1 **Neural dynamics at successive stages of the ventral visual stream are**
2 **consistent with hierarchical error signals**
3 Running title: Error signaling in the ventral stream

4      Elias B. Issa[1]\*, Charles F. Cadieu[2], and James J. DiCarlo
5
6      McGovern Institute for Brain Research and Department of Brain and Cognitive Sciences
7      Massachusetts Institute of Technology, Cambridge, Massachusetts 02139
8
9      *Correspondence: Elias Issa, elias.issa@columbia.edu
10
11     [1]Current address: Department of Neuroscience, Zuckerman Mind Brain Behavior Institute, Columbia
12     University, New York, NY 10027
13     [2] Current address: Bay Labs, Inc. San Francisco, CA 94102
14

15 **AUTHOR CONTRIBUTIONS**
16 Conceptualization, E.B.I and J.J.D.; Methodology, E.B.I., C.F.C, and J.J.D.; Investigation, E.B.I.; Formal
17 Analysis: E.B.I. and C.F.C; Writing – Original Draft, E.B.I., C.F.C., and J.J.D.; Writing – Review &
18 Editing, E.B.I. and J.J.D.; Funding Acquisition, E.B.I., C.F.C., and J.J.D.; Resources, J.J.D.;
19 Supervision, J.J.D.
20
21
22 **ABSTRACT**

23 Ventral visual stream neural responses are dynamic, even for static image presentations. However,

24 dynamical neural models of visual cortex are lacking as most progress has been made modeling static,

25 time-averaged responses. Here, we studied population neural dynamics during face detection across

26 three cortical processing stages. Remarkably, ~30 milliseconds after the initially evoked response, we

27 found that neurons in intermediate level areas decreased their preference for faces, becoming anti-face

28 preferring on average even while neurons in higher level areas achieved and maintained a face

29 preference. This pattern of hierarchical neural dynamics was inconsistent with extensions of standard

30 feedforward circuits that implemented recurrence within a cortical stage. Rather, recurrent models

31 computing errors between stages captured the observed temporal signatures. Without additional

32 parameter fitting, this model of neural dynamics, which simply augments the standard feedforward

33 model of online vision to encode errors, also explained seemingly disparate dynamical phenomena in

34 the ventral stream.

35

36    **INTRODUCTION**

37    The primate ventral visual stream is a hierarchically organized set of cortical areas beginning with the

38    primary visual cortex (V1) and culminating with distributed patterns of neural firing across the inferior

39    temporal cortex (IT) that explicitly encode objects (i.e. linearly decodable object identity) (Hung,

40    Kreiman, Poggio, & DiCarlo, 2005) and quantitatively account for core invariant object discrimination

41    behavior in primates (Majaj, Hong, Solomon, & DiCarlo, 2015). Formalizing object recognition as the

42    result of a series of feedforward computations yields models that achieve impressive performance in

43    object categorization (Krizhevsky, Sutskever, & Hinton, 2012)(Zeiler & Fergus, 2013) similar to the

44    absolute level of performance achieved by IT neural populations, and these models are the current best

45    predictors of neural responses in IT cortex and its primary input layer, V4 (Cadieu et al., 2014)(Yamins

46    et al., 2014). Thus, the feedforward inference perspective provides a simple but powerful, first-order

47    framework for the ventral stream and core invariant object recognition.

48

49         However, visual object recognition behavior may not be executed via a single feedforward neural

50    processing pass (a.k.a. feedforward inference) because IT neural responses are well-known to be

51    dynamic even in response to images without dynamic content (Brincat & Connor, 2006)(Sugase,

52    Yamane, Ueno, & Kawano, 1999)(Chen et al., 2014)(Meyer, Walker, Cho, & Olson, 2014), raising the

53    question of what computations those neural activity dynamics might reflect. Prior work has proposed

54    that such neuronal response dynamics could be the result of different types of circuits executing

55    different types of computation such as: 1) recurrent circuits within each ventral stream processing stage

56    implementing local normalization of the feedforward information as it passes through the stage

57    (Carandini, Heeger, & Movshon, 1997)(Schwartz & Simoncelli, 2001)(Carandini & Heeger, 2012), 2)

58    feedback circuits between each pair of ventral stream stages implementing the integration of top-down

59    with bottom-up information to improve the current (online) inference (Seung, 1997)(Lee, Yang, Romero,

60    & Mumford, 2002)(Zhang & Heydt, 2010)(Epshtein, Lifshitz, & Ullman, 2008), or 3) feedback circuits

61    between each pair of stages comparing top-down and bottom-up information to compute prediction

62    errors that guide changes in synaptic weights so that neurons are better tuned to features useful for

63    future feedforward behavior (learning) (Rao & Ballard, 1999). Thus, neural dynamics may reflect the

64    various adaptive computations (within-stage normalization, top-down Bayesian inference) or reflect the

65    underlying error intermediates that could be generated during those processes (e.g. predictive coding).

66

67    These computationally motivated ideas can each be implemented as neural circuits to ask which

68    idea best predicts response dynamics across the visual hierarchy. Here, our main goal was to look

69    beyond the initial, feedforward response edge to see if we could disambiguate among dynamics that

70    might result from stacked feedforward, lateral, and feedback operations. Rather than record from a

71    single processing level, we measured the dynamics of neural signals across three hierarchical levels

72    (pIT, cIT, aIT) within macaque IT. We focused on face processing subregions within each of these

73    levels for three reasons. First, prior evidence argues that these three subregions are tightly

74    anatomically and functionally connected and that the subregion in pIT is the dominant input to the

75    higher subregions (Grimaldi, Saleem, & Tsao, 2016)(Moeller, Freiwald, & Tsao, 2008). Second,

76    because prior work argues that a key behavioral function of these three subregions is to distinguish

77    faces from non-faces, this allowed us to focus our testing on a relatively small number of images

78    targeted to engage that processing function. Third, prior knowledge of pIT neural tuning properties (Issa

79    & DiCarlo, 2012) allowed us to design images that were quantitatively matched in their ability to drive

80    neurons in the pIT input subregion but that should ultimately be processed into two separate groups

81    (face vs. non-face). We reasoned that these images would force important computations for

82    disambiguation to occur somewhere between the pIT subregion and the higher level (cIT, aIT)

83    subregions. With this setup, our aim was to observe the dynamics at all three levels of the hierarchy in

84    response to that image processing challenge so that we might discover – or at least constrain -- which

85    type of computation is at work.

86

87      Consistent with the idea that the overall system performs face vs. non-face discrimination (i.e.

88      face detection), we found that in the highest face processing stage (aIT), neurons rapidly developed

89      and maintained a response preference for faces over non-faces even though our images were

90      designed to be challenging for frontal face detection. However, we found that many neurons in the early

91      (pIT) and intermediate (cIT) processing levels of IT had face selectivity that rapidly but paradoxically

92      *decreased*. That is, the responses of these neurons evolved to prefer images of non-faces over images

93      of faces within 30 milliseconds of their feedforward response. We found that standard feedforward

94      models that employ local recurrences such as adaptation, lateral inhibition, and normalization could not

95      capture this stage-wise pattern of image selectivity despite our best attempts. However, we found that

96      *decreasing* -- rather than increasing -- face preference in early and intermediate processing stages is a

97      natural dynamical signature of previously suggested "error coding" models (Rao & Ballard, 1999) in

98      which the neural spiking activity at each processing stage carries both an explicit representation of the

99      variables of interest (e.g. is a face present?) and an explicit encoding of errors computed between each

100     pair of stages in the hierarchy (e.g. a face was predicted, but a non-face was present leading to an

101     error).

102

103     **RESULTS**

104     We leveraged the hierarchically arranged face processing system in macaque ventral visual cortex to

105     study the dynamics of neural processing across a hierarchy (Tsao, Freiwald, Tootell, & Livingstone,

106     2006)(Tsao, Moeller, & Freiwald, 2008) (**Figure 1A**). The serially arranged posterior, central, and

107     anterior face-selective subregions of IT (pIT, cIT, and aIT) can be conceptualized as building increasing

108     selectivity for faces culminating in aIT representations (Freiwald & Tsao, 2010)(Chang & Tsao, 2017).

