicted visual word space**

FORXET

FORCET· FGROET
**FORGET**
FRGT·
FGEORT·
FOGRET

FGORET
·FROGET

OFRGET·

·TGRF

B

Predicted reading difficulty

HUAMN MIDN DOSE NOT RDEA YVERE TRETLE BY FSLTEI

HUAMN DNMI DEOS NOT DAER EVREY ETTELR BY ITSLEF

ANMHU DINM SOED NOT RDEA ERVEY LTETER BY ITSLEF

C

REAIDNG IS A RECNET CULTRUAL INVENITON

REAO1NG IS A R3C3NT CUL7UR4L INV3N710N

RDNIEAG IS A REENCT CLRTUAUL IONETNIVN

GDNREAI IS A CTNREE TARCLLUU OTIIVNNNE

Predicted reading difficulty

**Figure 6. Predicting reading difficulty using the letter model**

(A) Visual word space predicted by the letter model for a word (FORGET) and its jumbled versions. Letter model predictions were based on training the model on compound words (Experiment 4). The plot was obtained by performing multidimensional scaling on the pairwise dissimilarities between strings predicted by the letter model. It can be seen that classic features of orthographic processing are captured by the letter model, including priming effects such as FRGT (*green*) being more similar to FORGET than TGRF (*red*).

(B) The letter model can be used to sort jumbled words by their reading difficulty, allowing us to create any desired reading difficulty profile along a sentence. *Top row*: Sentence with increasing reading difficulty. *Middle row*: sentence with fluctuating reading difficulty. *Bottom row*: sentence with decreasing reading difficulty.

(C) The letter model yields a composite measure of reading difficulty that combines letter substitution and transposition effects. Sentences with digit substitutions (*second row*) can thus be placed along a continuum of reading difficulty relative to other sentences (*first, third and fourth rows*) with increasing degree of scrambling.

737 **METHODS**

738       All subjects had normal or corrected-to-normal vision and gave informed

739 consent to an experimental protocol approved by the Institutional Human Ethics

740 Committee of the Indian Institute of Science. All subjects were fluent English-speaking

741 students at the institute, where English is the medium of instruction. All subjects were

742 multi-lingual and knew at least one other Indian language apart from English.

743

744 **Experiment 1 – Single letter searches**

745 *Procedure.* A total of 16 subjects (8 males, 24.4 ± 2.5 years) participated in this

746 experiment. Subjects were seated comfortably in front of a computer monitor placed

747 ~60 cm away under the control of custom programs written in Psychtoolbox (Brainard,

748 1997) and MATLAB. In all experiments, we selected sample sizes based on our

749 previous studies which yielded highly consistent data (Agrawal et al., 2019).

750 *Stimuli.* Single letter images were created using the Arial font. There were 62 stimuli

751 in all comprising 26 uppercase letters (A-Z), 26 lowercase letters (a-z), and 10 digits

752 (0-9). Uppercase stimuli were scaled to have a height of 1°.

753 *Task.* Subjects were asked to perform an oddball search task without any constraints

754 on eye movements. Each trial began with a fixation cross shown for 0.5 s followed by

755 a 4x4 search array (measuring 40° by 25°). The search array always contained only

756 one oddball target with 15 identical distractors. Subject were instructed to locate the

757 oddball target as quickly and as accurately as possible, and respond with a key press

758 ('Z' for left, 'M' for right). A red line divided the screen in two halves. The search display

759 was turned off after the response or after 10 seconds, whichever was sooner. All

760 stimuli were presented in white against a black background. Incorrect or missed trials

761 were repeated after a random number of other trials. Subjects completed a total of

762 3,782 correct trials ($^{62}C_2$ letter pairs x 2 repetitions with either letter as target once).

763 For each search pair, the oddball target appeared equally often on the left and right

764 sides so as to avoid creating any response bias. Only correct responses were

765 considered for further analysis. The main experiment was preceded by 20 practice

766 trials involving unrelated stimuli.

767 *Data Analysis.* Subjects were highly accurate on this task (mean ± std: 98 ± 1%).

768 Outliers in the reaction times were removed using built-in routines in MATLAB (*isoutlier*

769 function, MATLAB R2018a). This function removes any value greater than three

770 scaled absolute deviations away from the median, and was applied to each search

771 pair separately. This step removed 6.8% of the response time data.

772

773 **Estimation of single letter tuning using multidimensional scaling**

774 To estimate neural responses to single letters from the visual search data, we

775 used a multidimensional scaling (MDS) analysis. We first calculated the average

776 search time for each letter pair by averaging across subjects and trials. We then

777 converted this search time (RT) into a distance measure by taking its reciprocal (1/RT).

778 This is a meaningful measure because it represents the underlying rate of evidence

779 accumulation in visual search (Sunder and Arun, 2016), behaves like a mathematical

780 distance metric (Arun, 2012) and combines linearly with a variety of factors (Pramod

781 and Arun, 2014, 2016; Sunder and Arun, 2016). Next we took all pairwise distances

782 between letters and performed MDS to embed letters into n dimensions, where we

783 varied n from 1 to 15. This yielded n-dimensional coordinates corresponding to each

784 letter, whose distances matched best with the observed distances. We then took the

785    activation of each letter along a given dimension as the response of a single neuron.

786    Throughout we performed MDS embedding into 10 dimensions, resulting in single

787    letter responses of 10 neurons. We obtained qualitatively similar results on varying

788    this number of dimensions.

789

790    **Estimation of data reliability**

791       To obtain upper bounds on model performance, we reasoned that any model

792    can predict the data as well as the consistency of the data itself. Thus, a model trained

793    on one half of the subjects can only predict the other half as well as the split-half

794    correlation $r_{sh}$. This process was repeated 100 times to obtain the mean and standard

795    deviation of the split-half correlation. However when a model is trained on all the data,

796    the upper bound will be larger than the split-half correlation. We obtained this upper

797    bound, which represents the reliability of the entire dataset ($r_{data}$) by applying a

798    Spearman-Brown correction on the split-half correlation, as given by $r_{data} = 2r_{sh}/(r_{sh}+1)$.

799

800    **Experiment 2 – Bigram searches**

801       A total of 8 subjects (5 male, aged 25.6 ± 2.9 years) took part in this experiment.

802    We chose seven uppercase letters (A, D, H, I, M, N, T) and combined them in all

803    possible ways to obtain 49 bigram stimuli. These letters were chosen to maximise the

804    number of two-letter words e.g.  HI, IT, IN, AN, AM, AT, AD, AH, and HA. Letters

805    measured 3° along the longer dimension. Subjects completed 2352 correct trials ($^{49}C_2$

806    search pairs x 2 repetitions). All other details were identical to Experiment 1.

807    Letter/Bigram frequencies were obtained from an online database

808    (http://norvig.com/mayzner.html).

Page 36 of 47

809    *Data Analysis.* Subjects were highly accurate on this task (mean ± std: 97.6 ± 1.8%).

810    Outliers in the reaction times were removed using built-in routines in MATLAB (*isoutlier*

811    function, MATLAB R2018a). This step removed 8% of the response time data.

812

813    **Estimating letter model parameters from observed dissimilarities**

814    The total dissimilarity between two bigrams in the letter model is calculated by

815    calculating the average dissimilarity across all neurons. For each neuron, the

816    dissimilarity between bigrams AB & CD is given by:

817    $$d(AB, CD) = |r_{AB} - r_{CD}| = |(w_1 r_A + w_2 r_B) - (w_1 r_C + w_2 r_D)|$$

818    where $r_A, r_B, r_C$ and $r_D$ are the responses of the neuron to individual letters A, B,

819    C and D respectively (derived from single letter dissimilarities), and $w_1, w_2$ are the

820    spatial summation weights for the first and second letters of the bigram. Note that

821    $w_1, w_2$ are the only free parameters for each neuron.

822    To estimate the spatial weights of each neuron, we adjusted them so as to

823    minimize the squared error between the observed and predicted dissimilarity. This

824    adjustment was done using standard gradient descent methods starting from randomly

825    initialized weights (*nlinfit* function, MATLAB R2018a). We followed a similar approach

826    for experiments involving longer strings.

827

828    **Experiment 3 – Upright and inverted bigrams**

829    *Methods.* A total of 8 subjects (6 males, aged 24 ± 1.5 years) participated in this

830    experiment. Six uppercase letters: A, L, N, R, S, and T were combined in all pairs to

831    form a total of 36 stimuli. These uppercase letters were chosen because their images

832    change when inverted (as opposed to letters like H that are unaffected by inversion),

833    and were chosen to maximize the occurrence of frequent bigrams. The same stimuli

834    were inverted to create another set of 36 stimuli. Detailed analyses for this experiment

835    are presented in Section S2.

836

837    **Experiment 4 – Compound words**

838    A total of 8 subjects (4 female, aged 25 ± 2.5 years) participated. Twelve 3-

839    letter words were chosen: ANY, FOR, TAR, KEY, SUN, TEA, ONE, MAT, GET, PAD,

840    DAY, POT. Each word was jumbled to obtain twelve 3-letter nonwords containing the

841    same letters. The 12 words were combined to form 36 compound words (shown in

842    Section S3), such that they appeared equally on the left and right half of the compound

843    words. Detailed analyses for this experiment are included in Section S3.

844    *Calculation of orthographic Levenshtein distance:* For each search pair, we estimated

845    the OLD metric using built-in MATLAB function "editdistance". This function estimates

846    the number of insertions, deletions, or substitutions are required to convert one string

847    to other. In this study, the substitution cost has a value of 2. We obtained qualitatively

848    similar results with other choices of substitution cost.

849

850    **Experiment 5 – Lexical decision task**

851    *Procedure.* A total of 16 subjects (9 male, aged 24.8 ± 2.1 years) participated in this

852    task as well as the jumbled word task.

853    *Stimuli.* The stimuli comprised 450 words + 450 nonwords. The nonwords were either

854    random strings or made by modifying the 450 words in some way, as detailed in the

855    table below.

| | Variations of word ABCDE | 4 letter words | 5 letter words | 6 letter words | Total |
|---|---|---|---|---|---|
| 1) | *Edge transpositions*: BACDE or ABCED | 15 | 15 | 20 | 50 |
| 2) | *Middle transposition*: ACBDE or ABDCE | 15 | 15 | 20 | 50 |
| 3) | *2 step edge transposition*: CBADE or ABEDC | 0 | 20 | 30 | 50 |
| 4) | *2 step middle transposition*: ADCBE | 0 | 20 | 30 | 50 |
| 5) | *Random transposition*: CDABE, ACDBE, etc. | 25 | 35 | 40 | 100 |
| 6) | *Edge Substitution*: MZCDE or ABCMZ | 15 | 15 | 20 | 50 |
| 7) | *Middle Substitution*: ABMZE | 15 | 15 | 20 | 50 |
| 8) | *Random substitution and permutation*: MACZE, AMDEZ, etc. | 15 | 15 | 20 | 50 |
| | Total | 100 | 150 | 200 | 450 |

**Table 1: Non-word stimuli in lexical decision task (Experiment 5).**

*Task.* Each trial began a fixation cross shown for 0.75 s followed by a letter string for 0.2 s after which the screen went blank. The trial ended either with the subject's response or after at most 3 s. Subjects were instructed to press 'Z' for words and 'M' for nonwords as quickly and accurately as possible. All stimuli were presented at the centre of the screen and were white letters against a black background. Before starting the main task, subjects were given 20 practice trials using other words and nonwords not included in the main experiment.

*Data Analysis.* Some nonwords were removed from further analysis due to low accuracy (n = 8, average accuracy <20%). Subjects made accurate responses for both words and nonwords (mean ± std of accuracy: 96 ± 2 % for words, 95 ± 3% for nonwords). Outliers in the reaction times were removed using built-in routines in MATLAB (*isoutlier* function, MATLAB R2018a).

**Experiment 6 (Lexical Decision Task – fMRI)**

A total of 17 subjects (10 males, 25 ± 4.2 years) participated in this experiment. All subjects were screened for safety and comfort beforehand to avoid adverse outcomes in the scanner.

*Stimuli:* The functional localizer block included English words, objects, scrambled words, and scrambled objects. In each run, 14 images were randomly selected from a pool of images. The English words list comprised of 90 five-letter words. Each word was divided into grids of dimension 9x3. Scrambled words were generated by randomly shuffling the grids. Object pool comprised of 80 man-made objects. To generate scrambled objects, the phase of the Fourier transformed images was scrambled and then reconstructed back using inverse Fourier transform. The object images were about 4.5° along the longer dimension and the height of the word stimuli subtended 2° of visual angle.

The event block consisted of 10 single letters and 64 five-letter strings (32 words and 32 nonwords formed using these single letters). The stimulus set comprised of 64 five-letter words and nonwords. The words were chosen from a wide range of frequency of occurrence and the nonwords were created by manipulating the chosen words i.e. They were: 1) 8-middle transposed version of words, 2) 8-edge transposed version of words, 3) 8-middle substituted version of words, and 4) 8-edge substituted version of words. The stimuli subtended 2° in height, which was the same as in the localizer block. All stimuli were presented as white against a black background.

*Procedure:* In the localizer block, a total of 16 images were presented for 0.8 s with an inter stimulus interval of 0.2 s. There were 14 unique stimuli and 2 of them repeated at random time point, in which subjects performed one-back task. Each block ended

900    with a blank screen with fixation cross present for 4 s. Thus, each block lasted 20 s.

901    Each block was repeated thrice in each run.

902        In the event-related design block, an image was presented at the centre of the

903    screen for 300ms followed by 3.7s of blank screen with a fixation cross. In a run, all

904    74 stimuli were presented once along with 16 trials of fixation cross to jitter inter

905    stimulus interval. Hence there were a total of 92 trials including 4s fixation trials at the

906    start and end of each run. Each run lasted 376 s. Subjects performed lexical decision

907    task only on strings and were instructed to not press any key for single letters. Overall,

908    subjects completed 2 runs of localizer block, 8 runs of event block and a structural

909    scan block.

910    *Data acquisition:* Subjects viewed images in a mirror-based projection system.

911    Functional MRI data was acquired using a 32-channel head coil on a 3T Siemens

912    Skyra scanner at HealthCare Global Hospital, Bengaluru. Functional scans were

913    performed using a T2*-weighted gradient-echo-planar imaging sequence with the

914    following parameters: TR = 2s, TE = 28ms, flip angle = 79º, voxel size = 3x3x3 mm$^3$,

915    field of view = 192x192 mm$^2$, and 33 axial-oblique slices covering the whole brain.

916    Anatomical scans were performed using T1-weighted images with the following

917    parameters: TR = 2.30s, TE = 1.99ms, flip angle = 9°, voxel size = 1x1x1 mm$^3$, field

918    of view = 256x256x176 mm$^3$.

919    *Data preprocessing:* All raw fMRI data were processed using custom built MATLAB

920    scripts        that        depended        on        SPM        12        toolbox

921    (https://www.fil.ion.ucl.ac.uk/spm/software/spm12/).  Raw  images  were  realigned,

922    slice-time corrected, co-registered with the anatomical image, segmented, and finally

923    normalized to the MNI305 anatomical template. The results were qualitatively similar

924    without normalization. Smoothing operation was performed only on functional localizer

925    blocks using a Gaussian kernel with FWHM of 5 mm. All SPM parameters were set to

926    default and the voxel size after normalization was set to 3x3x3 mm$^3$. Prior to

927    normalization, the data was preprocessed using GLMdenoise v1.4 (Kay et al., 2013).

928    This step improved the signal-to-noise ratio in the data by regressing out the noise

929    pattern common across all the voxels in the brain. The noise pattern is estimated from

930    voxels unrelated to the task. The activity corresponding to each condition was

931    estimated by modelling the denoised data using a generalized linear model (GLM) in

932    SPM after removing the low frequency drift using a high-pass filter with a cutoff at

933    128s. The event block data was modelled using 89 regressors (74 stimuli + 1 fixation

934    + 6 motion regressors + 8 runs). The localizer block data was modelled using 13

935    regressors (4 stimuli + 1 fixation + 6 motion regressors + 2 runs).

936    *ROI definitions:* All the regions of interest (ROI) were defined using functional localizer

937    while taking the anatomical location into consideration. Early visual area was defined

938    as the region that responds more to the scrambled object than fixation cross. This

939    functional region was further parsed into V1-V3 and V4 using an anatomical mask

940    from SPM anatomy toolbox (Eickhoff et al., 2005). Lateral Occipital (LO) region was

941    defined as a group of voxels that responded more to objects than scrambled objects.