109     Using serial, single electrode recording, we sampled neural sites across the posterior to anterior extent

110     of the IT hierarchy in the left hemispheres of two monkeys to generate neurophysiological maps

111    (**Figure 1A**; example neurophysiological map in one monkey using a faces versus non-face objects

112    screen set) (Issa, Papanastassiou, & DiCarlo, 2013). We localized the recording locations *in vivo* and

113    co-registered across all penetrations using a stereo microfocal x-ray system (~400 micron *in vivo*

114    resolution) (Cox, Papanastassiou, Oreper, Andken, & DiCarlo, 2008)(Issa, Papanastassiou, Andken, &

115    DiCarlo, 2010) allowing accurate assignment of sites to different face processing stages (n = 633 out of

116    1891 total sites recorded were assigned as belonging to a face-selective subregion based on their

117    spatial location; see **Methods**). Results are reported here for sites that were spatially located in a face-

118    selective subregion, that showed visual drive to any category in the screen set (see **Methods**), and that

119    were subsequently tested with our face versus non-face challenge set (**Figure 1B**, left panel) (n = 115

120    pIT, 70 cIT, and 40 aIT sites).

121

122        Our experimental design was intended to test previously proposed computational hypotheses of

123    hierarchical neural dynamics during visual face processing (**Figure 1B**). Briefly, these hypotheses

124    predict how stimulus preference (in this instance, for faces versus non-faces) might change over time in

125    a neural population (**Figure 1B**, middle panel): (1) simple spike-rate adaptation predicts that initial rank-

126    order selectivity (i.e. relative stimulus preference) will be largely preserved (**Figure 1B**, dashed line)

127    while neurons adapt their absolute response strength over time, (2) local normalization predicts that

128    stronger responses are in some cases normalized to match weaker responses based on population

129    activity to specific dimensions (Carandini et al., 1997); importantly, normalization is strongest for

130    nuisance (non-coding) dimensions (e.g. low versus high stimulus contrast) and in its idealized form

131    would not alter selectivity along coding dimensions (e.g. face versus non-face) (**Figure 1B**, dashed

132    line), (3) evidence accumulation through temporal integration, winner-take-all through recurrent

133    inhibition, or Bayesian inference through top-down feedback mechanisms all predict increasing

134    selectivity for faces over time (Lee & Mumford, 2003) (**Figure 1B**, light gray line), and (4) predictive

135    coding posits that, for neurons that are coding error, their responses would show increasing activity for

5

136    non-faces (images whose properties are inconsistent with predictions of a face) and thus decreasing

137    face selectivity over time (Rao & Ballard, 1999) (**Figure 1B**, black line). Note, that error signaling is a

138    qualitatively different computation than normalization, as error coding predicts a decreased response

139    along the coding dimension (face versus non-face) whereas normalization would ideally not affect

140    selectivity for faces versus non-faces and only affect variation along orthogonal, nuisance dimensions.

141    Properly testing these predictions (no change in face selectivity, increased face selectivity, decreased

142    selectivity to be anti-face preferring) requires measurements from the intermediate stages of the

143    hierarchy as all of these models operate under the premise that the system builds and maintains a

144    preference for faces at the top of the hierarchy (**Figure 1B**, right, and see Introduction). Thus, the

145    intermediate stages (here pIT, see **Figure 1B**) are most likely to be susceptible to face/non-face

146    confusion and thus be influenced by, for example, the top-down mechanisms posited in Bayesian

147    inference and predictive coding where higher areas encode the face predictions that directly influence

148    the responses of lower areas (Lee & Mumford, 2003)(Rao & Ballard, 1999).

149

150    **Face and non-face images driving similar initial responses in pIT**

151    Here, we chose to focus our key, controlled tests on pIT – an intermediate stage in the ventral stream

152    hierarchy, but the first stage within IT where neural specialization for face detection (i.e. face vs. non-

153    face) has been reported (Grimaldi et al., 2016). Consistent with its intermediate position in the ventral

154    visual system, we had previously found that pIT face-selective neurons are not truly selective for whole

155    faces but respond to local face features, specifically those in the eye region (Issa & DiCarlo, 2012).

156    Taking advantage of this prior result, we created both face and non-face stimuli that challenged the

157    face processing system by strongly driving pIT responses, thus forcing the higher IT stages to complete

158    the discrimination between face and challenging non-face images. To generate ambiguous face-like

159    images, we systematically varied the positions of parts, in particular the eye, within the face (Issa &

160    DiCarlo, 2012) (see **Methods**). This set included images that contained face parts in positions

6

161    consistent with a frontal view of a face or images that only differed in the relative spatial configuration of

162    the face parts within the face outline (**Figure 1B**, left). Of the 82 images screened, we identified 21 part

163    configurations that each drove the pIT population response to $\geq$90% of its response to a correctly

164    configured whole face. Of those 21 images, 13 images were inconsistent with the part configuration of a

165    frontal face (**Figure 1B**, black box). For the majority of the results that follow, we focus on comparing

166    the neural responses to these 13 pIT-matched images that could *not* have arisen from frontal faces

167    (referred to hereafter as "non-face images") with the 8 images that could have arisen from frontal faces

168    (referred to hereafter as "face images"). Again, we stress that these two groups of images were

169    selected to evoke closely matched initial pIT population activity.

170

171        Importantly, the pIT-matched images used here presented a more stringent test of face vs. non-

172    face discrimination than prior work. Specifically, most prior work used images of faces and non-face

173    objects ("classic images") that contain differences across multiple dimensions including local contrast,

174    spatial frequency, and types of features (Tsao et al., 2006)(Afraz, Kiani, & Esteky, 2006)(Moeller,

175    Crapse, Chang, & Tsao, 2017)(Sadagopan, Zarco, & Freiwald, 2017). Consistent with this, we found

176    that the population decoding classification accuracy of our recorded neural populations using these

177    classic images (faces versus non-face objects) is near perfect (>99% in pIT, cIT, and aIT, n=30 sites

178    per region). However, we found that population decoding classification accuracy for the pIT-matched

179    face vs. non-face images we used here was near chance level (50%) in pIT (**Figure 1C**, blue bar; by

180    comparison, classification accuracy for face versus non-face objects classification was 99.6% using the

181    same pIT sites). Further downstream in regions cIT and aIT, we found that the linear population

182    decoding classification of these pIT-matched face vs. non-face images was well above chance,

183    suggesting that our pIT-matched face detection challenge is largely solved somewhere between pIT

184    and aIT (**Figure 1C**).

185

**Time course of responses in pIT for images with face versus non-face arrangements of parts**

We next closely examined the pIT neural response dynamics. To do this, we defined a face preference

value (d'; see Methods) that measured each site's average selectivity for the face images relative to the

non-face images, and we asked how a given site's preference evolved over time (see alternative

hypotheses in **Figure 1B**). First, we present three example sites which were chosen based on having

the largest selectivity (absolute d') in the late phase (100-130 ms post image onset). In particular, most

standard interpretations of face processing would predict a late phase preference for faces (d' > 0).

However, all three sites with the largest absolute d' had evolved a strong late phase preference for the

non-face images (d' < 0) despite having had very similar rising edge responses to the face and non-

face stimulus classes (response in early phase from 60-90 ms) (**Figure 2**, left column). A late, non-face

preference was not restricted to the example sites as a majority of pIT sites (66%) preferred non-faces

over faces in the late response phase (prefer frontal face arrangement: 60-90 ms = 66% vs. 100-130

ms = 34%; p = 0.000, n = 115) (**Figure 3B**, blue bars).

Next, we examined the dynamics of face selectivity across the pIT population as this is key to

disambiguating among the competing models outlined earlier (**Figure 1B**). In the adaptation and

normalization models, we would expect no change in the average population face selectivity, and the

evidence accumulation, winner-take-all, or Bayesian inference models predict an increase in face

selectivity over the population over time. Instead, we found that many sites significantly decreased their

face preference over time similar to the three example sites. Of the 51 sites in our pIT sample that

showed a significantly changing preference over time (p < 0.01 criterion for significant change in d'),

84% of these sites showed a decreasing preference (n = 43 of 51 sites, $p < 10^{-6}$, binomial test, n =

115 total sites) (**Figure 3A**, left column, light gray vs. black lines). This surprising trend -- decreasing

face preference -- was strong enough that it erased any small, initial preference for images of frontally

arranged face parts over the population (median d': 60-90 ms = 0.11 $\pm$ 0.02 vs. 100-130 ms = -0.12 $\pm$

0.03, p = 0.000, n=115 sites), and this trend was observed in both monkeys when analyzed separately

8

211    ($p_{M1}$ = 0.000, $p_{M2}$ = 0.002, $n_{M1}$ = 43, $n_{M2}$ = 72 sites; **Figure 4A**). This decreasing face selectivity over

212    time was driven by decreasing firing rates to the face images containing normally arranged face parts.

213    Responses to these images were weaker by 18% on average in the late phase of the response

214    compared to the early phase ($\Delta$rate (60-90 vs 100-130 ms) = -18% $\pm$ 4%, p = 0.000; n = 7 images)

215    while responses to the non-face images with atypical spatial arrangements of face parts -- also capable

216    of driving high early phase responses -- did not experience any firing rate reduction in the late phase of

217    the response ($\Delta$rate (60-90 vs 100-130 ms) = 2 $\pm$ 1%, p = 0.467; n = 13 images).

218

219        The above observation of decreasing face preference over the pIT population seemed most

220    consistent with the prediction of error coding models, but one potential confound was that initial

221    responses to faces and challenging non-faces were not perfectly matched across the population (recall

222    that we only required face and non-face images to drive a response $\geq$90% of the whole face response).

223    As a result, initial selectivity was non-zero (d' = 0.11, n=115 sites). This residual face preference may

224    be small, but if this residual face selectivity is driven by nuisance dimensions, for example excess

225    stimulus contrast in the face class relative to the non-face class, then the face class may have

226    experienced stronger activity dependent adaptation or normalization resulting in a decreasing face

227    preference over time. To more adequately limit general activity dependent mechanisms leading to

228    decreasing face responses, we performed control analyses where initial activity was tightly matched per

229    site or where the number of parts were matched across images.

230

231    **Controls in pIT for firing rate and low-level image variation**

232    To strictly control for the possibility that simple initial firing rate differences could predict the observed

233    phenomenon, we re-computed selectivity after first matching initial responses site-by-site. For this

234    analysis, images were selected on a per site basis to evoke similar initial firing rates (images driving

9

235    initial response within 25% of synthetic whole face response for that site, at least 5 images required per

236    class). This image selection procedure virtually eliminated any differences in initial responses between

237    the face and non-face image classes and hence any firing rate difference driven by differences in

238    nuisance parameters between faces and challenge non-faces (**Figure 3C**, 60-90 ms), yet we still

239    observed a significant drop in preference for images with typical frontal face part arrangements versus

240    atypical face part arrangements in pIT ($\Delta_{d'}$ = -0.10 $\pm$ 0.03, p = 0.001, n = 77) (**Figure 3C**, blue line).

241    Thus, the remaining dependence of firing rate dynamics on the image class and not on initial response

242    strength argued against an exclusively activity based explanation to account for decreasing neural

243    responses to faces over time. Further arguing against this hypothesis, we found that the pattern of late

244    phase population firing rates in pIT across images could not be significantly predicted from early phase

245    pIT firing rates for each image ($\rho_{pIT\ early,\ pIT\ late}$ = 0.07 $\pm$ 0.17, p = 0.347; n = 20 images).

246         Thus far, we have performed analyses where images from the face and non-face class were

247    similar in their initially evoked response which equated images at the level of neural activity but

248    produced images varying in the number of parts. An alternative is to match the number of face parts

249    between the face and non-face classes as another means of limiting the differences in nuisance

250    dimensions such as the contrast, spatial frequency and retinal position of energy across images (see

251    examples in **Figure 4B**). When we recomputed selectivity across subsets of images containing a

252    matched number of one, two, or four parts (n=5, 30, and 3 images, respectively), we still observed that

253    pIT face selectivity decreased. For all three image subsets controlling the number of face parts, d'

254    began positive on average in the sampled pIT population (i.e. preferred frontal face part arrangements

255    in 60-90 ms post-image onset) (median d' for 60-90 ms = 0.13 $\pm$ 0.05, 0.05 $\pm$ 0.02, 0.33 $\pm$ 0.09 for one,

256    two, and four parts) and significantly decreased in the next phase of the response (100-130 ms post-

257    image onset) becoming negative on average (median d' for 100-130 ms: -0.27 $\pm$ 0.06, -0.14 $\pm$ 0.02, -

258    0.04 $\pm$ 0.12; one, two, four parts: p = 0.000, 0.000, 0.004, for d' comparisons between 60-90 ms and

10

259    100-130 ms, n = 115, 115, 76 sites) (**Figure 4B**). A similar decreasing face selectivity profile was

260    observed when we re-tested single part images at smaller ($3^o$) and larger ($12^o$) image sizes suggesting

261    a dependence on the relative configuration of the parts and not on their absolute retinal location or

262    absolute retinal size (median d' for 60-90 ms vs. 100-130 ms: three degrees = 0.51 $\pm$ 0.09 vs. -0.29 $\pm$

263    0.14, twelve degrees = 0.07 $\pm$ 0.14 vs. -0.11 $\pm$ 0.14; n = 15; p = 0.000, 0.025, 0.07) (**Figure 4C**). Thus,

264    we suggest that the decreasing population face selectivity dynamic in pIT is a fairly robust phenomenon

265    specific to the face versus non-face dimension as this dynamic persists even when limiting potential

266    variation across nuisance dimensions. This phenomenon can be distinguished from normalization

267    mechanisms that would operate most strongly to reduce selectivity along nuisance dimensions rather

268    than along dimensions that directly solve the face versus non-face discrimination challenge.

269

270    **Time course of responses in aIT and cIT for images with face versus non-face arrangements of**

271    **parts**

272    Under the possibility that decreasing face selectivity in intermediate stage pIT is a signature of error

273    signaling, we next asked whether the source of the prediction signal could be observed in higher

274    cortical areas. In the anterior face-selective regions of IT which are furthest downstream of pIT and

275    reflect additional stages of feedforward processing (see block diagram in **Figure 1B**), we did not

276    observe the strong decreases in the selectivity profile seen in pIT. Indeed, the three sites with the

277    greatest selectivity (absolute d') in the late response phase (100-130 ms) in our aIT sample all

278    displayed a preference for frontal face part arrangements (d' > 0) (**Figure 2**, right column). Also, in

279    contrast to the dynamic selectivity profiles observed in many pIT sites, 98% of aIT sites (39 of 40) did

280    not significantly change their relative preference for face vs. non-face arrangements of the parts (p <

281    0.01 criterion for significant change at the site level) (**Figure 3A,** right column, bottom row, dark gray

282    sites). Rather, we observed a stable selectivity profile over time in aIT (median d': 60-90 ms = 0.13 $\pm$

283    0.03 vs. 100-130 ms = 0.17 $\pm$ 0.03, p = 0.34, n=40 sites). As a result, the majority of anterior sites

11

284     preferred images with typical frontal arrangements of the face parts in the late phase of the response

285     (prefer face: 60-90 ms = 78% of sites vs. 100-130 ms = 78% of sites; p = 0.451, n = 40 sites; **Figure**

286     **3B**, red bars) despite only a minority (34%) of upstream sites in pIT preferring these images in their late

287     response. Thus, spiking responses of individual aIT sites were as expected from a computational

288     system whose purpose is to detect faces, as previously suggested (Freiwald & Tsao, 2010), and aIT

289     serves as one candidate for the putative prediction signal underlying face prediction errors.

290

291        In cIT whose anatomical location is intermediate to pIT and aIT, we observed many sites with

292     decreasing selectivity (**Figure 2C & 3A**, middle columns), a dynamic that persisted even when we

293     tightly matched initial responses on a site by site basis similar to pIT (**Figure 2C**, green line). The

294     overall stimulus preference in cIT was intermediate to that of pIT and aIT (**Figure 3B**) consistent with

295     the intermediate position of cIT in the IT hierarchy. Interestingly, we found that the patterns of

296     responses across images in the early phases of cIT and aIT activity were significant predictors of late

297     phase activity in pIT ($\rho_{cIT\ early,\ pIT\ late}$ = -0.52 $\pm$ 0.11, p = 0.000; $\rho_{aIT\ early,\ pIT\ late}$ = -0.36 $\pm$ 0.14, p = 0.012;

298     $n_{pIT}$=115, $n_{cIT}$=70, $n_{aIT}$=40 sites; n = 20 images), even better predictors than early phase activity in pIT

299     itself ($\rho_{pIT\ early,\ pIT\ late}$ = 0.07 $\pm$ 0.17, p = 0.347). That is, for images that produced high early phase

300     responses in cIT and aIT, the following later phase responses of units in the lower level area (pIT)

301     tended to be low, consistent with error coding models which posit that feedback from higher areas (in

302     the form of predictions) would contribute to the decreasing selectivity observed in lower areas.