942    The voxels in the LO region was restricted to Inferior Temporal Gyrus, Inferior Occipital

943    Gyrus, and Middle Occipital Gyrus. These anatomical regions were obtained from

944    Tissue Probability Map (TPM) labels in SPM 12. Visual Word Form Area (VWFA) was

945    defined as a region that responded more for words than scrambled words within

946    fusiform Gyrus. The activity for known words was also higher in Superior and Middle

947    Temporal regions. These groups of voxels were grouped under Temporal Gyrus (TG)

948    label. For each contrast, voxel-level threshold of $p < 0.001$ (uncorrected) or cluster

949  level threshold $p < 0.05$ (FWE correction) was used to obtain a contiguous region. For

950  one subject, very few VWFA voxels cross the pre-specified threshold. Hence, the

951  threshold was lowered to $p = 0.1$ (uncorrected).  The VWFA voxels were restricted to

952  top-40 voxels (based on T-value in the function localizer contrast). All these regions

953  were visualized on the cortical surface using BSPMVIEW toolbox

954  (http://www.bobspunt.com/bspmview/).

955

956  *Calculation of neural dissimilarity (fMRI).* For each ROI and subject, the pair-wise

957  dissimilarity between any two image pairs was computed using the cross-validated

958  Mahalanobis distance in the RSA toolbox (Nili et al., 2014). Outliers in dissimilarity

959  values across subjects were removed using built-in routines in MATLAB (*isoutlier*

960  function, MATLAB R2018a). The median dissimilarity across all the subjects was

961  considered for further analysis. We obtained qualitatively similar results for other

962  distance measures.

963

964  *Calculation of semantic dissimilarity.* The semantic distance between every pair of

965  words was computed as the cosine distance between the GloVe (Pennington et al.,

966  2014) feature vectors activated by the two words, using the MATLAB function

967  *word2vec.* These features are based on the co-occurrence statistics of words in a large

968  text corpus, and therefore reflect semantic dissimilarity rather than visual dissimilarity.

969

970  **Experiment 7 (5-letter string searches)**

971       A total of 11 subjects (6 males, 26 ± 2.7 years) participated in this experiment,

972  of which xx also participated in Experiment 6. Stimuli were identical to Experiment 6,

973    except that they were scaled down to a height of $1^0$ to allow placement in a visual

974    search array. Subjects performed a total of 2048 correct trials ($^{32}C_2$ search pairs x 2

975    conditions (words and nonwords) + 32 word-nonword pairs x 2 repetitions). All trials

976    were interleaved, and incorrect/missed trials appeared randomly later in the task but

977    were not analyzed. All other details were identical to Experiment 1.

978

979    *Data Analysis.* Subjects were highly accurate on this task (mean ± std: 98.6 ± 1%).

980    Outliers in the reaction times were removed using built-in routines in MATLAB (*isoutlier*

981    function, MATLAB R2018a). This step removed 7% of the response time data.

**REFERENCES**

Agrawal A, Hari KVS, Arun SP (2019) Reading Increases the Compositionality of Visual Word Representations. Psychol Sci 30:1707–1723.

Arun SP (2012) Turning visual search time on its head. Vision Res 74:86–92.

Baeck A, Kravitz D, Baker C, Op de Beeck HP (2015) Influence of lexical status and orthographic similarity on the multi-voxel response of the visual word form area. Neuroimage 111:321–328.

Bao P, Tsao DY (2018) Representation of multiple objects in macaque category-selective areas. Nat Commun 9:1774.

Brainard DH (1997) The Psychophysics Toolbox. Spat Vis 10:433–436.

Chanceaux M, Grainger J (2012) Serial position effects in the identification of letters, digits, symbols, and shapes in peripheral vision. Acta Psychol (Amst) 141:149–158.

Davis CJ (2010) The spatial coding model of visual word identification. Psychol Rev 117:713–758.

Dehaene S, Cohen L, Morais J, Kolinsky R (2015) Illiterate to literate: behavioural and cerebral changes induced by reading acquisition. Nat Rev Neurosci 16:234–244.

Dehaene S, Cohen L, Sigman M, Vinckier F (2005) The neural code for written words: a proposal. Trends Cogn Sci 9:335–341.

Dehaene S, Pegado F, Braga LW, Ventura P, Nunes Filho G, Jobert A, Dehaene-Lambertz G, Kolinsky R, Morais J, Cohen L (2010) How learning to read changes the cortical networks for vision and language. Science (80- ) 330:1359–1364.

Dufau S, Grainger J, Ziegler JC (2012) How to say "No" to a nonword: A leaky competing accumulator model of lexical decision. J Exp Psychol Learn Mem Cogn 38:1117–1128.

Eickhoff SB, Stephan KE, Mohlberg H, Grefkes C, Fink GR, Amunts K, Zilles K (2005) A new SPM toolbox for combining probabilistic cytoarchitectonic maps and functional imaging data. Neuroimage 25:1325–1335.

Ghose GM, Maunsell JH (2008) Spatial summation can explain the attentional modulation of neuronal responses to multiple stimuli in area v4. J Neurosci 28:5115–5126.

Gomez P, Ratcliff R, Perea M (2008) The overlap model: a model of letter position coding. Psychol Rev 115:577–600.

Grainger J (2018) Orthographic processing: A 'mid-level' vision of reading: The 44th Sir Frederic Bartlett Lecture. Q J Exp Psychol 71:335–359.

Grainger J, Dufau S, Montant M, Ziegler JC, Fagot J (2012) Orthographic processing in baboons (Papio papio). Science (80- ) 336:245–248.

Grainger J, Whitney C (2004) Does the huamn mnid raed wrods as a wlohe? Trends Cogn Sci 8:58–59.

Kay KN, Rokem A, Winawer J, Dougherty RF, Wandell BA (2013) GLMdenoise: A fast, automated technique for denoising task-based fMRI data. Front Neurosci 7:1–15.

Kronbichler M, Hutzler F, Wimmer H, Mair A, Staffen W, Ladurner G (2004) The visual word form area and the frequency with which words are encountered: evidence from a parametric fMRI study. Neuroimage 21:946–953.

Lehky SR, Tanaka K (2016) Neural representation for object recognition in inferotemporal cortex. Curr Opin Neurobiol 37:23–35.

Marcet A, Perea M (2017) Is nevtral NEUTRAL? Visual similarity effects in the early

1032    phases of written-word recognition. Psychon Bull Rev 24:1180–1185.
1033 McClelland JL, Rumelhart DE (1981) An interactive activation model of context
1034    effects in letter perception: I. An account of basic findings. Psychol Rev 88:375–
1035    407.
1036 Mruczek REB, Sheinberg DL (2005) Distractor familiarity leads to more efficient
1037    visual search for complex stimuli. Percept Psychophys 67:1016–1031.
1038 Nelson MJ, El Karoui I, Giber K, Yang X, Cohen L, Koopman H, Cash SS, Naccache
1039    L, Hale JT, Pallier C, Dehaene S (2017) Neurophysiological dynamics of
1040    phrase-structure building during sentence processing. Proc Natl Acad
1041    Sci:201701590.
1042 Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N (2014) A
1043    toolbox for representational similarity analysis. PLoS Comput Biol 10:e1003553.
1044 Norris D (2013) Models of visual word recognition. Trends Cogn Sci 17:517–524.
1045 Norris D, Kinoshita S (2012) Reading through a noisy channel: why there's nothing
1046    special about the perception of orthography. Psychol Rev 119:517–545.
1047 Op de Beeck H, Wagemans J, Vogels R (2001) Inferotemporal neurons represent
1048    low-dimensional configurations of parameterized shapes. Nat Neurosci 4:1244–
1049    1252.
1050 Pallier C, Devauchelle A-D, Dehaene S (2011) Cortical representation of the
1051    constituent structure of sentences. Proc Natl Acad Sci U S A 108:2522–2527.
1052 Pelli DG, Tillman K a (2008) The uncrowded window of object recognition. Nat
1053    Neurosci 11:1129–1135.
1054 Pelli DG, Tillman KA (2007) Parts, wholes and context in reading: A triple
1055    dissociation. PLoS One 2:e680.
1056 Pennington J, Socher R, Manning CD (2014) GloVe: Global Vectors for Word
1057    Representation. In: Empirical Methods in Natural Language Processing
1058    (EMNLP), pp 1532–1543.
1059 Perea M, Duñabeitia JA, Carreiras M (2008) R34D1NG W0RD5 W1TH NUMB3R5. J
1060    Exp Psychol Hum Percept Perform 34:237–241.
1061 Perea M, Panadero V (2014) Does viotin activate violin more than viocin? On the
1062    use of visual cues during visual-word recognition. Exp Psychol 61:23–29.
1063 Pramod RT, Arun SP (2014) Features in visual search combine linearly. J Vis 14:1–
1064    20.
1065 Pramod RT, Arun SP (2016) Object attributes combine additively in visual search. J
1066    Vis 16:8.
1067 Pramod RT, Arun SP (2018) Symmetric Objects Become Special in Perception
1068    Because of Generic Computations in Neurons. Psychol Sci 29:95–109.
1069 Ratan Murty NA, Arun SP (2015) Dynamics of 3D view invariance in monkey
1070    inferotemporal cortex. J Neurophysiol 113:2180–2194.
1071 Ratcliff R, Gomez P, McKoon G (2004) A diffusion model account of the lexical
1072    decision task. Psychol Rev 111:159–182.
1073 Ratcliff R, McKoon G (2008) The Diffusion Decision Model: Theory and Data for
1074    Two-Choice Decision Tasks. Neural Comput 20:873–922.
1075 Rawlinson GE (1976) The significance of letter position in word recognition.
1076 Rayner K (1998) Eye movements in reading and information processing: 20 years of
1077    research. Psychol Bull 124:372–422.
1078 Rayner K, White SJ, Johnson RL, Liversedge SP (2006) Raeding wrods with
1079    jubmled lettres: There is a cost. Psychol Sci 17:192–193.
1080 Rumelhart DE, McClelland JL (1982) An interactive activation model of context
1081    effects in letter perception: Part 2. The contextual enhancement effect and some

1082    tests and extensions of the model. Psychol Rev 89:60–94.

1083    Scaltritti M, Dufau S, Grainger J (2018) Stimulus orientation and the first-letter
1084    advantage. Acta Psychol (Amst) 183:37–42.

1085    Sripati AP, Olson CR (2010a) Global Image Dissimilarity in Macaque Inferotemporal
1086    Cortex Predicts Human Visual Search Efficiency. J Neurosci 30:1258–1269.

1087    Sripati AP, Olson CR (2010b) Responses to compound objects in monkey
1088    inferotemporal cortex: the whole is equal to the sum of the discrete parts. J
1089    Neurosci 30:7948–7960.

1090    Sunder S, Arun SP (2016) Look before you seek: Preview adds a fixed benefit to all
1091    searches. J Vis 16:3.

1092    Vighneshvel T, Arun SP (2015) Coding of relative size in monkey inferotemporal
1093    cortex. J Neurophysiol 113:2173–2179.

1094    Vinckier F, Dehaene S, Jobert A, Dubus JP, Sigman M, Cohen L (2007) Hierarchical
1095    Coding of Letter Strings in the Ventral Stream: Dissecting the Inner Organization
1096    of the Visual Word-Form System. Neuron 55:143–156.

1097    Yap MJ, Sibley DE, Balota DA, Ratcliff R, Rueckl J (2015) Responding to nonwords
1098    in the lexical decision task: Insights from the english Lexicon project. J Exp
1099    Psychol Learn Mem Cogn 41:597–613.

1100    Yarkoni T, Balota D, Yap M (2008) Moving beyond Coltheart's N: a new measure of
1101    orthographic similarity. Psychon Bull Rev 15:971–979.

1102    Zhivago KA, Arun SP (2014) Texture discriminability in monkey inferotemporal cortex
1103    predicts human texture perception. J Neurophysiol 112:2745–2755.

1104    Ziegler JC, Hannagan T, Dufau S, Montant M, Fagot J, Grainger J (2013)
1105    Transposed-Letter Effects Reveal Orthographic Processing in Baboons. Psychol
1106    Sci 24:1609–1611.

1107    Zoccolan D, Cox DD, DiCarlo JJ (2005) Multiple Object Response Normalization in
1108    Monkey Inferotemporal Cortex. J Neurosci 25:8150–8164.

1109    Zorzi M, Barbiero C, Facoetti A, Lonciari I, Carrozzi M, Montico M, Bravar L, George
1110    F, Pech-Georgel C, Ziegler JC (2012) Extra-large letter spacing improves
1111    reading in dyslexia. Proc Natl Acad Sci U S A 109:11455–11459.

1112

## ACKNOWLEDGEMENTS

1  **APPENDIX**

2

3  **For**

4

5  **A compositional letter code in high-level visual cortex**
6  **explains how we read jumbled words**

7

8

9  **CONTENTS**

19

20 ## SECTION A1. ADDITIONAL ANALYSIS FOR EXPERIMENT 1
21
22    The results in the main text were presented for uppercase English letters
23 (Figure 2), but in Experiment 1 we also collected visual search data for all pairs of
24 English letters and numbers (n = 62 characters in all, comprising 26 uppercase + 26
25 lowercase + 10 numbers). We did so in order to predict the visual dissimilarity between
26 letter strings containing both mixed case letters as well as numbers.
27    To visualize the dissimilarity relations between the 62 characters used, we
28 performed multidimensional scaling. In the resulting plot (Figure S1A), nearby
29 characters represent hard searches. A number of interesting patterns can be seen:
30 letters like C, G, Q, O are nearby which is expected given their shared curvatures.
31 Letter pairs such as (M,W) and number pairs such as (6,9) are similar due to mirror
32 confusion (Vighneshvel and Arun, 2013).
33    Next, we investigated the degree to which the observed pairwise dissimilarities
34 are captured by the multidimensional embedding as a function of the number of
35 dimensions. In the resulting plot (Figure S1B), it can be seen that nearly 89% of the
36 variance is captured by 10 dimensions as before, which reaches roughly the reliability
37 of the dissimilarity data itself. For the analyses involving mixed case searches or fewer
38 searches, we took a total of 6 neurons for the letter model, which explain 87.7% of the
39 variance in the pairwise dissimilarities.
40
41



42
43 **Figure S1. Visual search space for letters and digits**
44    (A) Visual search space for letters (uppercase and lowercase) and digits obtained
45       by multidimensional scaling of observed dissimilarities. Nearby letters
46       represent hard searches. Distances in this 2D plot are highly correlated with the
47       observed distances (r = 0.79, p < 0.00005).
48    (B) Correlation between observed distances and MDS embedding as a function of
49       number of MDS dimensions. The horizontal line represents the split-half
50       correlation with error bars representing s.d calculated across 100 random splits.
51
52 **Can letter dissimilarity be predicted using low-level visual features?**
53    To investigate whether single letter dissimilarity can be predicted using low-
54 level visual features, we attempted to predict letter dissimilarities using two models. In
55 the first model, which we call the pixel model, we calculated the dissimilarity between
56 letters to be the absolute difference in pixel intensities between the images of the two

57    letters. This pixel-based model showed a significant correlation (r = 0.50, p < 0.00005)
58    but was far from the reliability of the data itself ($r_{sh}$ = 0.90; Figure S1B). In the second
59    model, we calculated the dissimilarity between two letters as the vector distance
60    between the responses evoked by a population of simulated V1 neurons (Ratan Murty
61    and Arun, 2015). This V1 model also showed a significant correlation (r = 0.44, p <
62    0.00005) but again far from the reliability of the data itself). We conclude that single
63    letter dissimilarity can only be partially predicted by low-level visual features.
64
65    **Is visual search dissimilarity related to subjective dissimilarity?**
66          In this study, we have used visual search as a natural and objective measure
67    for visual dissimilarity. However previous studies have measured letter dissimilarity
68    either through confusions in letter recognition, or through subjective dissimilarity
69    ratings (Mueller and Weidemann, 2012; Simpson et al., 2013). We have previously
70    shown that subjective dissimilarity for abstract silhouettes is strongly correlated with
71    visual search dissimilarity (Pramod and Arun, 2016). This may not hold for letters since
72    subjects can activate letter representations that are modified through extensive
73    familiarity. To investigate how visual search dissimilarity compares with subjective
74    similarity ratings for letters, we compared search dissimilarities for uppercase letters
75    against two sets of previously reported similarity data. First, we compared visual
76    search dissimilarities with subjective dissimilarity ratings (Simpson et al., 2013). This
77    revealed a significant positive correlation (r = 0.69, p < 0.0005). Second, we compared
78    visual search dissimilarities with letter confusion data (*3*). To convert letter confusion
79    response times, which are a measure of similarity, into dissimilarities, we took their
80    reciprocals, and then compared them with visual search dissimilarities. This revealed
81    a significant positive, albeit weaker correlation (r = 0.34, p< 0.0005).
82

83        **SECTION A2. UPRIGHT AND INVERTED BIGRAMS AND TRIGRAMS**

84

85        It has been observed that readers are more sensitive to letter transpositions for
86    letters of their familiar script. Since discrimination of letter transpositions in the letter
87    model is a direct consequence of asymmetric spatial summation (main text, Figure 3),
88    we predicted that readers should show more asymmetric spatial summation for familiar
89    letters compared to unfamiliar letters. As a strong test of this prediction, we compared
90    visual search performance on upright letters (which are highly familiar) with inverted
91    letters (which are unfamiliar) across two experiments, one on bigrams and the other
92    on trigrams.