303

304     **Computational models of neural dynamics in IT**

305     We next proceeded to formalize the conceptual ideas introduced in **Figure 1B** and build neurally

306     mechanistic, dynamical models of gradually increasing complexity to determine the minimal set of

307     assumptions that could capture our empirical findings of non-trivial, dynamic selectivity changes during

308    face detection across face-selective subregions in IT. Previous functional and anatomical data show

309    that the face-selective subregions in IT are connected forming an anterior to posterior hierarchy and

310    show that pIT serves as the primary input into this hierarchy (Moeller et al., 2008)(Freiwald & Tsao,

311    2010)(Grimaldi et al., 2016). Thus, we evaluated dynamics in different hierarchical architectures using a

312    linear dynamical systems modeling framework where pIT, cIT, and aIT act as sequential stages of

313    processing (**Figure 5** and see **Methods**). A core principle of feedforward ventral stream models is that

314    object selectivity is built by stage-wise feature integration in a manner that leads to relatively low

315    dimensional representations at the top of the hierarchy abstracted from the high-dimensional input

316    layer. We were interested in how signals temporally evolve across a similar architectural layout. We

317    used the simplest feature integration architecture where a unit in a downstream area linearly sums the

318    input from units in an upstream area, and we stacked this computation to form three layer networks

319    (**Figure 5B**). This simple, generic feedforward encoding model conceptualizes the idea that different

320    types of evidence, local and global (i.e. information about the parts and the relative spatial arrangement

321    of parts), have to converge and be integrated to separate face from non-face images in our image set.

322    We used linear networks as monotonic nonlinearities can be readily accommodated in our framework

323    (Seung, 1997)(Rao & Ballard, 1999)(also see **Figure 7C**). Importantly, we used a simple encoding

324    scheme as our goal was not to build full-scale deep neural network encoding models of image

325    representation (Yamins et al., 2014) but to bring focus to an important biological property that is often

326    not considered in deep nets, neural dynamics.

327        We implemented a range of ideas previously proposed in the literature. The functional

328    implications of these ideas were highlighted in **Figure 1B**, but at a mechanistic level, these functional

329    properties can be directly realized via different recurrent processing motifs between neurons (**Figure 5**,

330    base feedforward model (first column) is augmented with recurrent connections to form new models

331    (remaining columns)). For example, self-connections can be viewed as implementing spike rate

332    adaptation in a feedforward architecture (**Figure 5A**, top row), lateral connections support winner-take-

333    all or normalization mechanisms (Carandini et al., 1997) (**Figure 5A**, second, third, and fourth

334    columns), and top-down connections implement Bayesian inference (Seung, 1997) (**Figure 5A**, fifth

335    and sixth columns). Here, normalization is implemented by a leak term that scales adaptation (the

336    degree of decay in the response) and is controlled recursively by the summed activity of the network

337    (Carandini et al., 1997), ultimately scaling down responses to strong driving stimuli over time. To

338    constrain our choice of a feedback-based model, we took a normative approach minimizing a quadratic

339    reconstruction cost between stages as the classical reconstruction cost is at the core of an array of

340    hierarchical generative models including hierarchical Bayesian inference (Lee & Mumford, 2003),

341    Boltzmann machines (Ackley, Hinton, & Sejnowski, 1985), analysis-by-synthesis networks (Seung,

342    1997), sparse coding (Olshausen & Field, 1996), predictive coding (Rao & Ballard, 1999), and

343    autoencoders in general (Rifai, Vincent, Muller, Glorot, & Bengio, 2011). Optimizing a quadratic loss

344    results in feedforward and feedback connections that are symmetric -- reducing the number of free

345    parameters -- such that inference on the represented variables at any intermediate stage is influenced

346    by both bottom-up sensory evidence and current top-down interpretations. Critically, a common feature

347    of this large model family is the computation of between-stage error signals via feedback, which is

348    distinct from state-estimating model classes (i.e. feedforward models) that do not compute or propagate

349    errors. A dynamical implementation of such a network uses leaky integration of error signals which, as

350    shared computational intermediates, guide gradient descent of the values of the represented variables

351    to a previously learned target value (Δactivity of each neuron => online inference) or descend the

352    connection weights to values that give the best future behavior (Δsynaptic strengths => offline learning),

353    here defined as an unsupervised reconstruction goal (similar results were found using other goals and

354    networks such as supervised discriminative networks; see **Figure 7C**).

355        When we fit each of the models to our neural data, they generally produced an increase in

356    selectivity from the first stage of the network to the later stages of the network. This increase is not

357    surprising because the models had built-in converging feedforward connections from the first to second

358    to third stages (**Figure 6A**, first five columns, compare blue to red curves). However, discrepancies

359    between models emerged in the lower stages where we found that neither the lateral inhibition model,

360    nor the normalization model, could capture the decreasing selectivity phenomenon observed in pIT.

361    Instead, the selectivity of these models simply increased to a saturation level set by the leak term

362    (shunting inhibition) in the system (**Figure 6A**, first five columns). Similar behavior was present when

363    we tried a nonlinear implementation of the normalization model that more powerfully modulated

364    shunting inhibition (Carandini et al., 1997). That normalization proved insufficient to generate the

365    observed neural dynamics can be explained by the fact that the normalized response to a stimulus

366    cannot easily fall below the response to a stimulus that was initially similar in strength. Thus, a

367    decreasing average preference for a stimulus across a population of cells (i.e. **Figures 2-4**, pIT data)

368    for similar levels of average input is difficult when only using a basic normalization model mediated by

369    surround (within-stage) suppression.

370

371       In contrast to the above models, we found that the feedback model capable of computing

372    hierarchical error signals naturally displayed a strong decrease of selectivity in a sub-component of its

373    first processing stage -- qualitatively similar behavior to the selectivity decrease that we observed in

374    many pIT and cIT neural sites. Specifically, this model displayed these dynamics in the magnitude of its

375    reconstruction error signals but not in its state signals (the feature values) (**Figure 6A,** compare fifth

376    and sixth columns). These error signals integrate converging state signals from two stages -- one

377    above (prediction) and one below (sensory evidence). The term "error" is thus meaningful in the hidden

378    processing stages where state signals from two stages can converge. The top nodes of a hierarchy

379    receive little descending input and hence do not carry additional errors with respect to the desired

380    computation; rather, top nodes convey the face predictions that influence errors in bottom nodes. This

381    behavior in the higher processing stages is consistent with our observation of explicit representation of

382    faces in aIT in all phases of the response (**Figures 2-3**) and with similar observations of decodable

383    identity signals by others in all phases of aIT responses for faces (Meyers, Borzello, Freiwald, & Tsao,

384    2015) and objects (Hung et al., 2005)(Majaj et al., 2015). We also found similar error dynamics when

385    using a simpler two-layer network as opposed to three layers suggesting that these error signal

386    dynamics along with prediction signals emerge even in the simplest cascaded architecture (**Figure 7A**).

387

388    **Predictions of an error coding hierarchical model**

389    While our neural observations at multiple stages of the IT hierarchy led us to the error coding

390    hierarchical model above, a stronger test of the idea of error signaling is whether it predicts other IT

391    neural phenomena. To identify stimulus regimes that would lead to insightful predictions, we asked in

392    what way would the behavior of error-estimating hierarchical models differ most from the behavior of

393    generic feedforward state-estimating models. Because our feedback-based model uses feedforward

394    inference at its core, it behaves similarly to a state-estimating hierarchical feedforward model when the

395    statistics of inputs match the learned feedforward weight pattern of the network (i.e. 'natural' images

396    drawn from everyday objects and scenes) since for these inputs, feedforward inferences derived from

397    the sensory data are aligned with top-down expectations. Thus, predictions of feedback-based models

398    that could distinguish them from feedforward-only models are produced when the natural statistics of

399    images are altered so that they differ from the feedforward patterns previously learned by the network

400    and hence differ from the predictions generated in the network. We have (above) considered one such

401    type of alteration: images where local face features are present but altered from their naturally

402    occurring (i.e. statistically most likely) arrangement in frontal faces. Next, we tested two other image

403    manipulations from recent physiology studies which yielded novel neural phenomena that lacked a

404    principled, model-based explanation (Freiwald, Tsao, & Livingstone, 2009)(Meyer et al., 2014). To test

405    whether the error coding hierarchical model family displays these behaviors, we fixed the architectural

406    parameters derived from our fitting procedure in **Figure 6** and simply varied the input to this network,

407    specifically the correlation between the inputs and the network's feedforward weight pattern, in order to

408    match the nature of the image manipulations performed in prior experiments.