93        The comparison of upright and inverted letter strings is also interesting for a
94    second reason. If reading or familiarity with upright letters led to the formation of
95    specialized detectors for longer strings, then we predict that the letter model (which
96    assumes responses to be driven by single letters only) should yield worse fits for
97    upright compared to inverted letters.

98        We tested the above two predictions in the following two experiments.

99

100   **Experiment 3: Upright vs inverted bigrams**

101

102   *Methods.* A total of 8 subjects (6 males, aged 24 ± 1.5 years) participated in this
103   experiment. Six uppercase letters: A, L, N, R, S, and T were combined in all pairs to
104   form a total of 36 stimuli. These uppercase letters were chosen because their images
105   change when inverted (as opposed to letters like H that are unaffected by inversion),
106   and were chosen to maximize the occurrence of frequent bigrams. The same stimuli
107   were inverted to create another set of 36 stimuli. Stimuli subtended ~4° along the
108   longer dimension. Subjects performed all possible searches among the upright letters
109   ($^{36}C_2$ = 630 searches) with two repetitions and likewise for inverted letters. All trials
110   were interleaved. All other details were exactly as in Experiment 2.

111

112   **Results**

113       We observed interesting differences in search difficulty depending on the nature
114   of the bigrams. This pattern is illustrated in Figure S2A-B. When the target and
115   distractors consisted of repeated letters (e.g. TT among AA in Figure S2A), search is
116   equally easy when the array is upright or inverted. In contrast if the target and
117   distractors are transposed versions of each other (e.g. TA among AT in Figure S2B),
118   search is easier in the upright array compared to when it is inverted.

119       To confirm that this effect is present across all such pairs, we compared
120   observed response times for these two types of searches between upright and
121   inverted conditions (Figure S2C). Response times for the AA-BB searches were
122   comparable for upright and inverted conditions (mean ± sd of RT: 0.66 ± 0.09 s for
123   upright, 0.67 ± 0.1 s for inverted). To assess the statistical significance of this
124   difference, we performed an ANOVA with subject (8 levels), bigram (15 pairs) and
125   orientation (upright vs inverted) as factors. We observed no significant difference in
126   the response times between upright and inverted conditions for AA-BB searches ($p$ =
127   0.65 for main effect of orientation; $p < 0.00005$ for subject and bigram factors, $p > 0.05$
128   for all interactions).

129       Next we compared transposed letter (AB-BA) searches. Here, subjects were
130   clearly faster on the upright searches compared to inverted searches (mean ± sd of
131   RT: 1.58 ± 0.25 s for upright, 3.12 ± 0.76 s for inverted). This difference was statistically
132   significant ($p < 0.00005$ for main effect of orientation; $p < 0.0005$ for subject and $p <$

133  .05 for bigram factors, p < 0.05 for interactions between pairs and orientation. Other
134  interaction effects were not significant).
135       To compare bigram dissimilarity between upright and inverted bigrams, we
136  plotted one against the other. This revealed a highly significant correlation (r = 0.80, p
137  < 0.00005; Figure S2D). Here too it can be seen that the transposed letter searches
138  are clearly faster when they are upright whereas the repeated letter searches show no
139  such difference.
140       Thus, inversion slows down transposed letter searches but not repeated letter
141  searches.
142
143  **Explaining upright and inverted bigram dissimilarity using the letter model**
144       We fit the letter model to both upright and inverted bigram searches using a
145  total of 10 neurons with single letter responses derived from Experiment 1. The letter
146  model yielded excellent fits on both upright and inverted bigrams. In both cases, the
147  model fits approached the data consistency (Figure S2E), implying that the model
148  explained nearly all the explainable variance in the data.
149       To compare these model fits for upright vs inverted statistically, we performed
150  a bootstrap analysis. Each time, we selected subjects with replacement and fit the
151  letter model to the average dissimilarity computed for this random pool of subjects.
152  Each time we calculated a normalized correlation measure that takes into account the
153  difference in data reliability between upright and inverted trigram searches. This
154  normalized correlation is simply the model correlation divided by the data consistency.
155  To assess statistical significance, we calculated the fraction of times the normalized
156  correlation in the upright samples was larger than the inverted samples. This analysis
157  revealed significant difference in model performance between upright and inverted
158  searches, but in the opposite direction (average model correlation: r = 0.92 for upright,
159  0.9 for inverted; fraction of upright < inverted normalized model correlation: p = 0).
160  Thus, upright searches are more predictable than inverted searches using the letter
161  model.
162       Next we asked whether the letter model can explain the intriguing observation
163  that inversion affects transposed letter searches but not repeated letter searches. This
164  is easy to explain in the letter model: The response to repeated letter bigrams such as
165  AA is unaltered (Figure 3B), and therefore the dissimilarity between AA and TT is
166  unaffected by the asymmetry in spatial summation. By contrast, the dissimilarity
167  between transposed letter pairs like AT & TA is directly driven by the asymmetry in
168  spatial summation. We also note that the search TT among AA is much easier than
169  the search for TA among AT. This is also explained by the letter model by the fact that
170  the response to repeated letters is the same as the response to individual letters,
171  leaving their discrimination unaltered. By contrast transposed letters are much more
172  similar since their neural responses are much closer (Figure 3B).
173       To be sure that letter model predictions show the same pattern, we plotted the
174  average response time predicted by the letter model for repeated letter (AA-BB) and
175  transposed letter (AB-BA) searches. To assess the statistical significance, we
176  performed a sign-rank test on the predicted RT. The letter model predictions were
177  exactly as expected (Figure S2F).
178       Next we analysed the model parameters in the letter model to ascertain whether
179  the spatial summation in the neurons was indeed different for upright and inverted
180  bigrams. To quantify the degree of asymmetry, we calculated for each neuron a spatial
181  modulation index of the form $MI = abs(w1-w2)/(w1+w2)$ where w1 and w2 are the
182  estimated weights for each letter in the bigram. To avoid unnaturally large modulation

183   indices, w1 and w2 values smaller than 0.01 were set to 0.01. The spatial modulation
184   index for all 10 neurons for upright and inverted bigrams is shown in Figure S2G. It
185   can be seen that the modulation index is larger in most cases for the upright bigrams.
186   This difference was statistically significant, as assessed using a sign-rank test on the
187   spatial modulation indices (Figure S2H).
188
189



190
191   **Figure S2. Letter model fits for upright and inverted bigrams**
192   (A) Example oddball search array for a repeated letter target (TT) among identical
193        repeated-letter distractors (AA). It can be seen that inverting this search array
194        does not affect search difficulty.
195   (B) Example oddball search array for transposed letters (TA among AT). It can be
196        seen by inverting this search array makes the search substantially more
197        difficult.
198   (C) Average search times in the oddball search task for repeated-letter searches
199        (AA-BB) and transposed letter (AB-BA) searches. Error bars represent s.e.m
200        calculated across subjects. Asterisks represent statistical significance (**** is p
201        < 0.00005), as obtained using an ANOVA on the response times with subject,
202        bigram and orientation as factors (see text).
203   (D) Dissimilarity of inverted bigram pairs plotted against the dissimilarity of upright
204        bigram pairs. Correlation is shown at the top left. Asterisks indicate statistical
205        significance of the correlations (**** is p < 0.00005).
206   (E) Cross-validated model correlation of the letter model for upright bigrams and
207        inverted bigrams. *Shaded gray bars* represent the upper bound achievable in
208        each case given the consistency of the data, calculated using the split-half
209        correlation $r_{sh}$.
210   (F) Predicted RT from the letter model for repeated letter pairs and transposed
211        letter pairs. Asterisks denote statistical significance as obtained using a sign-
212        rank test on the predicted RTs between upright and inverted conditions.
213   (G) Spatial modulation index for each neuron in the letter model for upright and
214        inverted bigrams.

215     (H) Average spatial modulation index for upright and inverted bigrams. Asterisks
216           represent statistical significance (* is $p < 0.05$) obtained using a sign-rank test
217           on the spatial modulation index across the 10 neurons.
218
219 **Experiment S1: Upright and inverted trigrams**
220     Here, we asked whether the above results would extend to trigrams. We tested
221 two predictions. First, we predicted greater spatial modulation for upright compared to
222 inverted trigrams, on the premise that better discrimination of trigram transpositions
223 should be driven by asymmetric spatial summation. Second, if repeated viewing of a
224 trigram or word led to the formation of specialized trigram detectors, then the letter
225 model (which is based only on knowledge of single letters) should produce larger
226 errors compared to other trigrams. We tested this prediction by comparing model fits
227 for searches involving frequent trigrams and words compared to other searches.
228
229 *Methods.* A total of 9 subjects (6 females, aged 24.5 ± 2.3 years) participated in the
230 experiment. Six uppercase letters: A, G, N, R, T and Y were combined in all possible
231 3-letter combination to form a total of 216 stimuli. These letters were chosen to include
232 as many three-letter words as possible. In all, 15 three-letter words could be created
233 using these letters (ANT, ANY, ART, GAG, GAY, NAG, NAY, RAG, RAN, RAT, RAY,
234 TAG, TAN, TAR, and TRY).
235     Since the total number of possible search pairs is large ($^{216}C_2$ = 23,220 pairs),
236 we chose 500 search pairs such that the regression matrix of the part-sum model had
237 full rank i.e. all the model parameters can be estimated reliably using linear regression.
238 These 500 searches consisted of 368 random search pairs, 105 ($^{15}C_2$) word-word
239 pairs, 15 ($^{3!}C_2$) transposed pairs of nonword comprised of letters G,N, and R. Further,
240 another set of 15 ($^{3!}C_2$) transposed pairs were created using the word TAR. The search
241 pairs formed using the words TAR, ART and RAT were presented only once (although
242 they were counted as both word-word pairs and transposed pairs in the main analysis).
243     Subjects performed the same searches using upright and inverted trigrams.
244 Stimuli subtended ~5° along the longer dimension. All subjects completed 2000
245 correct trials (500 searches x 2 orientations x 2 repetitions). All other details were
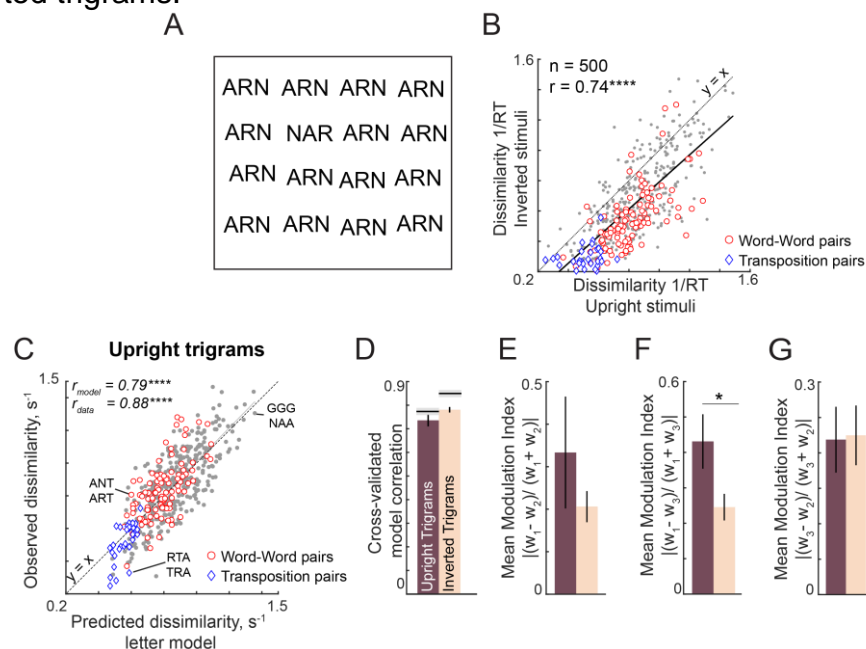246 identical to Experiment 1.
247
248 **Results**
249     An example oddball array in the trigram experiment is shown in Figure S3A.
250 Note that it is no longer meaningful to compare repeated letter trigrams (AAA-BBB)
251 with transposed trigrams (ABC-BCA) because the repeated letter pairs contain two
252 unique letters whereas the transposed trigrams contain three unique letters. Subjects
253 were highly consistent in both upright and inverted searches (split-half correlation
254 between even and odd- subjects: r = 0.76 & 0.80, $p < 0.00005$). Upright and inverted
255 dissimilarities were highly correlated (r = 0.80, $p < 0.00005$; Figure S3B), although
256 upright searches had higher dissimilarity compared to inverted searches.
257     Next we asked whether the letter model can predict dissimilarities between
258 upright trigrams. As before, letter model predictions were highly correlated with the
259 observed data (r = 0.79, $p < 0.00005$; Figure S3C) and this model fit approached the
260 data consistency itself ($r_{data}$ = 0.88). Model fits errors were acctually lower for
261 transposed pairs compared to word-word pairs and other pairs (mean ± sd error: 0.1
262 ± 0.08 for word pairs; 0.07 ± 0.06 for transposed pairs; 0.11 ± 0.08 for other pairs; p =
263 0.02, rank-sum test). The letter model was also able to predict dissimilarities between
264 various trigram transpositions (r = 0.69, $p < 0.00005$; Figure S3C). Thus, trigram

265 dissimilarities can be predicted by the letter model regardless of word status or trigram
266 frequency.
267      We then compared model fits for upright and inverted trigrams. In both cases,
268 the letter model predictions (r = 0.78 & 0.73 for upright and inverted) were close to the
269 consistency of the data ($r_{data}$ = 0.85 & 0.78; Figure S3D). To compare these model fits
270 for upright vs inverted statistically, we performed a bootstrap analysis as before
271 (Experiment 3). This analysis revealed no significant difference in model performance
272 between upright and inverted searches (fraction of upright < inverted normalized
273 model correlation: p = 0.07).
274      Finally we asked whether the spatial summation weights of the letter model
275 were systematically different between upright and inverted trigrams. Since there are
276 three spatial modulation weights for each neuron, we calculated the spatial modulation
277 index for all possible pairs of weights (Figure S3 E,F,G). The spatial modulation ratio
278 was larger for upright compared to inverted trigrams in two of the three pairs, and this
279 difference attained statistical significance for the first and third letters in the trigram
280 (Figure S3F). We conclude that the spatial modulation is stronger for upright compared
281 to inverted trigrams.



282
283 **Figure S3. Letter model fits for upright and inverted trigrams**
284      (A) Example trigram search array containing letter transpositions, with oddball
285          target (NAR) among distractors (ARN). It can be seen that this search is
286          substantially harder when inverted compared to upright.
287      (B) Dissimilarity for inverted trigram searches (1/RT) plotted against dissimilarity for
288          upright trigram searches for word-word pairs (red circles, n = 105), transposed
289          letter pairs (blue diamonds, n = 30), and other pairs (gray circles, n = 365).
290      (C) Observed dissimilarity for upright trigrams plotted against the predicted
291          dissimilarity from the letter model with symbol conventions as in (B).
292      (D) Cross-validated letter model correlation for upright and inverted trigrams.
293      (E) Average spatial modulation index (across 10 neurons) for the first and second
294          letters in the trigram.
295      (F) Same as (E) but for the first and third letters.
296      (G) Same as (E) but for the second and third letters.

## SECTION A3: COMPOUND WORDS

Here we created compound words by combining two valid words such as FORGET from FOR and GET (Figure S5A). This resulted in some valid words (e.g. FORGET, TEAPOT) and many invalid words (e.g. FORPOT and TEAGET). The full stimulus set is shown in Figure S4.

If valid words are driven by specialized detectors, responses to valid words should be less predictable by the single letter model. We formulated two specific predictions. First, we hypothesize that the dissimilarity between valid words (e.g. FORMAT vs TEAPOT) would yield larger model errors compared to invalid word pairs (e.g. DAYFOR vs ANYMAT). Second, we predicted that the dissimilarity between two invalid compound words (e.g. DAYFOR vs ANYMAT) should be explained better by their constituent trigrams (DAY, FOR, ANY, MAT) rather than by their constituent letters (Figure S5B).


**METHODS**

A total of 8 subjects (4 female, aged 25 ± 2.5 years) participated in the experiment. Twelve 3-letter words were chosen: ANY, FOR, TAR, KEY, SUN, TEA, ONE, MAT, GET, PAD, DAY, POT. Each word was scrambled to obtain twelve 3-letter nonwords containing the same letters. The 12 words were combined to form 36 compound words (Figure S4), such that they appeared equally on the left and right half of the compound words. It can be seen that there are seven valid words, whereas the other compound words are pseudowords that carry no meaning. The compound words measured 6° along the longer dimension. Subjects completed 1260 correct trials ($^{36}C_2$ search pairs x 2 repetitions). Additionally, subjects also performed visual search on 3-letter words (n = 132, $^{12}C_2$ x 2 repetitions) and their jumbled versions (n = 132). Trials timed out after 15 seconds. All other details were identical to Experiment 1.