409

410    *Sublinear integration of the face features.* In the face-selective subregion of cIT, the sum of the

411    responses to the face parts exceeds the response to the whole face when averaging firing rates over a

412    200 ms window (Freiwald et al., 2009; see their Figure 2C), and we observed a similar phenomenon in

413    the face-selective subregion in pIT (Issa & DiCarlo, 2012). Interestingly, examining the dynamics of part

414    integration more closely within the first 100ms of the response reveals that the response to the whole

415    face does begin at a level similar to the sum of the responses to the face parts but becomes sublinear

416    in the late phase (whole face response relatively low compared to linear prediction from parts

417    responses) (ratio of sum of responses to parts vs. response to whole: 60-90 ms = $1.5 \pm 0.1$, 100-130

418    ms = $4.6 \pm 0.3$; p = 0.000, n = 33 sites) (**Figure 8A**, left panel). This result runs counter to what would

419    be expected in a model where selectivity for the whole face is built from the conjunction of the parts. In

420    such a model, the population response to the whole face would be at least as large if not greater

421    (superlinear) than the summed responses to the individual features. To test whether an error coding

422    model exhibited the phenomenon of sublinear feature integration at the population level, we compared

423    the response with all inputs active (co-occurring features) to the sum of the responses when each input

424    was activated independently (individual features). The reconstruction errors in our feedback-based

425    model showed a strong degree of sublinear integration of the inputs such that the response to the

426    simultaneous inputs (whole) was much smaller than what would be predicted by a linear sum of the

427    responses to each input alone (parts), and the model's sublinear integration behavior qualitatively

428    replicated the time course observed in pIT without any additional fitting of parameters (**Figure 8A**, right

429    panel). Although we certainly expect that similar sublinear integration may also be observed in a

430    normalization model, this would require a particularly strong form of normalization since population

431    activity to the whole face would have to be normalized to nearly the same level as that for an individual

432    part in the late response phase (ratio of response to single part vs. response to whole = $0.92 \pm 0.06$,

433    100-130 ms, n = 33 sites) despite being three times larger during the early response phase (ratio of

434    response to single part vs. response to whole = 0.3 $\pm$ 0.02, 60-90 ms, n = 33 sites). Furthermore, the

435    stacked normalization model could not account for the main dynamical phenomenon that we observed

436    (**Figure 6**); therefore, an error coding perspective may provide a more parsimonious account of the

437    observed set of dynamical phenomena.

438

439    *Evolution of neural signals across time.* Neural responses to familiar images are known to rapidly

440    attenuate in IT when compared to responses to novel images (Freedman, Riesenhuber, Poggio, &

441    Miller, 2006)(Woloszyn & Sheinberg, 2012)(Meyer & Olson, 2011). This observation seems to

442    contradict what would be predicted by simple Hebbian potentiation for the more exposed stimuli.

443    Furthermore, familiar image responses show much sharper temporal dynamics than responses to novel

444    images when presented repeatedly (Meyer et al., 2014). These qualitatively different dynamics for

445    familiar versus novel images are surprising given that stimuli are drawn from the same distribution of

446    natural images and are thus matched in their average image-level statistical properties (color, spatial

447    frequency, contrast). To test whether our network displayed these different dynamical behaviors, we

448    simulated familiar inputs as those that match the learned weight pattern of a high-level detector and

449    novel inputs as those with the same overall input level but with weak correlation to the learned network

450    weights (here, we have extended the network to include two units in the output stage corresponding to

451    storage of the two familiarized input patterns to be alternated; conceptually, we consider these familiar

452    pattern detectors as existing downstream of IT in a region such as perirhinal cortex which has been

453    shown to code familiarized image statistics and memory-based object signals (Murray, Bussey, &

454    Saksida, 2007)). We repeatedly alternated two familiar inputs or two novel inputs and found that error

455    coding model responses in the hidden processing stage were temporally sharper for familiar inputs that

456    matched the network's feedforward weight patterns compared to novel patterns of input (pseudo-

457    randomly drawn; see **Methods**), consistent with the previously observed phenomenon (**Figure 8B**; data

18

458    reproduced with permission from Meyer et al., 2014). Model responses reproduced additional details of

459    the neural dynamics including a large initial peak followed by smaller peaks for responses to novel

460    inputs and a phase delay in the oscillations of responses to novel inputs compared to familiar inputs.

461    Intuitively, these dynamics are composed of two phases after the initial response transient to the onset

462    of the image sequence. In the first phase, familiar patterns lead to lower errors and hence lower neural

463    responses than random patterns (**Figure 8B**, red curve drops below the black curve after the onset

464    response), similar to the observed weaker response to more familiar face-like images present in our

465    data (**Figure 2**, red curves drop below black curves in pIT example sites). When the familiar pattern *A*

466    is switched to another familiar pattern *B*, this induces a short-term error in adjusting to the new pattern

467    (**Figure 8B**, red curve briefly goes above the black curve during pattern switch and then decreases). In

468    contrast, two unfamiliar patterns are closer together in the high-level encoding space than two learned

469    patterns (**Figure 8B**, inset at right), and the switch between two unlearned patterns introduces relatively

470    less shift in top-down signals and hence a smaller dynamical change in error signals. This result

471    demonstrates that our model, derived from fitting only the first 70 ms (60-130 ms post image onset) of

472    IT responses to face images, can extend to much longer timescales and may generalize to studies of

473    images besides face images.

474

475    **Dynamical properties of neurons across cortical lamina**

476    In the large family of state-error coding hierarchical networks, a number of different cortical circuits are

477    possible. A key distinction of two such circuit mapping hypotheses (predictive coding versus error

478    backpropagation) is the expected laminar location of state coding neurons transmitting information

479    about features in the image. In typical neural network implementations, the feedforward projecting

480    neurons in superficial lamina are presumed to encode estimates about states of the visual world (e.g.

481    presence of a face). In contrast, predictive coding posits that superficial layers contain error units and

482    that errors are projected forward to the next cortical level (Rao & Ballard, 1999)(Friston & Kiebel,

19

483    2009)(Hyvärinen, Hurri, & Hoyer, 2009). State and error signals can be distinguished by their dynamical

484    signatures in our leading model, which was fit on error signals but produces predictions of the

485    corresponding state signals underlying the generation of errors. Since state units are integrators (see

486    **Methods**), they have slower dynamics than error units leading to longer response latencies and a

487    milder decay in responses (**Figure 8C**, red curve in right panel). To test this prediction, we localized our

488    recordings relative to the cortical mantle by co-registering the x-ray determined locations of our

489    electrode (~400 micron *in vivo* accuracy; (Issa et al., 2010)) to structural MRI data (see **Methods**).

490    When we separated units into those at superficial depths closer to the pial surface (1/3 of our sites;

491    corresponds to approximately 0 to 1 mm in depth) versus those in the deeper layers (remaining 2/3 of

492    sites, ~1 to 2.5 mm in depth), we found a longer latency and less response decay in superficial units

493    consistent with the expected profile of state units (**Figure 8C**, left panel). Thus, the trend toward state-

494    like signals in superficial layers is more consistent with typical error backpropagation models (states fed

495    forward, errors fed backward) than with predictive coding proposals. In fact, the latency difference

496    between cortical lamina in pIT (deep vs superficial: $66.0 \pm 1.7$ vs $76.0 \pm 1.8$ ms, p = 0.002) was greater

497    than the conduction delay from pIT to cIT (i.e. from superficial layers of pIT to the deeper layers of cIT)

498    (superficial pIT vs deep cIT: $76.0 \pm 1.8$ vs $75.5 \pm 2.0$ ms, p = 0.15) even though laminar distances

499    within pIT are smaller than the distance traveled between cortical stages pIT and cIT. Thus, instead of a

500    simple conduction delay accounting for latency differences across lamina, our model suggests that

501    temporal integration of inputs in superficial lamina, more consistent with the behavior of state units as

502    opposed to error units, may drive the lagged dynamical properties of neurons in superficial lamina.