Subjects were highly accurate on this task (mean ± std: 98 ± 1%). Outliers in the reaction times were removed using built-in routines in MATLAB (isoutlier function, MATLAB R2018a). This step removed 6.4% of the response time data.

|  | **ANY** | **FOR** | **TAR** | **KEY** | **SUN** | **TEA** |
|---|---|---|---|---|---|---|
| **ONE** | ANYONE | ONEFOR | ONETAR | KEYONE | ONESUN | TEAONE |
| **MAT** | MATANY | FORMAT | MATTAR | MATKEY | SUNMAT | TEAMAT |
| **GET** | GETANY | FORGET | TARGET | KEYGET | GETSUN | GETTEA |
| **PAD** | PADANY | FORPAD | TARPAD | KEYPAD | PADSUN | PADTEA |
| **DAY** | ANYDAY | DAYFOR | TARDAY | DAYKEY | SUNDAY | DAYTEA |
| **POT** | ANYPOT | POTFOR | POTTAR | POTKEY | SUNPOT | TEAPOT |

**Figure S4. Stimulus set used for Experiment 4 (Compound Words).** The left and the right 3 letters words were combined to form a 6-letter string. The strings that formed compound words are highlighted in red.


**RESULTS**

We recruited 8 subjects to perform oddball search involving pairs of trigrams as well as 6-letter strings. In all there were 12 three-letter words which resulted in $^{12}C_2 = 66$ searches and 36 compound 6-letter strings which resulted in $^{36}C_2 = 630$ searches. We also included 12 three-letter nonwords created by transposing each three-letter

340  words, resulting in an additional $^{12}C_2 = 66$ searches. As before, subjects were highly
341  consistent in their responses (split-half correlation between odd and even subjects: r
342  = 0.54, p < 0.00005 for 3-letter words; r = 0.46, p < 0.00005 for 3-letter nonwords; r =
343  0.65, p < 0.00005 for 6-letter words).
344      We started by using the single letter model as before to predict compound word
345  responses. We took single neuron responses as before from Experiment 1, and took
346  the response of each neuron to a compound word to be a weighted sum of its
347  responses to the individual letters. Using these compound word responses, we
348  calculated the dissimilarity between pairs of compound words, and used nonlinear
349  fitting to obtain the best model parameters. The single letter model yielded excellent
350  fits to the data (r = 0.68, p < 0.00005; Figure S5C). This performance was comparable
351  to the data consistency estimated as before ($r_{data}$ = 0.72).
352      Next we asked whether discrimination between compound words can be
353  explained better as a combination of two valid three-letter words, or as a combination
354  of all the constituent six letters. To address this question we constructed a new
355  compositional model based on trigrams, and asked if its performance was better than
356  the single letter model (Figure S5D). The trigram-based letter model used trigram
357  dissimilarity to construct neurons with trigram tuning, and spatial summation over the
358  two trigrams to predict the 6-gram responses. To compare the performance of both
359  models even though they have different numbers of free parameters, we used cross-
360  validation: we fit both models on half the subjects and tested their performance on the
361  other half. The letter model outperformed the trigram model (Figure S5D). Because
362  both models were trained on half the subjects and tested on the other half, the upper
363  bound on their performance is simply the split-half correlation between the two halves
364  of the data (denoted by $r_{sh}$). Indeed the letter model performance was close to this
365  upper bound ($r_{sh}$ = 0.56; Figure S5D), suggesting that it explained nearly all the
366  explainable variance in the data. Finally, the letter model outperformed a widely used
367  model for orthographic distance – the Orthographic Levenshtein Distance (OLD)
368  (Figure S5D). Thus, compound word discrimination can be understood from single
369  letters.
370      Finally, the letter model fits for word-word pairs and nonword-nonword pairs
371  were not significantly different (Figure S5E). This further validates the absence of local
372  combination detectors (Dehaene et al., 2005) in perception.
373
**Three-letter word and nonword dissimilarities**
375      To investigate whether the letter model can predict dissimilarities between
376  three-letter words and non-words, we fit a separate letter model with 6 neurons as
377  before to the word and non-word dissimilarities. If frequent viewing of words led to the
378  formation of specialized word detectors, the letter model would show worse model fits
379  compared to nonwords. However, we observed no such pattern: the letter model fits
380  were equivalent for words (r = 0.69, p < 0.00005; Figure S5F) and nonwords (r = 0.57,
381  p < 0.00005; Figure S5F) – and these fits approached the respective data
382  consistencies ($r_{data}$ = 0.67 for words, 0.68 for nonwords). We conclude that three-letter
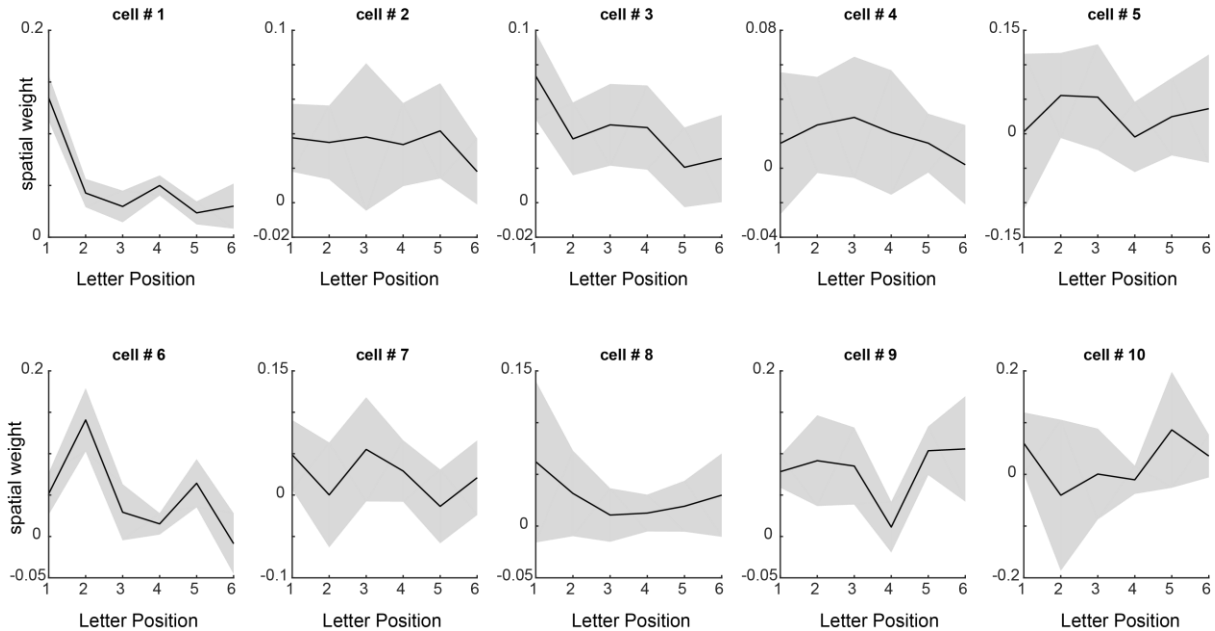383  string dissimilarities can be predicted by the letter model regardless of word status.
384

**Figure S5. Discrimination of compound words in visual search (Experiment 4).**

(A) 3-letter words (*top*) used to create compound words (*bottom*).

(B) Illustration of letter and trigram models. In the letter model, the response to a compound word is a weighted sum of responses to the six single letters. In the trigram model, the response to a compound word is a weighted sum of its two trigrams.

(C) Observed dissimilarity for compound words plotted against predicted dissimilarity from the letter model for word pairs (*red*) and other pairs (*gray*).

(D) Cross-validated model correlations for the letter model, trigram model and the Orthographic Levenshtein distance (OLD) model. The upper bound on model fits is the split-half correlation ($r_{sh}$), shown in black with shaded error bars representing standard deviation across 30 random splits. Horizontal lines above shaded error bar depicts significant difference across different models.

(E) Cross-validated model fits of the letter model for word-words pairs and nonword-nonword pairs.

(F) Observed dissimilarities for 3-letter words (*black*) and nonwords (*red*) plotted against letter model predictions.

**Spatial summation weights**

To investigate the spatial summation weights for each neuron, we plotted the estimated spatial summation weights separately (Figure S6). It can be seen that spatial summation is heterogeneous across neurons, but the spatial summation of the first neuron follows the characteristic W-shaped curve for letter position observed in studies of reading.



**Figure S6. Spatial summation weights for each neuron.** Estimated spatial summation weights (mean ± std across many random starting points of the nonlinear model fit algorithm) for each neuron in the letter model.

418 ## SECTION A4. EXPERIMENTS WITH LONGER STRINGS

419

420 In the main text, we showed that bigram dissimilarity in visual search can be
421 explained using a simple letter model with single letter responses that match
422 perception, and a compositional spatial summation rule that predicts responses to
423 bigrams. Here we asked whether this approach would generalize to longer strings of
424 letters.

425 To this end, we performed four additional experiments on longer strings. In
426 Experiment S2, we created trigrams with a fixed middle letter and all possible
427 combinations of flanking letters, to create multiple three-letter words. In Experiment
428 S3, subjects performed searches involving 3, 4, 5 and 6-letter searches with
429 uppercase, lowercase and mixed case strings. In Experiments S4 & S5, we attempted
430 to optimize the search pairs used to estimate model parameters.

431

432 **METHODS**

433

434 *Experiment S2: Trigrams with fixed middle letter.* A total of 8 subjects (5 males, aged
435 $23.9 \pm 1.8$ years) participated in this experiment. Seven uppercase letters: A, E, I, P,
436 S, T and Y were combined (around the stem R i.e. xRx) in all pairs to form a total of
437 49 stimuli. These letters were chosen to maximize the occurrence of 3-letter words
438 and psuedowords in the stimulus set. The longer dimension of the stimuli was ~5°.
439 Each subject completed searches corresponding to all possible pairs of stimuli ($^{49}C_2 =$
440 1176) with two trials for each search. All other details were identical to Experiment 2.

441

442 *Experiment S3: Random string searches.* A total of 12 subjects (9 female, aged 24.8
443 $\pm 1.64$ years) participated in this experiment. All 26 uppercase and lowercase letters
444 were used to create 1800 stimuli, which were organized into 900 stimulus pairs with
445 varying string length. These 900 pairs comprised 300 6-gram uppercase pairs, 100 6-
446 gram lowercase pairs, 100 6-gram mixed-case pairs, 100 5-gram uppercase pairs, 50
447 4-gram uppercase pairs, 50 3-gram uppercase pairs and 200 pairs with uppercase
448 strings of differing lengths (50 pairs each of 6- vs 5-grams, 6- vs 4-grams, 5- vs 4-
449 grams, 5- vs 3-grams = 200 pairs total). For each string length, letters were randomly
450 combined to form strings with a constraint that all 26 letters should appear at least
451 once at each location. Each stimulus pair was shown in two searches (with either item
452 as target, and either on the left or right side). The trial timed out at 15 seconds for all
453 searches.

454

455 *Experiment S4 – Optimized 4-letter searches.* In all, 8 subjects (5 females, aged 23.5
456 $\pm 2.3$ years) participated in this experiment. To maximize the importance of each
457 spatial location in a 4-letter uppercase string, stimuli were created such that there were
458 at least 75 search pairs with the same letter at either of the corresponding locations.
459 Further, to reliably estimate the model parameters, the randomly chosen letters were
460 arranged to minimize the condition number of the linear regression matrix X of the ISI
461 model described below. In all there were 300 search pairs. The trial timed out after 15
462 seconds. All other details were similar to Experiment 2.

463

464 *Experiment S5 – Optimized 6-letter searches.* A total of 9 subjects (5 males, aged 24.1
465 $\pm 2.2$ years) participated in this experiment. We chose 300 search pairs with 6-letter
466 strings, according to the same criteria as in Experiment S4. All other details were the
467 same as in Experiment S4.

468 **RESULTS**

469    Cross-validated model fits across all experiments are shown in Figure S7. It
470 can be seen that the letter model fit is close to the split-half consistency of the data.
471 Thus, visual discrimination of longer strings can be explained using a compositional
472 neural code. Below we discuss some experiment-specific findings of interest.

473

474 *Lowercase and mixed-case strings*

475    Word shape is thought to play a role in reading lowercase letters, because of
476 the upward deflection (e.g. l, d) and downward deflections (e.g. p, g) of letters which
477 might confer a specific overall shape to a word. To conclusively establish this would
478 require factoring out the contribution of individual letters to word discrimination, as with
479 the letter model. We were therefore particularly interested in whether the letter model
480 would predict the dissimilarity between lowercase and mixed-case strings where word
481 shape might potentially play a role. As can be seen in Figure S7, cross-validated model
482 predictions for lowercase letters were highly correlated with the observed data (r =
483 0.59, p < 0.00005). This correlation approached the upper bound given by the split-
484 half reliability itself ($r_{sh} = 0.64$). Likewise, model predictions for mixed-case letters were
485 also highly correlated with the observed data (r = 0.59, p < 0.00005; Figure S7).
486 However in this case model fits were well below the split-half consistency ($r_{sh} = 0.72$),
487 suggesting that there is still some systematic unexplained variance in mixed-case
488 strings. This gap in model fit could be simply due to the relatively few mixed-case
489 searches used in this experiment (n = 100), or because of unaccounted factors like
490 word shape. Nonetheless, the letter model explains a substantial fraction of variation
491 in both lowercase and mixed case strings, suggesting that it can be used as a powerful
492 baseline to elucidate the contribution of word shape to reading.

493

494 *Unequal length strings*

495    The letter model can be used to calculate responses to any string length,
496 provided the spatial summation weights are known. Given the relatively few searches
497 for unequal lengths in our data, we fit the letter model to unequal length strings using
498 6 neurons. Doing so still raised a fundamental issue: which subset of the 6 spatial
499 summation weights for each neuron should be used to calculate the response to a 4-
500 letter string? This requires aligning the 4-letter string to the 6-letter string in some
501 manner.

502    To address this issue, we evaluated the letter model fit on four possible
503 alignments between longer and shorter strings, and asked whether model predictions
504 were better for any one alignment compared to others. We aligned the smaller length
505 string to either the left, right, centre or edge of the longer string. Model performance
506 for these different variations is shown in Table S1. It can be seen that the model fits
507 are comparable across different choices. However, edge alignment is slightly but not
508 significantly better than other choices. We therefore used edge alignment for all
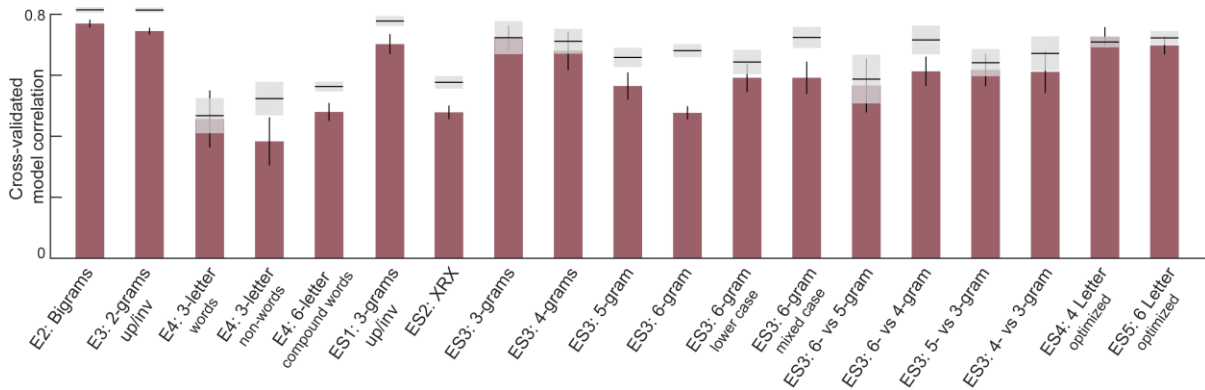509 subsequent model predictions.

510

511

| Alignment | Letter model correlation | | | |
|---|---|---|---|---|
| | 6 vs 5 | 6 vs 4 | 5 vs 3 | 4 vs 3 |
| Left: ABCDEF vs EFGHxx | 0.54 | 0.66 | 0.58 | 0.57 |
| Right: ABCDEF vs xxEFGH | 0.51 | 0.66 | 0.57 | 0.58 |
| Centre: ABCDEF vs xEFGHx | - | **0.68** | 0.58 | - |
| Edge: ABCDEF vs EFxxGH | **0.55** | 0.63 | **0.60** | **0.59** |

512 **Table S1: Model fits for various choices of string alignment.** In each case we fit
513 the letter model with unknown weights corresponding to the longer length. The
514 alignment is indicated by the position of "x"s in the string. For instance, "Left" alignment
515 means that a 6-letter string ABCDEF is matched to a 4-letter string EFGH by assuming
516 that the response to EFGH is created using the first four weights of spatial summation.
517 Likewise, right alignment means that EFGH is aligned to the right, and therefore its
518 response is created using the last four weights in the 6-letter letter model. The best
519 alignment is highlighted for each column in **bold**. None of the correlation coefficient
520 differences were statistically significant (p > 0.05, Fisher's z-test).

521
522



523
524 **Figure S7. Letter model performance for varying length strings.** For each
525 experiment, we obtained a cross-validated measure of model performance using 6
526 neurons as follows: each time we divided the subjects randomly into two halves, and
527 trained the letter model on one half of the subjects and tested it on the other half. This
528 was repeated for 30 random splits. The correlation between the model predictions and
529 the average dissimilarity from the held-out half of the data was taken to be the model
530 fit. The correlation between the observed dissimilarity between the two random splits
531 of subjects is then the upper bound on model performance (mean ± std shown as *gray*
532 *shaded bars*).

533
534

## SECTION A5. ESTIMATING LETTER DISSIMILARITY FROM BIGRAMS

**Part-sum model**

The letter model described in the text has many desirable features but requires as input the responses to single letters, which were obtained from searches involving single isolated letters. However, it could be that bigram representations can be understood in terms of component letter responses that are different from the responses of letters seen in isolation. It could also be that letter responses are different at each location.

To address these issues, we developed an alternate model in which bigram dissimilarities can be written in terms of unknown single letter dissimilarities. These single letter dissimilarities can be estimated in the model. In this model, which we call the part-sum model, the dissimilarity between two bigrams AB & CD is written as the sum of all pairs of part dissimilarities in the two bigrams (Figure S8A). Specifically:

$$d(AB,CD) = CL_{AC} + CR_{BD} + X_{AD} + X_{BC} + W_{AB} + W_{CD} + constant$$

where $CL_{AC}$ is the dissimilarity between letters at Corresponding Left (CL) locations (A & C), $CR_{BD}$ is the dissimilarity between letters at the Corresponding Right (CR) locations (B & D), $X_{AD}$ & $X_{BC}$ are the dissimilarities between letters across locations in the two bigrams (A & D, B & C), and $W_{AB}$ & $W_{CD}$ are the dissimilarities of letters within each bigram.

The part-sum model works because a given letter dissimilarity $CL_{AC}$ will occur in the dissimilarity of many bigram pairs (e.g. in the pair AB-CD and in AE-CF) thereby allowing us to estimate its unique contribution. Since there are 7 parts, there are $^7C_2$ = 21 possible part-pairs of each type (i.e. for CL, CR, X and W terms), resulting in 21 x 4 = 84 unknown part dissimilarities. Since a given bigram experiment contains all possible $^{49}C_2$ = 1176 bigram searches, there are many more observations than unknowns. The combined set of bigram dissimilarities can be written in the form of a matrix equation $\mathbf{y} = \mathbf{Xb}$ where $\mathbf{y}$ is a 1176x1 vector of observed bigram dissimilarities, $\mathbf{X}$ is a 1176 x 85 matrix containing the number of times (0, 1 or 2) a given letter-pair of each type (CL, CR, X & W) contributes to the overall dissimilarity, and b is a 85 x 1 vector of unknown letter dissimilarities of each type (21 each of CL, CR, X & W and one constant term). The unknown letter dissimilarities of each type was estimated using standard linear regression (*regress* function, MATLAB).

The part sum model has several advantages over the letter model: (1) It is linear which means that its parameters can be uniquely estimated; (2) it is compositional in that the net dissimilarity between two bigrams is explained using the constituent parts without invoking more complex interactions; (3) it can account for potentially different part relations at each location in the two bigrams. We have previously shown that the part-sum model can explain the dissimilarities between a variety of objects (Pramod and Arun, 2016).

The part sum model yielded excellent fits to the data (r = 0.88, p < 0.00005; Figure S8B) that were close to the reliability of the data ($r_{data}$ = 0.90). As before, we observed no systematic deviations between model fits for frequent bigrams compared to infrequent bigrams (Figure S8B; average absolute residual error for the top 20 bigram pairs with highest mean bigram frequency: 0.09 ± 0.1 s$^{-1}$; for the bottom-20 bigram pairs: 0.11 ± 0.08 s$^{-1}$; p = 0.42, rank-sum test). To assess whether the part dissimilarities of each type (CL, CR, X and W) were related to each other, we plotted each of CR, X and W terms against the CL terms (Figure S8C). The CR and X terms

585  were highly positively correlated (Figure S8C), whereas the W terms were negative in
586  sign and negatively correlated (Figure S8C). The negative values of the W terms
587  means that bigrams with dissimilar letters become less dissimilar, an effect akin to
588  distractor heterogeneity in visual search (Duncan and Humphreys, 1989; Vighneshvel
589  and Arun, 2013). We conclude that the CL, CR, X and W terms in the part-sum model
590  are driven by a common part representation.
591          To visualize this underlying letter representation, we performed
592  multidimensional scaling on the estimated part dissimilarities of the CL terms. In the
593  resulting plot, nearby letters represent similar letters (Figure S8D). It can be seen that
594  I & T, M & N are similar as in the single-letter representation (Figure S1A). These
595  single letter dissimilarities estimated from bigrams using the part-sum model were
596  highly correlated with the single-letter dissimilarities directly observed from visual
597  search with isolated letters (Figure S8D).
598          We conclude that bigram dissimilarities can be predicted from a common
599  underlying letter representation that is identical to that of single isolated letters.
600
601  **Equivalence between part-sum and letter model**
602          Given that the part-sum model and letter model both give equivalent fits to the
603  data, we investigated how they are related. Consider a single neuron whose response
604  to a bigram AB is given by: $r_{AB} = \alpha r_A + r_B$, where $r_A$ and $r_B$ are its responses to A & B,
605  and $\alpha$ is the spatial weight of A relative to B. Similarly its response to the bigram CD
606  can be written as $r_{CD} = \alpha r_C + r_D$. Then the dissimilarity between AB and CD can be
607  written as
608
609  $d(AB, CD)^2$
610  $= (r_{AB} - r_{CD})^2 = (\alpha r_A + r_B - \alpha r_C - r_D)^2$
611  $= \left( \alpha(r_A - r_C) + (r_B - r_D) \right)^2$
612  $= \alpha^2 (r_A - r_C)^2 + (r_B - r_D)^2 + 2\alpha(r_A - r_C)(r_B - r_D)$
613  $= \alpha^2 (r_A - r_C)^2 + (r_B - r_D)^2 + 2\alpha(r_A r_B + r_C r_D - r_A r_D - r_B r_C)$
614  $= \alpha^2 (r_A - r_C)^2 + (r_B - r_D)^2 + \alpha[(r_A - r_D)^2 + (r_B - r_C)^2 - (r_A - r_B)^2 - (r_C - r_D)^2]$
615  $= \alpha^2 d_{AC}^2 + d_{BD}^2 + \alpha(d_{AD}^2 + d_{BC}^2 - d_{AB}^2 - d_{CD}^2)$
616  $= \alpha^2 d_{AC}^2 + d_{BD}^2 + \alpha(d_{AD}^2 + d_{BC}^2) - \alpha(d_{AB}^2 + d_{CD}^2)$
617
618  Thus, the squared dissimilarity between AB & CD can be written as a weighted sum
619  of squared dissimilarities between parts at corresponding locations (A-C & B-D), parts
620  at opposite locations (A-D & B-C) and between parts within each bigram (A-B & C-D),
621  which is essentially the same as the part-sum model. The same argument extends to
622  multiple neurons because the total bigram dissimilarity will be the sum of bigram
623  dissimilarities across all neurons.
624          There are however two important differences. First, the part sum model is
625  written in terms of a weighted sum of part dissimilarities, whereas the above equation
626  refers to a weighted sum of squared dissimilarities. However, the squared sum of
627  distances and a weighted sum of distances are highly correlated, so the essential
628  relation will still hold. Second, the letter model predicts that the across-bigram terms
629  ($X_{AD}$, $X_{BC}$) should be similar in magnitude but opposite in sign to the within-bigram
630  terms ($W_{AB}$, $W_{CD}$). These weights are similar in magnitude but not exactly equal, as
631  can be seen in Fig S8C. The part-sum model thus allows for greater flexibility in part
632  interactions compared to the letter model.
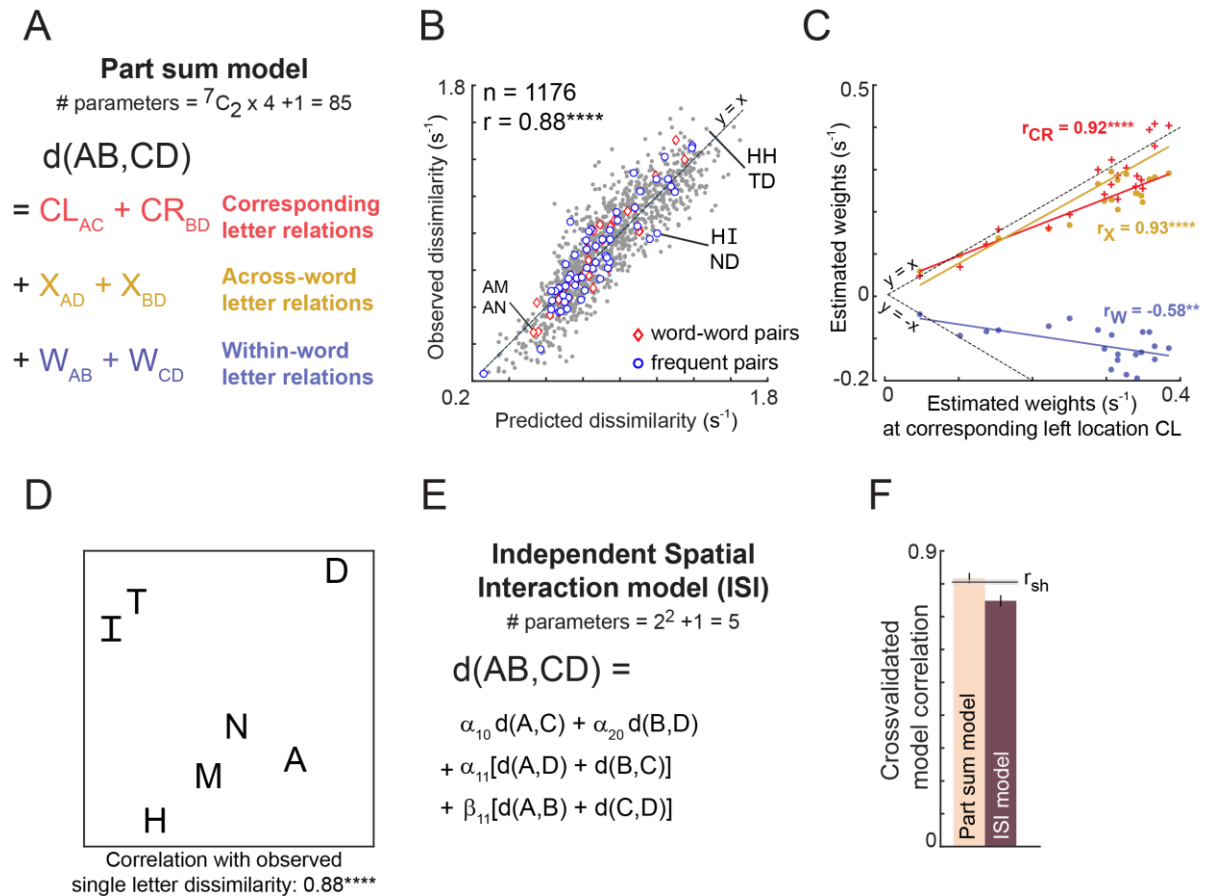
**Reducing part-sum model complexity (ISI model)**

The observation that a common set of letter dissimilarities drive the part-sum model suggests that the part-sum model can be simplified. We therefore devised a reduced version of the part-sum model – called the Independent Spatial Interaction (ISI) model – in which the CL, CR, X and W terms are scaled versions of the single letter dissimilarities (Figure S8E). Specifically, the dissimilarity between bigrams AB & CD is:

$$d(AB, CD) = \alpha_{10}d_{AC} + \alpha_{20}d_{BD} + \alpha_{11}(d_{AD} + d_{BC}) + \beta_{11}(d_{AB} + d_{CD}) + c$$

where $d_{AC}$ is the observed dissimilarity between the left letters A & C from visual search and $\alpha_{10}$ is an unknown scaling term, $d_{BD}$ is the observed dissimilarity between the right letters B & D, and $\alpha_{20}$ is an unknown scaling term. Likewise, $\alpha_{11}$ is an unknown scaling term for the net dissimilarity ($d_{AD} + d_{BC}$) between letters across locations, $\beta_{11}$ is the unknown scaling term for the net dissimilarity ($d_{AB} + d_{CD}$) between letters within the two bigrams and $c$ is a constant. Thus, the ISI model has only 5 free parameters: $\alpha_{10}, \alpha_{20}, \alpha_{11}, \beta_{11}$ and $c$. These parameters can be estimated by solving the matrix equation **y** = **Xb** where **y** is a 1176x1 vector of observed bigram dissimilarities, **X** is a 1176 x 5 matrix containing the net single dissimilarity of each type (CL, CR, X & W) that contributes to the total dissimilarity, and **b** is a 5 x 1 vector of unknown weights corresponding to the contribution of each type of dissimilarity (plus a constant).

The performance of the ISI model is summarized in Figure S8F. It can be seen that, despite having only 5 free parameters compared to 85 parameters of the part-sum model, the ISI model yields comparable fits to the data (Figure S8F).

658



659

**Figure S8. Predicting bigram dissimilarity using part-sum model**

(A) Schematic of the part sum model. According to this model, the dissimilarity (1/RT) between bigrams 'AB' and 'CD' is written as a linear sum of dissimilarities of its corresponding part terms (AC and BD, shown in red), across part terms (AD and BC, shown in yellow), and within part terms (AB and CD, shown in blue).

(B) Correlation between the observed and predicted dissimilarities (1/seconds). Each point represents one search pair (n = $^{49}C_2$ = 1176). Word-word pairs are highlighted using red diamonds, and frequent bigram pairs are highlighted using blue circles. Dotted lines represent unity slope line.

(C) Correlation between the estimated weights at corresponding location left with estimated weights at 1) corresponding location right (red), 2) across location (yellow), and 3) within location (blue). Each point represents one letter pair (n = $^{7}C_2$ = 21). Dotted lines represent positive and negative unity slope line.

(D) Perceptual space of the single letter dissimilarities, that are the model coefficients of part terms at left corresponding location

(E) Schematic of the Independent Spatial Interaction model. In this model, we use the observed letter-pair dissimilarities and only estimate the weights of these letter-pair dissimilarities across different locations.

(F) Comparing part-sum and ISI model fits. Bar plots represents mean correlation coefficient between the observed and predicted dissimilarities. Error bars represent one standard deviation across 30 splits. Black horizontal line represents mean split-half correlation ($r_{sh}$) and the shaded error bar represents one standard deviation around the mean. (****, p < 0.00005, **, p < 0.005).

683
684

685 **ISI model performance across all experiments**
686       Next we asked whether the ISI model can be generalized to explain
687 dissimilarities between longer strings. Consider two n-letter strings $u_1 u_2 u_3 u_4 \ldots u_n$ and
688 $v_1 v_2 v_3 v_4 \ldots v_n$. The net dissimilarity between the two strings can be written as:
689

690
$$d(u_1 u_2 \ldots u_n, v_1 v_2 \ldots v_n) = \sum_{i=0}^{n} \sum_{k=0}^{n-i} \alpha_{ik}\big(d(u_i, v_{i+k}) + d(v_i, u_{i+k})\big) - \sum_{i=0}^{n} \sum_{k=1}^{n-i} \beta_{ik}\big(d(u_i, u_{i+k}) + d(v_i, v_{i+k})\big) + c$$

691
692 where $\alpha_{ik}$ are the unknown weights corresponding to pairs of letters across the two n-
693 grams separated by "k" positions starting from 0, and $\beta_{ik}$ are the unknown weights
694 corresponding to pairs of letters separated by "k" positions within the two n-grams.
695 Written in this manner, the total number of unknowns in the n-gram ISI model is $n^2+1$,
696 which can be estimated using standard linear regression as before. For instance, for
697 the 6-gram ISI model, there are $6^2+1 = 37$ free parameters.
698       In this manner, we fit the ISI model to all experiments. The resulting cross-
699 validated model fits are shown together with the letter model in Figure S9. It can be
700 seen that the ISI model performance is comparable to that of the letter model across
701 all experiments.
702



703
704 **Figure S9. ISI & letter model performance across all experiments**
705 For each experiment, we obtained a cross-validated measure of both neural and ISI
706 model performance as follows: each time we divided the subjects randomly into two
707 halves, and trained the letter model on one half of the subjects and tested it on the
708 other half. This was repeated for 30 random splits. The correlation between the model
709 predictions and the average dissimilarity from the held-out half of the data was taken
710 to be the model fit. The correlation between the observed dissimilarity between the
711 two random splits of subjects is then the upper bound on model performance (mean ±
712 std shown as *gray shaded bars*).
713

**Reducing the complexity of the ISI model**

According to the ISI model, the net dissimilarity between two n-grams can be written as a weighted sum of dissimilarities between letter pairs that are varying distances apart. We wondered if the ISI model can be simplified further if there is a systematic pattern whereby these weight corresponding to a given letter pair varies systematically with letter position and distance between the letters.