503

504    **DISCUSSION**

505    We have measured neural responses during a difficult frontal face detection task across the IT

506    hierarchy and demonstrated that the population preference for faces in the intermediate (a.k.a hidden)

507    processing stages decreases over time – that is population responses at lower levels of the hierarchy

508    (pIT and cIT) rapidly evolved toward preferring unnatural, non-face part arrangements whereas the top

509    level (aIT) rapidly developed and then maintained a preference for natural, frontal face part

510    arrangements. The relative speed of selectivity changes in pIT (~30 ms) makes high-level explanations

511    based on fixational eye movements or shifts in attention (e.g. from behavioral surprise to unnatural

512    arrangements of face parts) unlikely as saccades and attention shifts occur on slower timescales

513    (hundreds of milliseconds) (Egeth & Yantis, 1997). The presence of stronger responses to face images

514    than to non-face images in aIT further argues against general arousal effects as these would have been

515    expected to cause stronger responses to challenging non-face images in aIT. Rather, the rapid

516    propagation of neural signals over tens of milliseconds suggested intracortical processing within the

517    ventral visual stream in a manner that was not entirely consistent with a pure feedforward model, even

518    when we included strong nonlinearities in these models such as normalization and even when we

519    stacked these operations to form more complex three stage models. However, augmenting the

520    feedforward model so that it represented the errors generated during hierarchical processing produced

521    the observed neural dynamics and hierarchical signal propagation (**Figures 6-7**). This view argues that

522    many IT neurons code error signals. Using this new modeling perspective, we went on to generate

523    predictions of previously observed IT neural phenomena (**Figure 8**).

524

525    **Comparison to previous neurophysiology studies in IT**

526    Our suggestion that many IT neurons code errors is consistent with the observation of strong

527    responses to extremes in face space (Leopold, Bondar, & Giese, 2006) providing an alternative

528    interpretation to the prior suggestion that cIT neurons are not tuned for typical faces but are instead

529    tuned for atypical face features (i.e. extreme feature tuning) (Freiwald et al., 2009). In that prior work,

530    the response preference of each neuron was determined by averaging over a long time window (~200

531    ms). By looking more closely at the fine time scale dynamics of the IT response, we suggest that this

532    same extreme coding phenomenon can instead be interpreted as a natural consequence of networks

533 that have an actual tuning preference for typical faces (as evidenced by an initial response preference

534 for typical frontal faces in pIT, cIT, and aIT; **Figure 3B**) but that also compute error signals with respect

535 to that preference. Under the present hypothesis, some IT neurons are preferentially tuned to typical

536 spatial arrangements of face features, and other IT neurons are involved in coding errors with respect

537 to those typical arrangements. We speculate that these intermixed state estimating and error coding

538 neuron populations are both sampled in standard neural recordings of IT, even though only state

539 estimating neurons are truly reflective of the tuning preferences of that IT processing stage. This

540 intermixed view of IT neural signaling is further supported by recent studies demonstrating correlates of

541 temporal prediction errors to image sequences in IT (Meyer & Olson, 2011) including the face-selective

542 subregion of cIT (Schwiedrzik & Freiwald, 2017).

543

544 The precise fractional contribution of errors to total neural activity is difficult to estimate from our

545 data. Under the primary image condition tested, not all sites significantly decreased their selectivity

546 (~60% did not change their selectivity). We currently interpret these sites as coding state (feature)

547 estimates (**Figure 3A**, light and dark gray lines in top and bottom rows, respectively), and we did

548 observe evidence of emergence of state-like signals in our superficial neural recordings (**Figure 8C**).

549 Alternatively, at least some of the non-reversing sites might be found to code errors under other image

550 conditions than the one that we tested. Furthermore, while in our primary image condition selectivity

551 decreases only accounted for ~15% of the overall spiking modulation (**Figure 6A**, data panel), larger

552 modulations in late phase neural firing (50-100%) are possible under other image conditions (**Figure

553 8A,B**). At a computational level, the absolute contribution of error signals to spiking may not be the

554 critical factor as even a small relative contribution may have important consequences in the network.

555

556 **Comparison across dynamical models of neural processing**

557 Our goal was to test a range of existing recurrent models by recording neural dynamics across multiple

22

558   cortical stages which provided stronger constraints on computational models than fitting neural

559   responses from only one area as in prior work (Carandini et al., 1997)(Rao & Ballard, 1999). Crucially,

560   we found that the multi-stage neural dynamics observed in our data could not be adequately fit by only

561   using recurrences within a stage such as adaptation, lateral inhibition, and standard forms of

562   normalization (**Figure 6**). These results did not change when we made our simple networks more

563   complex by adding more stages (compare **Figure 6** versus **Figure 7**) or by using more realistic model

564   units with monotonic nonlinearities similar to a spiking nonlinearity (data not shown). Indeed, we

565   specifically chose our stimuli to evoke similar levels of within stage neural activity to limit the effects of

566   known mechanisms that depend on activity levels within an area (e.g. adaptation, normalization), and

567   we fully expect that these activity dependent mechanisms would operate in parallel to top-down,

568   recurrent processes during general visual processing. We emphasize that we only tested the standard

569   form of normalization as originally proposed, using within stage pooling and divisive mechanisms

570   (Carandini et al., 1997). Since that original mechanistic formulation, normalization has evolved to

571   become a term that broadly encapsulates many forms of suppression phenomena and can include both

572   lateral interactions within an area and feedback interactions from other areas (Nassi, Gómez-Laberge,

573   Kreiman, & Born, 2014)(Coen-Cagli, Kohn, & Schwartz, 2015). Thus, while our results do not follow

574   from the original mechanistic form of normalization, they may yet fall under normalization more broadly

575   construed as a term for suppression phenomena (error coding would require a similar suppressive

576   component). Here, we have provided a normative model for how top-down suppression would follow

577   from the well-defined computational goals of many hierarchical neural network models.

578

579   **Error signals generated across different hierarchical inference and learning models**

580   The notion of error is inherent to many existing models in the literature that go beyond the basic

581   feedforward, feature estimation class. Example models use errors for guiding top-down inference by

582   computing errors implicitly (hierarchical Bayesian inference (Seung, 1997)(Lee & Mumford, 2003)) or by

583    representing errors explicitly (predictive coding (Rao & Ballard, 1999)). In addition, errors can be used

584    specifically for driving unsupervised learning (autoencoder (Rifai et al., 2011)) or supervised learning

585    (classic error backpropagation (Williams & Hinton, 1986)). Recently, models have incorporated aspects

586    of both inference and learning (Salakhutdinov & Hinton, 2012)(Patel, Nguyen, & Baraniuk, 2015). A

587    key, unifying feature across inference and learning models is the need to compute an error signal

588    between processing stages. This error signal can be in the form of a generative, reconstruction cost

589    (stage *n* predicting stage *n-1*) or a discriminative, construction cost (stage *n-1* predicting stage *n*).

590    Across-stage "performance" error terms are used in all model cost functions, are typically the only term

591    combining signals from different model stages, and are distinct from within-stage "regularization" terms

592    (i.e. sparseness or weight decay) in driving network behavior (Marblestone, Wayne, & Kording, 2016).

593    The present study provides evidence that such errors are not only computed, but that they are explicitly

594    encoded in spiking rates. We emphasize that this result at the level of population dynamics was robust

595    across choices of cost function such as those used in the literature; we tested models with different

596    unsupervised and supervised performance errors (reconstruction, nonlinear reconstruction, and

597    discriminative) and found similar population level error signals across these networks in the basic two-

598    layer implementation (**Figure 7C**). Thus, errors as generally instantiated in the state-error coding

599    hierarchical model family provide a good approximation to IT population neural dynamics.

600

601    **Computational utility of coding errors in addition to states**

602    In error-computing networks, errors provide control signals for guiding learning giving these networks

603    additional adaptive power over basic feature estimation networks. This property helps augment the

604    classical, feature coding view of neurons which, with only feature activations and Hebbian operations,

605    does not lead to efficient learning in the manner produced by gradient descent using error

606    backpropagation (Williams & Hinton, 1986). Observation of error signals may provide insight into how

607    more intelligent unsupervised and supervised learning algorithms such as backpropagation could be

608    plausibly implemented in the brain. A potentially important contribution of this work is the suggestion

609    that gradient descent algorithms are facilitated by using an error code so that efficient learning is

610    reduced to a simple Hebbian operation at synapses and efficient inference is simply integration of

611    inputs at the cell body (see *eqn 10* and text in **Methods**). This representational choice, to code the

612    computational primitives of gradient descent in spiking activity, would simply leverage the existing

613    biophysical machinery of neurons for inference and learning.