To assess this possibility, we plotted model coefficients of the ISI model estimated from Experiment S3 along two dimensions. First, we asked if the contribution of letter pairs at corresponding locations in the two n-grams varies with letter position. For varying string lengths (3-, 4-, 5- and 6-letter strings) we observed a characteristic U-shaped function whereby the edge letters contribute more to the net dissimilarity compared to the middle letters (Figure S10A). Second, we asked if model weights decrease systematically with inter-letter distance. This was indeed the case regardless of the starting letter in the pair (Figure S10B). Finally, we note that across and within part terms are roughly equal in magnitude but opposite in sign (Figure S8C).

The above pattern of weights in the ISI model suggest that we can make two simplifying assumptions. First, the weight of the starting letter is a U-shaped function when the inter-letter distance is zero ($\alpha_{i0}$). Second, weights decrease exponentially thereafter with increasing inter-letter distance. Specifically:

$$\alpha_{i0} = ai^2 + bi + c \ for \ i \ = \ 1,2,\ldots n$$
$$\alpha_{ik} = \alpha_{i0} e^{-k/\tau} \ for \ k \ \geq \ 1$$
$$\beta_{ik} = -\alpha_{ik} \ for \ k \ \geq \ 1$$

where $a, b, c \ and \ \tau$ are the free parameters in this model. This simplified model, which we call the Spatial Interaction Decay (SID) model has only 5 parameters and can be used to predict the dissimilarities between strings of arbitrary length. The model parameters are obtained using nonlinear gradient descent methods (*nlinfit* function, MATLAB).

To illustrate the performance of the SID model in comparison to the ISI model, we fit the model to 6-letter compound words (Experiment 4). To compare the two models, we plotted the ISI model terms directly estimated from the search data against the ISI model terms predicted from the SID model. This yielded a strong positive correlation (Figure S10C). The SID model also yielded excellent fits to the data (Figure S10D), and both models yielded comparable fits (Figure S10E).

To evaluate this pattern across all experiments, we fit both SID and ISI models to all experiments. Here too we obtained qualitatively similar fits for the two models (Figure S11). To confirm whether the SID model trained on one experiment can capture the variations in another, we trained the SID model on data from Experiment S5 and evaluated it on all other experiments. This too yielded largely similar but smaller predictions (Figure S11). This decrease in model fit suggests that model parameters are somewhat dependent on the search pairs chosen.

We conclude that dissimilarities between arbitrary letter strings can be predicted using highly simplified models that operate on single letter dissimilarities and simple compositional rules.
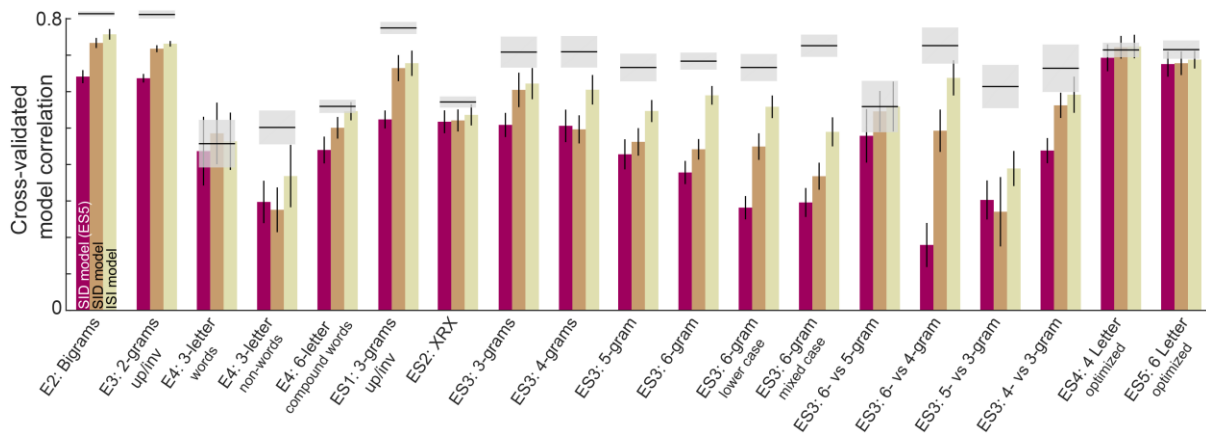
758
759 **Figure S10. Reducing the ISI model**
760     (A) ISI model coefficients $\alpha_{i0}$ as a function of starting letter position $i$, for Experiment
761         S3, for varying string lengths.
762     (B) ISI model coefficients $\alpha_{1k}$ as a function of inter-letter distance $k$ for Experiment
763         S3, for varying string lengths.
764     (C) ISI model coefficients (both $\alpha_{ik}$ and $\beta_{ik}$) plotted against the predicted ISI model
765         coefficients from the SID model. Both models are fitted to data from Experiment
766         4 (compound words).
767     (D) Observed dissimilarity in Experiment 4 plotted against predicted dissimilarity
768         from the SID model.
769     (E) Cross-validated model correlation for ISI & SID models.
770

771



772

773 **Figure S11. ISI and SID model fits across all experiments.** Cross-validated model
774 fits for the ISI and SID models across all experiments. In each case the SID and ISI
775 models were fit on a randomly chosen half of the subjects and tested on the other half.
776 The SID (ES5) bars refer to the SID model trained on Experiment S5 and tested on
777 data from a randomly chosen half of subjects in each experiment.
778

**Comparing upright and inverted bigrams using part-sum model**

The results in Section A2 were based on fitting the letter model to upright and inverted bigrams but assuming a fixed set of single letter responses derived from uppercase letters. The fact that the letter model yielded excellent fits to both upright and inverted bigrams validates this assumption. Nonetheless, we wondered whether differences between upright and inverted bigram searches can be explained solely by different letter representations or by differences in letter interactions.

To investigate this possibility, we fit the part-sum model to upright and inverted bigram searches (Figure S12A). The part-sum model also yielded equivalent fits to both upright and inverted searches (Figure S12B). If model predictions were similar, we reasoned that the difference between upright and inverted searches must be explained by differences in model parameters. To this end, we compared the estimated letter dissimilarities of each type (CL, CR, X and W) in the upright and inverted searches (Figure S12C). Model terms were comparable in magnitude for the CL terms, but were systematically weaker for both CR, X and W terms for inverted compared to upright searches (Figure S12C). However in all cases, the recovered letter dissimilarities were correlated between upright and inverted conditions (correlation between upright and inverted model terms: $r = 0.93, 0.91, 0.97$ & $0.87$ for CL, CR, X & W terms; all correlations $p < 0.00005$).



**Figure S12. Part-sum model fits for upright and inverted bigrams**

(A) Schematic of the part-sum model, in which the net dissimilarity between two bigrams is given as a linear sum of letter dissimilarities at corresponding locations (CL & CR), across-bigrams (X) and within-bigrams (W).

(B) Cross-validated model correlation of the part sum model for upright and inverted bigrams.

(C) Average model coefficients (mean ± sem) of each type for upright and inverted bigrams. Asterisks denote statistical significance (**** is $p < 0.00005$) obtained on a sign-rank test comparing 15 letter dissimilarities between upright and inverted conditions).
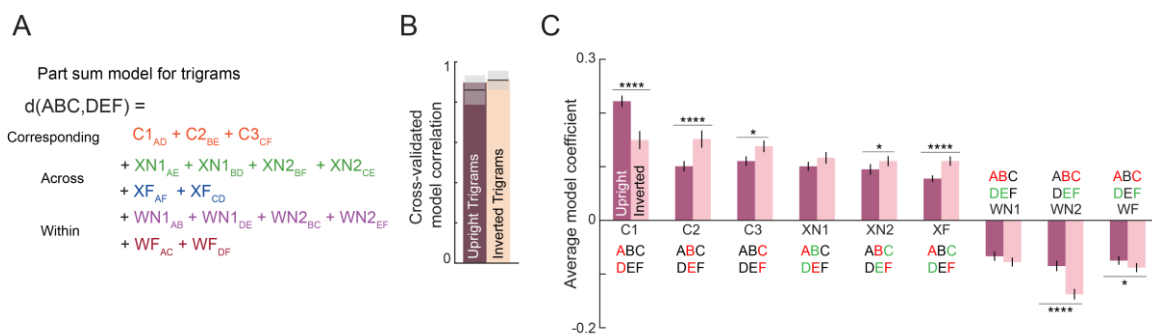
**Comparing upright and inverted trigrams using part-sum model**

813
814    The part sum model applied to trigrams is depicted in Figure S13A. In this
815  model, the net dissimilarity between two trigrams can be written as a sum of single
816  letter dissimilarities at every possible pair of locations. These locations are grouped as
817  corresponding letters at left (C1), middle (C2) and right (C3) locations, letters across
818  trigrams that are one letter apart starting from the left letter (XN1) or the middle letter
819  (XN2), letters across trigrams that are two letters apart (XF), letters within each trigram
820  that are one letter apart starting from the left letter (WN1) or middle letter (WN2), and
821  letters within each trigram that are two letters apart (WF). Thus the full part-sum model
822  has 9 groups of letter dissimilarities (C1, C2, C3, XN1, XN2, XF, WN1, WN2, WF) each
823  having $^6C_2 = 15$ unknown single letter dissimilarities. Together with a constant term,
824  this part-sum model has 9 x 15 + 1 = 136 free parameters. Since we have 500
825  searches each for upright and inverted trigrams, the part-sum model can be fit to this
826  data to estimate these free parameters using standard linear regression.
827    Cross-validated model fits for the part-sum model are shown in Figure S13B. It
828  can be seen that the part-sum model explains nearly all the explainable variance in
829  the data for both upright and inverted trigrams (Figure S13B). This in turn means that
830  differences between upright and inverted trigrams can be explained using differences
831  in model parameters. This was indeed the case: on plotting the strength of model terms
832  of each type it was clear that 7 of the 9 types of model terms (C1, C2, C3, XN2, XF,
833  WN2, WF) were systematically larger for upright trigrams compared to inverted
834  trigrams (Figure S13C). Finally we confirmed that model terms for upright and inverted
835  trigrams were highly correlated (correlation between upright and inverted model terms,
836  averaged across 9 model term types: r = 0.65 ± 0.1, p < 0.05 in all cases).
837    We conclude that upright and inverted trigram searches can be explained using
838  the part-sum model driven by a common single letter representation.
839



840
841  **Figure S13. Part-sum model fits for upright and inverted trigrams**
842      (A) Schematic of part-sum model for trigrams.
843      (B) Cross-validated model correlation of part-sum model for upright and inverted
844          trigrams.
845      (C) Average model coefficient (averaged across $^6C_2 = 15$ terms) of each type for
846          upright and inverted trigrams. Asterisks indicate statistical significance (* is p <
847          0.05, ** is p < 0.005, etc) calculated using a sign-rank test comparing the upright
848          and inverted model terms.
849

850 **SECTION A6. JUMBLED WORD READING (EXPT S6)**

851

852       Here, in Experiment S6, we tested subjects on a jumbled word reading task,
853 where they had to view a jumbled word and recognize the original word.

854

855 **METHODS**
856 *Procedure.* A total of 16 subjects (9 male, aged 24.8 ± 2.1 years) participated in the
857 task. Other details were similar to Experiment 5.
858 *Stimuli.* We chose 300 words such that no two words were anagrams of each other.
859 These comprised 75 four-letter words, 150 five-letter words and 75 six-letter words.
860 Jumbled words were created by shuffling 2, 3, or 4 letters of each word. There were
861 an equal proportion of 2, 3, and 4 letter transpositions. All stimuli were presented in
862 uppercase against a black background.
863 *Task.* Each trial began with a fixation cross shown for 0.5 s followed by a jumbled word
864 that appeared for 5 seconds (for the first 6 subjects) and 7 seconds (for the rest), or
865 until the subject made a response by pressing the space bar on the keyboard. Subjects
866 were asked to press a key as soon as they could recognize the unjumbled word. To
867 ensure that subjects correctly recognized the unjumbled word, they were asked to type
868 the unjumbled word within 10 seconds of pressing the space bar. The response time
869 was taken as the time at which the subject pressed the space bar. To avoid any
870 memory effects, the same set of jumbled words were shown to all subjects exactly
871 once. We analysed response times only on trials in which the subject subsequently
872 entered the correct word.
873 *Data Analysis.* Subjects were reasonably accurate on this task (average accuracy:
874 59.5 ± 8% across 300 words). Response times for wrongly typed words were
875 discarded. Words correctly solved by more than 6 subjects (n = 238) were included for
876 further analysis. Since trials were self-paced, we did not remove any outliers in the
877 reaction times. Lexical properties were obtained from the English Lexicon Project
878 (Balota et al., 2007).

879

880 **RESULTS**
881       Of a total of 300 jumbled words tested, we selected for further analysis 238
882 words that were correctly unjumbled by more than two-thirds of the subjects. Subjects
883 responded quickly and accurately to these words (mean ± std of accuracy: 71 ± 9%;
884 response time: 2.13 ± 0.33 s across 238 words). Subjects took longer to respond to
885 some jumbled words (e.g. REHID) compared to others (e.g. DBTOU), as seen in the
886 sorted response times (Figure S14A). These patterns were consistent across subjects,
887 as evidenced by a significant split-half correlation (r = 0.55, p < 0.00005 between odd-
888 and even-numbered subjects).
889       Can these patterns in unscrambling time be explained using the letter model?
890 To do so, we reasoned that jumbled words with large dissimilarity to the original word
891 will take longer to elicit a response (Figure S14B). Accordingly, we took the average
892 response times to each jumbled word and asked whether it can be predicted using the
893 single letter model described previously. For each word length, we optimized the

894  weights of the single letter model to find the best fit to this data, and then combined
895  the predictions across all word lengths to obtain a composite measure of performance.
896  The single letter model yielded excellent fits to the data (r = 0.76, p < 0.00005; Figure
897  S14C). This model fit was comparable to the data consistency ($r_{data}$ = 0.70). An
898  alternate distance model - Orthographic Levenshtein (OL) distance (Levenshtein,
899  1966) – calculates the number of edits required to transform one string to other. This
900  model neither accounts for letter similarity nor the position of edit. Hence, it fails to
901  account for all the variance in the data (r = 0.44, p < 0.00005; Figure S14D).

902          The above finding shows that human performance on unscrambling words is
903  driven primarily by the visual dissimilarity between the jumbled and original word.
904  However, it does not rule out the presence of lexical factors. To assess this possibility
905  we formulated a model to predict the unscrambling time as a linear sum of many lexical
906  factors. We used five lexical properties: log word frequency, log mean letter frequency,
907  log mean bigram frequency of the jumbled word, log mean bigram frequency of the
908  unjumbled i.e. original word, and the number of orthographic neighbours (see
909  Methods). To avoid overfitting by either model, we trained both models on one-half of
910  the subjects and tested it on the other half. This lexical model yielded relatively poor
911  fits (r = 0.30, p < 0.00005, Figure S14E) compared to visual dissimilarity from both
912  single letter model and OL distance model. The difference in model fits was statistically
913  significant (p < 0.05, Fisher's z-test).  Among the lexical factors, word frequency and
914  letter frequency contributed the most compared to the others (partial correlation of
915  each lexical factor after accounting for all others: r = -0.23, p < 0.0005 for log word
916  frequency, r = 0.18, p < 0.05 for log mean letter frequency; r = .05, p = 0.49 for log
917  mean bigram frequency of jumbled word; r = -0.02, p = 0.77 for log mean bigram
918  frequency in original word; r = 0.04, p = 0.58 for number of orthographic neighbours).

919          To assess the extent of shared variance in the two models, we calculated the
920  partial correlation between the observed data and the lexical model predictions after
921  factoring out the contribution from visual dissimilarity. This revealed a small partial
922  correlation (r = 0.31, p < 0.00005). Conversely, the partial correlation for the single
923  letter model after factoring out the lexical model was much higher (r = 0.75, p <
924  0.00005). Thus, visual dissimilarity from the single letter model dominates jumbled
925  word reading.