614

615    **MATERIALS & METHODS**

616    **Animals and surgery.** All surgery, behavioral training, imaging, and neurophysiological techniques are

617    identical to those described in detail in previous work (Issa & DiCarlo, 2012). Two rhesus macaque

618    monkeys (*Macaca mulatta*) weighing 6 kg (Monkey 1, female) and 7 kg (Monkey 2, male) were used. A

619    surgery using sterile technique was performed to implant a plastic fMRI compatible headpost prior to

620    behavioral training and scanning. Following scanning, a second surgery was performed to implant a

621    plastic chamber positioned to allow targeting of physiological recordings to posterior, middle, and

622    anterior face patches in both animals. All procedures were performed in compliance with National

623    Institutes of Health guidelines and the standards of the MIT Committee on Animal Care and the

624    American Physiological Society.

625

626    **Behavioral training and image presentation.** Subjects were trained to fixate a central white fixation

627    dot during serial visual presentation of images at a natural saccade-driven rate (one image every 200

628    ms). Although a $4^o$ fixation window was enforced, subjects generally fixated a much smaller region of

629    the image ($<1^o$) (Issa & DiCarlo, 2012). Images were presented at a size of $6^o$ except for control tests at

630    $3^o$ and $12^o$ sizes (**Figure 4C**), and all images were presented for 100 ms duration with 100 ms gap

631    (background gray screen) between each image. Up to 15 images were presented during a single

632    fixation trial, and the first image presentation in each trial was discarded from later analyses. Five

633   repetitions of each image in the general screen set were presented, and ten repetitions of each image

634   were collected for all other image sets. The screen set consisted of a total of 40 images drawn from

635   four categories (faces, bodies, objects, and places; 10 exemplars each) and was used to derive a

636   measure of face versus non-face object selectivity. Following the screen set testing, some sites were

637   tested using an image set containing images of face parts presented in different combinations and

638   positions (**Figure 1B**, left panel). We first segmented the face parts (eye, nose, mouth) from a monkey

639   face image. These parts were then blended using a Gaussian window, and the face outline was filled

640   with pink noise to create a continuous background texture. A face part could appear on the outline at

641   any one of nine positions on an evenly spaced 3x3 grid. Although the number of possible images is

642   large ($4^9$ = 262,144 images), we chose a subset of these images for testing neural sites (n = 82

643   images). Specifically, we tested the following images: the original whole face image, the noise-filled

644   outline, the whole face reconstructed by blending the four face parts with the outline, all possible single

645   part images where the eye, nose, or mouth could be at one of nine positions on the outline (n = 3x9 =

646   27 images), all two part images containing a nose, mouth, left eye, or right eye at the correct outline-

647   centered position and an eye tested at all remaining positions (n = 4*8-1 = 31 images), all two part

648   images containing a correctly positioned contralateral eye while placing the nose or mouth at all other

649   positions (n = 2*8-2 = 14 images), and all correctly configured faces but with one or two parts missing

650   besides those already counted above (n = 4+3 = 7 images). The particular two-part combinations

651   tested were motivated by prior work demonstrating the importance of the eye in early face processing

652   (Issa & DiCarlo, 2012), and we sought to determine how the position of the eye relative to the outline

653   and other face parts was encoded in neural responses. The three and four part combinations were

654   designed to manipulate the presence or absence of a face part for testing the integration of face parts,

655   and in these images, we did not vary the positions of the parts from those in a naturally occurring face.

656   In a follow-up test on a subset of sites, we permuted the position of the four face parts under the

657   constraint that they still formed the configuration of a naturally occurring face (i.e. preserve the 'T'

658    configuration, n = 10 images; **Figure 4B**). We tested single part images at $3^o$ and $12^o$ sizes in a subset

659    of sites (n = 27 images at each size; **Figure 4C**). Finally, we measured the responses to the individual

660    face parts in the absence of the outline (n = 4 images; **Figure 8A**).

661

662    **MR Imaging and neurophysiological recordings.** Both structural and functional MRI scans were

663    collected in each monkey. Putative face patches were identified in fMRI maps of face versus non-face

664    object selectivity in each subject. A stereo microfocal x-ray system (Cox et al., 2008) was used to guide

665    electrode penetrations in and around the fMRI defined face-selective subregions of IT. X-ray based

666    electrode localization was critical for making laminar assignments since electrode penetrations are

667    often not perpendicular to the cortical lamina when taking a dorsal-ventral approach to IT face patches.

668    Laminar assignments of recordings were made by co-registering x-ray determined electrode

669    coordinates to MRI where the pial-to-gray matter border and the gray-to-white matter border were

670    defined. Based on our prior work estimating sources of error (e.g. error from electrode tip localization

671    and brain movement), registration of electrode tip locations to MRI brain volumes has a total of <400

672    micron error which is sufficient to distinguish deep from superficial layers (Issa et al., 2013). Multi-unit

673    activity (MUA) was systematically recorded at 300 micron intervals starting from penetration of the

674    superior temporal sulcus such that all sites at these regular intervals were tested with a screen set

675    containing both faces and non-face objects, and a subset of sites that were visually driven were further

676    tested with our main image set manipulating the position of face parts. Although we did not record

677    single-unit activity, our previous work showed similar responses between single-units and multi-units on

678    images of the type presented here (Issa & DiCarlo, 2012), and our results are consistent with

679    observations in previous single-unit work in IT (Freiwald et al., 2009). Recordings were made from PL,

680    ML, and AM in the left hemisphere of monkeys 1 and 2 and additionally from AL in monkey 2. AM and

681    AL are pooled together in our analyses forming the aIT sample while PL and ML correspond to the pIT

682    and cIT samples, respectively.

27

683

684    **Neural data analysis.** The face patches were physiologically defined in the same manner as in our

685    previous study (Issa & DiCarlo, 2012). Briefly, we fit a graded 3D sphere model (linear profile of

686    selectivity that rises from a baseline value toward the maximum at the center of the sphere) to the

687    spatial profile of face versus non-face object selectivity across our sites. We tested spherical regions

688    with radii from 1.5 to 10 mm and center positions within a 5 mm radius of the fMRI-based centers of the

689    face patches. The resulting physiologically defined regions were 1.5 to 3 mm in diameter. Sites which

690    passed a visual response screen (mean response in a 60-160 ms window >2*SEM above baseline for

691    at least one of the four categories in the screen set) were included in further analysis. All firing rates

692    were baseline subtracted using the activity in a 25-50 ms window following image onset averaged

693    across all repetitions of an image. Finally, given that the visual response latencies in monkey 2 were on

694    average 13 ms slower than those in monkey 1 for corresponding face-selective regions, we applied a

695    single latency correction (13 ms shift to align monkey 1 and monkey 2's data) prior to averaging across

696    monkeys. This was done so as not to wash out any fine timescale dynamics by averaging. Similar

697    results were obtained without using this latency correction as dynamics occurred at longer timescales

698    (~30 ms). This single absolute adjustment was more straightforward than the site-by-site adjustment

699    used in our previous work (Issa & DiCarlo, 2012) (though similar results were obtained using this

700    alternative latency correction); even when each monkey was analyzed separately, we still observed pIT

701    selectivity dynamics (**Figure 4A**). Images that produced an average population response $\geq$ 0.9 of the

702    initial response (60-100 ms) to a face image with all face parts arranged in their typical positions in a

703    frontal face were analyzed further (**Figures 2 and 3**). Stimulus selection was intended to limit

704    potentially confounding differences in visual drive between image classes. In a control test, we also

705    repeated our analysis by selecting images on a site-by-site basis where images with frontal face and

706    non-face arrangements of parts were chosen to be within 0.75x to 1.25x of the initial response to the

707    complete face image (minimum of five face and five non-face images in this response range for

28

708    inclusion of site in analysis). In follow-up analyses of population responses, we specifically limited

709    comparison to images with the same number of parts (**Figure 4B,C**). For example, for single part

710    images, we used the image with the eye in the upper, contralateral region of the outline as a reference

711    requiring a response $\geq$ 0.9 of the initial population response to this reference for inclusion of the images

712    in this analysis. We found that four other images of the 27 single-part images elicited a response at

713    least as large as 90% of the response to this standard image. For images containing all four face parts,

714    we used the complete, frontal face as the standard and found non-face arrangements of the four face

715    parts that drove at least 90% of the early response to the whole face (2 images out of 10 tested). To

716    compute individual site d' for each of these stimulus partitions (e.g. typical versus atypical

717    arrangements of 1 face part), we combined all presentations of images with frontal face arrangements

718    and compared these responses to responses from all presentations of images with non-face

719    arrangements using $d' = (u_1 - u_2)/((var_1 + var_2)/2)^{1/2}$ where variance was computed across all trials for that

720    image class (e.g. all presentations of all typical face images); this was identical to the d' measure used

721    in previous work for computing selectivity for faces versus non-face objects (Aparicio, Issa, & DiCarlo,

722    2016; Ohayon, Freiwald, & Tsao, 2012). For example, for the main image set (**Figure 2A**), we

723    compared all presentations of frontal face arrangements (8 images x 10 presentations/image = 80 total

724    presentations) to all presentations of non-face arrangements (13 images x 10 presentations/image =

725    130 total presentations) to compute the d' values for each site in two time windows (60-90 ms and 100-

726    130 ms) as shown in **Figure 3A**. A positive d' implies a stronger response to more naturally occurring

727    frontal face arrangements of face parts while a negative d' indicates a preference for unnatural non-

728    face arrangements of the face parts.