926          Finally we asked whether both visual dissimilarity and lexical factors contribute
927  to the jumbled word task. We created a combined model in which the jumbled word
928  response times were a linear combination of the predictions of both models. This
929  combined model yielded better predictions than either model by itself (r = 0.78, p <
930  0.00005, Figure S14E). To assess the statistical significance of these results, we
931  performed a bootstrap analysis. On each trial, we trained three models on the
932  dissimilarity obtained from considering only one randomly chosen half of subjects: the
933  visual dissimilarity model, the lexical model and the combined model. We calculated
934  the correlation between all three model predictions on the other half of the data, and
935  repeated this procedure 1000 times. The OL distance model does not have any free
936  parameters, hence the distances were directly correlated with the other half of the
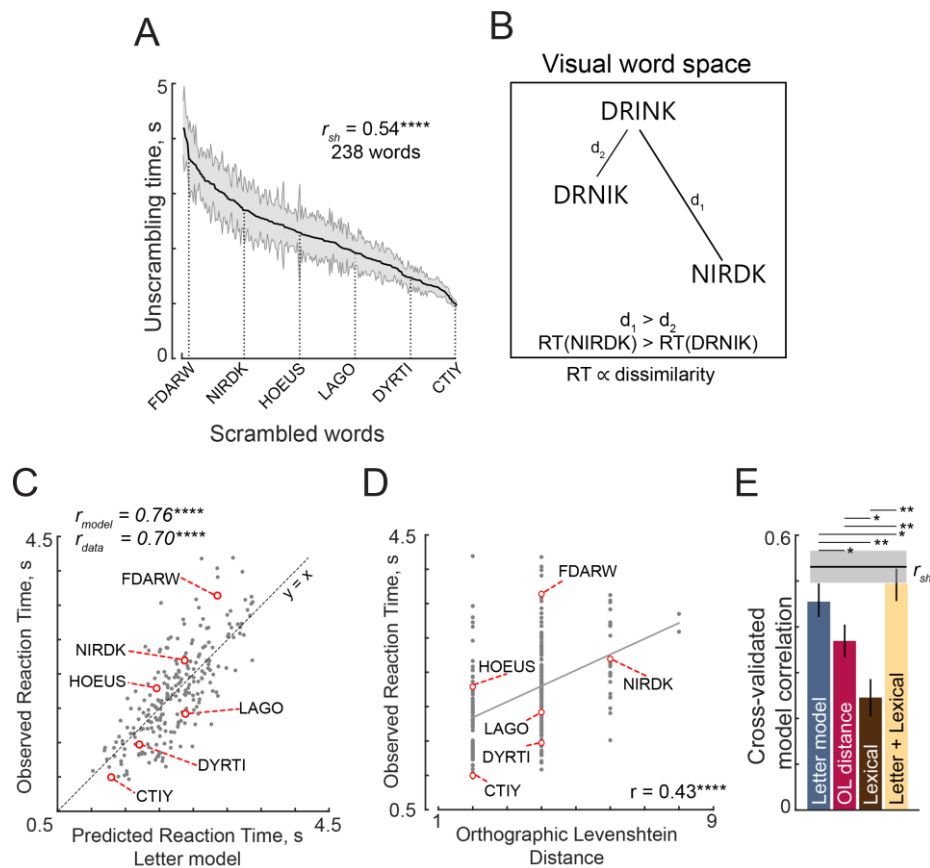937  data. Across these samples, the lexical model fits never exceeded the visual

938    dissimilarity model, suggesting that the visual dissimilarity model was significantly
939    better ($p < 0.05$). Likewise, the combined model was only marginally better than the
940    visual letter model (fraction of combined < visual: $p = 0.07$) but was significantly better
941    than the lexical model (fraction of combined < lexical: $p = 0$).

942          We conclude that performance on the jumbled word task relies primarily on
943    visual dissimilarity. We propose that this initial visual representation of a word allows
944    the subject to make a quick guess at the correct word without explicit symbolic
945    manipulation.

946

947

## Scrambled word task



**Figure S14. Jumbled word task (Experiment S7).**

(A) Response times in the jumbled word task sorted in descending order. Shaded error bars represent s.e.m. Some example words are indicated using dotted lines. The split-half correlation between subjects ($r_{sh}$) is indicated on the top left.

(B) Schematic of visual word space, with one stored word (DRINK) and two jumbled versions (DRNIK & NIRDK). We predicted that the time taken by subjects to unscramble a jumbled word would be proportional to its dissimilarity to the stored word. Thus, subjects would take longer to unscramble NIRDK compared to DRNIK.

(C) Observed response times in the jumbled word task plotted against predictions from the letter model based on single letters with spatial summation. Each point represents one word. Asterisks indicate statistical significance (**** is p < 0.00005).

(D) Observed response times in the jumbled word task plotted against Orthographic Levenshtein (OL) distance. Each point represents one word. Asterisks indicate statistical significance (**** is p < 0.00005).

(E) Cross-validated model correlations for the letter model, OLD model, lexical model and the neural+lexical model. Model correlations were obtained by training each model on one half of subjects, and evaluating the correlation on the other half (error bars represent standard deviation across 1000 random splits). The upper bound on model fits is the split-half correlation ($r_{sh}$), shown in black with shaded error bars representing standard deviation across the same random splits. All correlations were individually statistically significant (p < 0.00005). Horizontal lines above shaded error bar depicts significant difference across different models i.e. the fraction of splits in which the observed difference was violated. All significant comparisons are indicated.

974 **SECTION A7. ADDITIONAL ANALYSES FOR EXPERIMENTS 6 & 7**

975 **Stimulus set**
976 32 words were chosen of varying frequency of occurrence and the nonwords were
977 created by either transposition or substitution of middle or edge letters. 10 single
978 letters: E, S, A, R, O, L, I, T, N, and D were used to form words. The full set of strings
979 used experiments 6 and 7 is shown below.
980

| Middle Letter Transposition | | Edge Letter Transposition | | Middle Letter Substitution | | Edge Letter Substitution | |
|---|---|---|---|---|---|---|---|
| Words | Nonwords | Words | Nonwords | Words | Nonwords | Words | Nonwords |
| *AORTA* | AROTA | *STOLE* | TSOLE | *NOISE* | NANSE | *ONION* | ESION |
| *DRAIN* | DARIN | *OASIS* | AOSIS | *ERROR* | EDLOR | *RADIO* | EEDIO |
| *TREND* | TERND | *SOLID* | OSLID | *DRILL* | DTELL | *ASSET* | EESET |
| *ATLAS* | ALTAS | *TRAIN* | RTAIN | *ARISE* | AOESE | *TEASE* | RDASE |
| *DRONE* | DRNOE | *ORDER* | ORDRE | *LITRE* | LINOE | *ENTER* | ENTRO |
| *LEARN* | LERAN | *INDIA* | INDAI | *SLIDE* | SLONE | *IDEAL* | IDEDI |
| *SANTA* | SATNA | *RINSE* | RINES | *NASAL* | NATDL | *ADORE* | ADODI |
| *INSET* | INEST | *SNAIL* | SNALI | *ALIEN* | ALOTN | *LASER* | LASRO |

981 **Table S3: List of 32 words and 32 nonwords used in Experiment 6 & 7.** All words
982 and nonwords were created from 10 single letters whose activations were also
983 measured in the experiment.
984
985 **ROI definitions**
986

| ROI | Definition | #voxels (mean ± sd) | ROI peak location |
|---|---|---|---|
| V1-V3 | Voxels activated for scrambled > fixation overlaid with anatomical mask of V1-V3 | 398 ± 131 | X: 8 ± 17 Y: -96 ± 5 Z: 6 ± 9 |
| V4 | Voxels activated for scrambled > fixation overlaid with anatomical mask of V4 | 185 ± 63 | X: 5 ± 26 Y: -88 ± 3 Z: 27 ± 11 |
| LO | Voxels activated for object > scrambled and not in other ROIs | 371 ± 115 | X: -17 ± 43 Y: -66 ± 15 Z: -19 ± 5 |
| VWFA | Voxels with known words > scrambled word in a contiguous region in fusiform gyrus | 52 ± 15 | X: -44 ± 4 Y: -50 ± 5 Z: -17 ± 5 |
| TG | Voxels with native words > scrambled word in a contiguous region in temporal gyrus | 289 ± 182 | X: -44 ± 39 Y: -43 ± 18 Z: 3 ± 9 |

987 **Table S4. Variability in ROI definitions across subjects.** For each ROI we report
988 the mean and standard deviation across subjects of the number of voxels, and the
989 XYZ location of the voxel with peak T-value in the normalized brain.
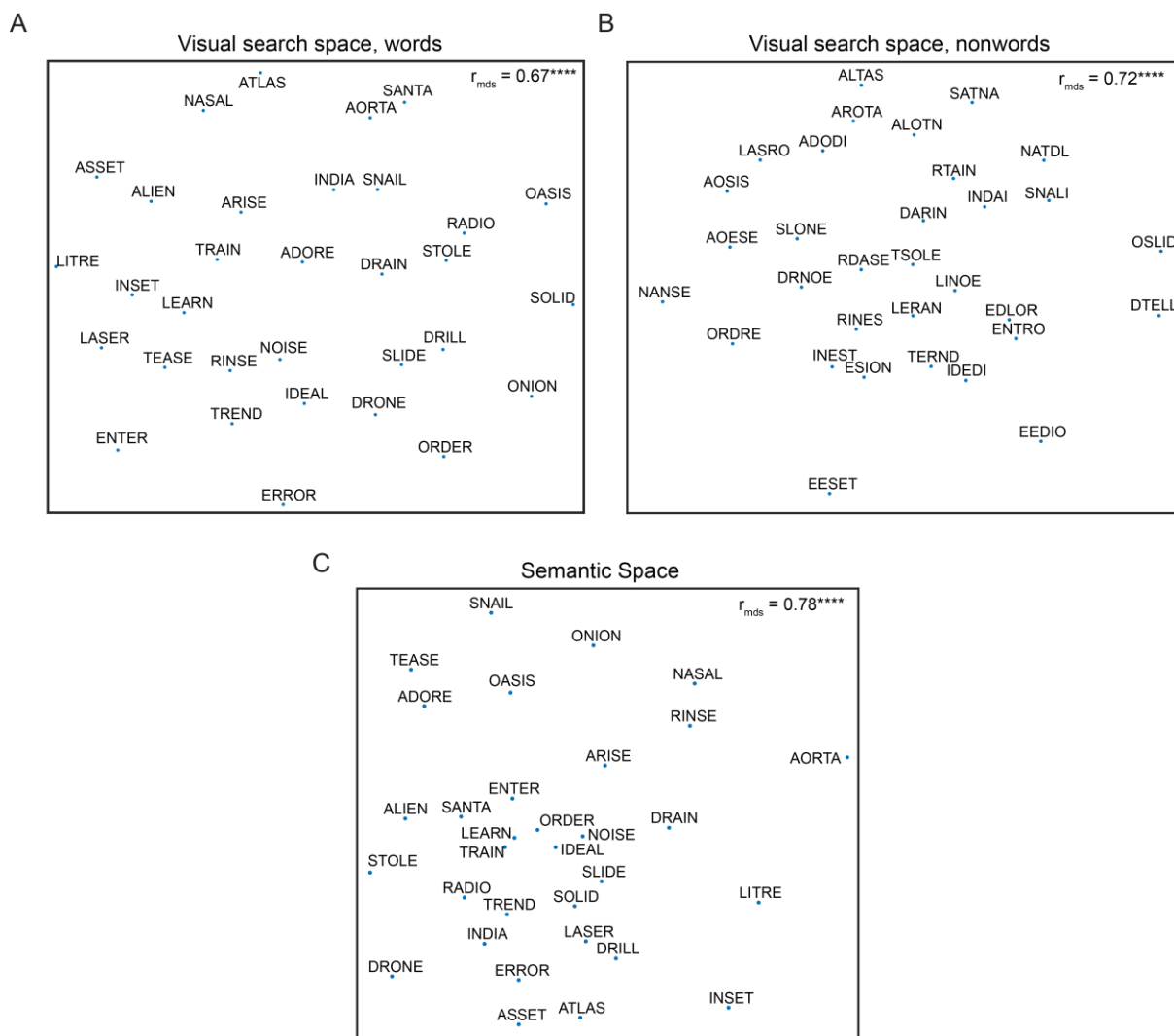990
991
992
993

**Visualization of perceptual and semantic space**

To visualize words and nonwords in perceptual space, we performed a multidimensional scaling (MDS) analysis of the visual search data (Experiment 7). Briefly, MDS finds the best-fitting 2D coordinates that best match with the observed distances. In the resulting plot, nearby stimuli correspond to hard searches. The perceptual space for words and nonwords is shown in Figure S15 A-B. It can be seen that stimuli with common first letters are grouped together. MDS coordinates for nonwords was rotated without altering their overall configuration so as to best match the MDS coordinates for words.

The semantic dissimilarities were estimated using the GloVe features (Pennington et al., 2014), and visualized using MDS analysis (Figure S15C). In the resulting plot, semantically related words/ frequently cooccurring words are closer to each other.



**Figure S15: Multi-dimensional representation of words and nonwords.**

A. Perceptual space for words. we used multidimensional scaling to find the 2D coordinates of all words that best match the observed distances. In the resulting plot, nearby words indicate hard searches. The correlation coefficient between dissimilarities in 2D plane and the observed data is shown. Asterisks indicate significant correlation (**** is $p < 0.00005$).

B. Same as (A) but for nonwords.

C. Same as (A) but for semantic space of words.

**Neural activity corresponding to words, nonwords, and letters**

For each category of stimuli i.e. words, nonwords, and letters, we averaged the activity values across voxels and subjects within each ROI. The mean activity values are shown in Figure S16A-E.

**Word vs nonword classification**

For each ROI and subject, we built linear classifier to discriminate between words and nonwords using built-in MATLAB routine "fitcdiscr". We built separate classifiers to distinguish the activity pattern of transposed and substituted nonwords from their corresponding word activity patterns. The resulting decoding accuracy is shown in Figure S16F.

**Can string responses be predicted from single letters?**

We modelled the response of each voxel across the 64 strings (32 words, 32 nonwords) as a linear combination of the single letter activations (Figure S16G). We evaluated model fits by comparing model correlations separately for words and nonwords. If string responses were driven by specialized detectors for letter combinations (such as those present in words), then we reasoned that model correlations would be worse for words compared to nonwords. By contrast, if there are no specialized detectors of this kind, model fits would be equivalent for words and nonwords.

We calculated cross-validated model fits by training the model on half the trials and testing it on the other half of the trials. Since voxels could vary widely in their reliability of responses to the stimuli, we normalized the model fit of each voxel by its split-half reliability. The average noise-corrected model fit (averaged across voxels and subjects) is shown in Figure S16H. This revealed no systematic difference in model performance for words and nonwords in any of the ROIs (Figure S16H). We obtained qualitatively similar results using a searchlight, where there were no clear regions in which model fits differed for words and nonwords (Figure S17D).

To further validate the letter model, we compared the single letter tuning along each MDS dimension with the observed single letter tuning in each ROI (Figure S18A). For each ROI, we grouped voxels with similar response profile and matched it to the MDS dimension (Figure S18A). We obtained similar single letter tuning and weight profiles for voxels across different ROIs. However this analysis is inconclusive because there is no systematic way to compare a small set of neurons inferred from behaviour with the much larger, possibly overcomplete set of voxel activations observed in brain imaging. Likewise, we grouped voxels with similar summation weights to compare the weight profiles in behaviour and brain imaging. However this analysis is also inconclusive because different MDS-derived neurons might contribute differently towards behaviour, so the summation weights cannot be directly averaged to make overall comparisons between ROI activations and behaviour. Despite these caveats, there is a general match between tuning profiles and summation weights observed in behaviour with those observed in different brain regions.
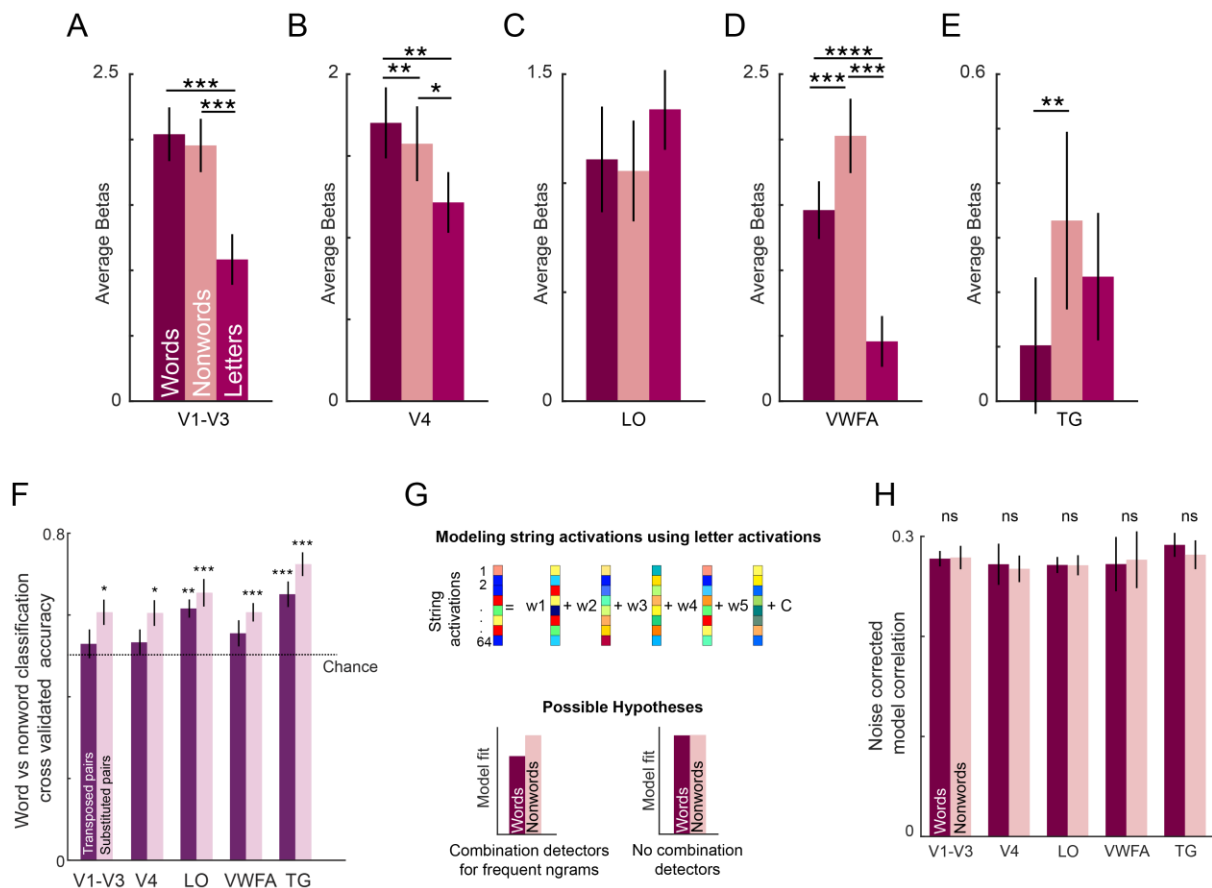
**Figure S16: Neural activity**

(A) Average activation levels for words, nonwords, and letters. Error bar indicate ±1 s.e.m. across subjects. Asterisks indicate statistical significance (* is $p < 0.05$, ** is $p < 0.005$, etc. in a sign-rank test comparing subject-wise average activations).

(B)-(E). Same as in A but for V4, Lateral Occipital areas, Visual Word Form Area, and Temporal Gyri respectively.

(F) Cross-validated classification accuracy for transposed word-nonword pairs (*dark*) and substituted word-nonword pairs (*light*). Error bars indicate s.e.m. across subjects. Asterisks indicate statistical significance (* is $p < 0.05$, ** is $p < 0.005$, etc. in a sign-rank test comparing subject-wise accuracy w.r.t. chance level).

(G) Schematic of the voxel model. The response of each voxel across strings is modelled as a linear combination of the constituent letter responses. Bottom: Hypothetical model fits based on the presence (right) or absence (left) of local combination detectors. Predicted responses for words will deviate from the observed responses under the influence of LCD.

(H) Average model correlation (normalized using split-half correlation) for each ROI for words (dark) and nonwords (light). Error bar indicate s.e.m. across subject.

**Searchlight analyses**

1088
1089       To identify other brain regions that might show the effects observed in the
1090 individual ROIs, we performed a whole-brain searchlight analysis. Specifically, for
1091 each voxel in a given subjects' brain, we considered a local neighbourhood of 27
1092 voxels (3x3x3 voxels) and performed the following analyses of interest. We obtained
1093 similar results for larger searchlight volumes. The resulting maps were smoothed using
1094 a Gaussian filter with FWHM of 3 mm
1095
1096 *Searchlight for regions that match lexical decision time*
1097       For each voxel, its activity across strings is correlated with mean lexical
1098 decision time. The resulting whole brain correlation map is averaged across subjects.
1099 Overall, activity in VWFA, Superior Parietal Lobe (SPL), Pre-Frontal and motor cortex
1100 is correlated with lexical decision time. This correlation map was visualized on the
1101 brain surface (Figure S17A).
1102
1103 *Searchlight for regions that match perceptual space*
1104       For the neighbourhood of each voxel, we calculated the pairwise neural
1105 dissimilarity for all word-word, nonword-nonword, and word-nonword pairs for a given
1106 subject, and averaged this across subjects. We then calculated the correlation
1107 between this local neural dissimilarity and the corresponding string dissimilarities
1108 estimated using experiment 7. This correlation map was visualized on the brain
1109 surface (Figure S17B).
1110
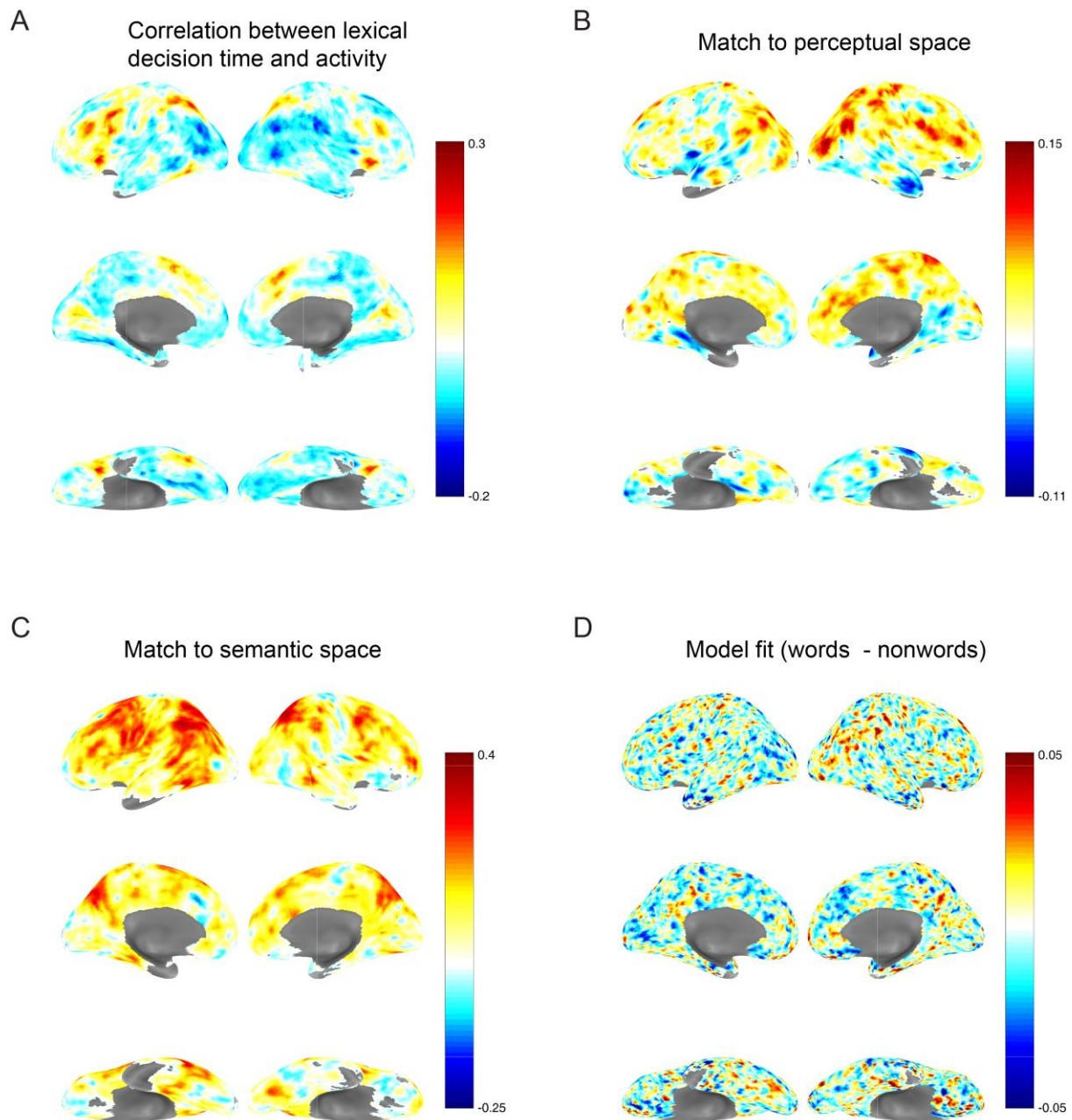1111 *Searchlight for regions that match semantic space*
1112       For the neighbourhood of each voxel, we calculated the pairwise neural
1113 dissimilarity for all word-word pairs for a given subject and averaged this across
1114 subjects. We then calculated the correlation between this local neural dissimilarity and
1115 the corresponding semantic dissimilarities. This correlation map was visualized on the
1116 brain surface (Figure S17C).
1117
1118 *Searchlight for comparing linear model fits between words and nonwords*
1119       For each subject and voxel, we modelled the response to strings as a linear
1120 combination of its single letter responses. The model fits (correlation between
1121 observed and predicted string responses) was evaluated separately for words and
1122 nonwords. The difference in the mean model fits between words and nonword is
1123 visualized on the brain surface (Figure S17D).
1124

A

Correlation between lexical decision time and activity

B

Match to perceptual space

C

Match to semantic space

D

Model fit (words - nonwords)

**Figure S17: Searchlight analysis**

A. Searchlight map of correlation between neural activity and lexical decision time for each voxel.

B. Searchlight map of correlation between neural dissimilarity and search dissimilarities in behaviour.

C. Searchlight map of correlation between neural dissimilarity and semantic dissimilarities.

D. Searchlight map depicting the difference in model fit for words versus nonwords for each voxel, averaged across subjects.
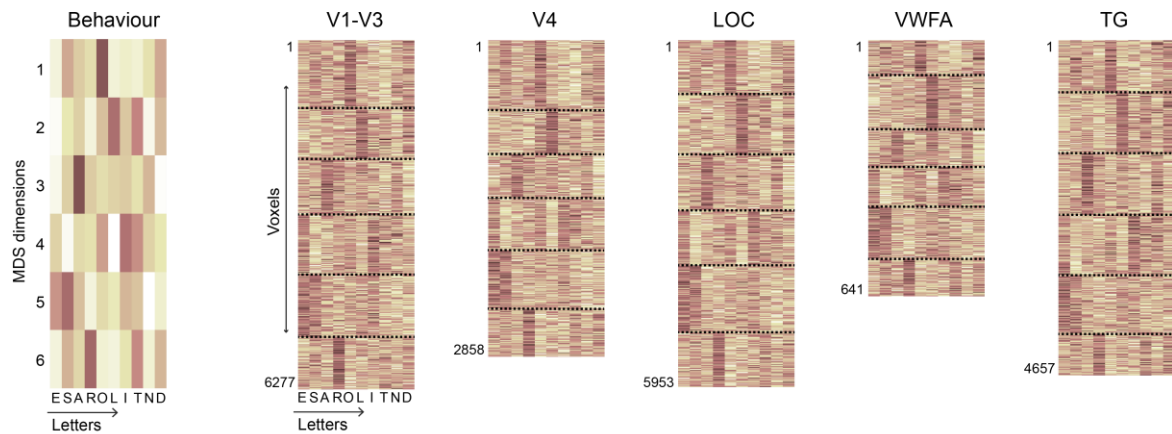
**Match between letter model and fMRI data**

The letter model described throughout the study is derived from dissimilarities measured in behaviour in two steps. First, the dissimilarities between single letters were used to construct single neurons tuned to letter shape, whose activity predicts these dissimilarities. Second, the summation weights of each neuron were adjusted so that they match the dissimilarities between longer strings.

Given that we recorded responses to single letters as well as strings in fMRI, we wondered whether these can be matched in some manner to the letter tuning and summation weights derived from behaviour in the letter model. Any direct comparison is fraught with the difficulty that many single letter tuning functions could produce the same behaviour. For instance, simply rotating the MDS-derived tuning functions could yield another set of neurons that match the observed letter dissimilarities. This is further compounded by the fact that the MDS-derived neurons contribute unequally to behaviour, and by the fact that this mapping could change completely with increasing numbers of neurons. Thus it is unreasonable to expect voxel tuning for single letters or the summation weights to match exactly with the behaviourally derived tuning.
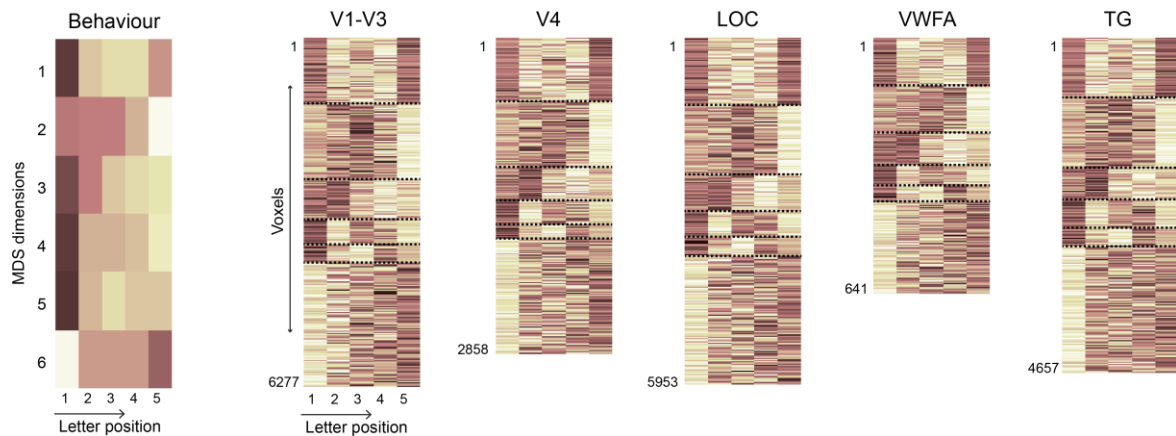
Nonetheless, we attempted to find a broad link between the single letter tuning and summation weights observed in behaviour with those observed in each ROI. The results are summarized in Figure S18. Since there are only 10 single letters, 6 MDS neurons were sufficient to explain > 95% of the variance of the pair-wise single letter dissimilarities observed in Experiment 1. For each MDS neuron, we identified the voxels whose activity for single letters had the least residual error compared to other MDS neurons. In this manner, we sorted the voxels into six groups corresponding to each MDS neuron. The resulting plots are shown in Figure S18A. It can be see that all ROIs show single letter tuning profiles similar to the behaviourally derived single letter tuning profiles. The corresponding summation weights for these voxels are shown in Figure S18B. Once again, it can be seen that many ROIs show similar summation weights as those observed in behaviour.

1165



1166
**Figure S18: Comparison of letter tuning and summation weights**
A. (*Left*) Response of 6 MDS neurons for all the 10 letters. (*Right*) Single letters response across all the voxels (concatenated across subjects) within a given ROI. Each voxels is sorted into one of 6 groups depending on which MDS neuron it matches best. The height of each ROI plot is logarithmically scaled to match the number of voxels across all subjects. Black dashed lines are used to separate the clusters corresponding to each MDS neuron.
B. Same as (A) but showing the summation weights corresponding to each MDS neuron or ROI voxel.

1176

## SUPPLEMENTARY REFERENCES

Balota DA et al. (2007) The English Lexicon Project. Behav Res Methods 39:445–459.

Dehaene S, Cohen L, Sigman M, Vinckier F (2005) The neural code for written words: a proposal. Trends Cogn Sci 9:335–341.

Duncan J, Humphreys GW (1989) Visual search and stimulus similarity. Psychol Rev 96:433–458.

Levenshtein V (1966) Binary codes capable of correcting deletions, insertions, and reversals. Sov Phys Dokl 10:707–710.

Mueller ST, Weidemann CT (2012) Alphabetic letter identification: Effects of perceivability, similarity, and bias. Acta Psychol (Amst) 139:19–37.

Pennington J, Socher R, Manning CD (2014) GloVe: Global Vectors for Word Representation. In: Empirical Methods in Natural Language Processing (EMNLP), pp 1532–1543.

Pramod RT, Arun SP (2016) Object attributes combine additively in visual search. J Vis 16:8.

Ratan Murty NA, Arun SP (2015) Dynamics of 3D view invariance in monkey inferotemporal cortex. J Neurophysiol 113:2180–2194.

Simpson IC, Mousikou P, Montoya JM, Defior S (2013) A letter visual-similarity matrix for Latin-based alphabets. Behav Res Methods 45:431–439.

Vighneshvel T, Arun SP (2013) Does linear separability really matter? Complex visual search is explained by simple search. J Vis 13:1–24.