729

730    **Dynamical models**

731    *Modeling framework and equations.* To model the dynamics of neural response rates in a hierarchy, we

732    start with the simplest possible model that might capture those dynamics: a model architecture

733 consisting of a hidden stage of processing containing two units that linearly converge onto a single

734 output unit. We use this two-stage cascade for illustration of the basic concepts which can be easily

735 extended to longer cascades with additional stages, and we ultimately used a three-stage version of the

736 model to fit our neural data collected from three cortical stages (**Figures 5B & 6**).

737       An external input is applied separately to each hidden stage unit, which can be viewed as

738 representing different features for downstream integration. We vary the connections between the two

739 hidden units within the hidden processing stage (lateral connections) or between hidden and output

740 stage units (feedforward and feedback connections) to instantiate different model families. The details

741 of the different architectures specified by each model class can be visualized by their equivalent neural

742 network diagrams (**Figure 5**). Here, we provide a basic description for each model tested. All two-stage

743 models utilized a 2x2 feedforward identity matrix $A$ that simply transfers inputs $\boldsymbol{u}$ (2x1) to hidden layer

744 units $\boldsymbol{x}$ (2x1) and a 1x2 feedforward vector $B$ that integrates hidden layer activations $\boldsymbol{x}$ into a single

745 output unit $y$.

746

747     $A = aI,\ B = b[1,1]$                      (1)

748

749 By simply substituting in the appropriate unit vector and weight matrix transforming inputs from one

750 layer to the next for the desired network architecture, this simple two-stage architecture can be

751 extended to larger networks (e.g. see three-stage network diagrams in **Figure 5B**). To generate

752 dynamics in the simple networks below, we assumed that neurons act as leaky integrators of their total

753 synaptic input, a standard rate-based model of a neuron used in previous work (Seung, 1997)'(Rao &

754 Ballard, 1999).

755

756 *Pure feedforward.* In the purely feedforward family, connections are exclusively from hidden to output

757    stages through feedforward matrices *A* and *B*.

758

759    $$\dot{\mathbf{x}} = A\mathbf{u} - \mathbf{x}/\tau \qquad \dot{y} = B\mathbf{x} - y/\tau \qquad (2)$$

760

761    where $\tau$ is the time constant of the leak current which can be seen as reflecting the biophysical

762    limitations of neurons (a perfect integrator with large $\tau$ would have almost no leak and hence infinite

763    memory).

764

765    *Lateral inhibition.* Lateral connections (matrix with off-diagonal terms) are included and are inhibitory.

766    The scalar $k_l$ sets the relative strength of lateral inhibition versus bottom-up input.

767

768    $$\dot{\mathbf{x}} = A\mathbf{u} - \begin{bmatrix} 0 & k_l \\ k_l & 0 \end{bmatrix}\mathbf{x} - \mathbf{x}/\tau \qquad\qquad \dot{y} = B\mathbf{x} - y/\tau \qquad (3)$$

769

770    *Normalization.* An inhibitory term that scales with the summed activity of units within a stage is

771    included. The scalar $k_s$ sets the relative strength of normalization versus bottom-up input.

772

773    $$\dot{\mathbf{x}} = A\mathbf{u} - k_s \sum x \cdot \mathbf{x} - \mathbf{x}/\tau \qquad\qquad \dot{y} = B\mathbf{x} - k_s y \cdot y - y/\tau \qquad (4)$$

774

775    *Normalization (nonlinear)* (Carandini et al., 1997). The summed activity of units within a stage is used to

776    nonlinearly scale shunting inhibition.

777

778    $$\dot{\mathbf{x}} = A\mathbf{u} - \frac{\mathbf{x}}{\tau\sqrt{1 - k_s \sum x}} \qquad\qquad \dot{y} = B\mathbf{x} - \frac{y}{\tau\sqrt{1 - k_s y}} \qquad (5)$$

779

780     Note that this is technically a nonlinear dynamical system, and since the normalization term in equation

781     (5) is not continuously differentiable, we used the fourth-order Taylor approximation around zero in the

782     simulations of equation (5).

783

784     *Feedback (linear reconstruction).* The feedback-based model is derived using a normative framework

785     that performs optimal inference in the linear case (Seung, 1997) (unlike the networks in equations (2)-

786     (5) which are motivated from a mechanistic perspective but do not directly optimize a squared error

787     performance loss). The feedback network minimizes the cost *C* of reconstructing the inputs of each

788     stage (i.e. mean squared error of layer *n* predicting layer *n-1*).

789

790
$$C = \frac{1}{2}\left(\mathbf{u} - A^T\mathbf{x}\right)^2 + \frac{1}{2}\left(\mathbf{x} - B^T y\right)^2$$
(6)

791

792     Differentiating this coding cost with respect to the encoding variables in each layer *x*, *y* yields:

793

794
$$\frac{\partial C}{\partial \mathbf{x}} = -A(\mathbf{u} - A^T\mathbf{x}) + (\mathbf{x} - B^T\mathbf{y}) \qquad\qquad \frac{\partial C}{\partial \mathbf{y}} = -B(\mathbf{x} - B^T\mathbf{y})$$
(7)

795

796     The cost function *C* can be minimized by descending these gradients over time to optimize the values

797     of *x* and *y*:

798

799
$$\frac{\mathrm{d}\mathbf{x}}{\mathrm{d}t} = -\frac{\partial C}{\partial \mathbf{x}} = A(\mathbf{u} - A^T\mathbf{x}) - (\mathbf{x} - B^T\mathbf{y}) - \mathbf{x}/\tau$$

800
(8)

$$\frac{d\mathbf{y}}{dt} = -\frac{\partial C}{\partial \mathbf{y}} = B(\mathbf{x} - B^T\mathbf{y}) - \mathbf{y}/\tau$$

801

802

803    The above dynamical equations are equivalent to a linear network with a connection matrix containing

804    symmetric feedforward ($B$) and feedback ($B^T$) weights between stages $x$ and $y$ as well as within-stage

805    pooling followed by recurrent inhibition ($-AA^T x$ and $-BB^T y$) that resembles normalization. The property

806    that symmetric connections minimize the cost function $C$ generalizes to a feedforward network of any

807    size or number of hidden processing stages (i.e. holds for arbitrary lower triangular network connection

808    matrix). The final activation states ($x,y$) of the hierarchical generative network are optimal in the sense

809    that the bottom-up activations (implemented through feedforward connections) are balanced by the top-

810    down expectations (implemented by feedback connections) which is equivalent to a Bayesian network

811    combining bottom-up likelihoods with top-down priors to compute the maximum *a posteriori* (MAP)

812    estimate. Here, the priors are embedded in the weight structure of the network. In simulations, we

813    include an additional scalar $k_{td}$ that sets the relative weighting of bottom-up versus top-down signals.

814

815    $$\dot{\mathbf{x}} = A(\mathbf{u} - A^T\mathbf{x}) - k_{td}(\mathbf{x} - B^T y) - \mathbf{x}/\tau \qquad\qquad (9)$$

816

817    *Error signals computed in the feedback model.* In equation (9), inference can be thought of as

818    proceeding through integration of inputs on the dendrites of neuron population $x$. In this scenario, all

819    computations are implicit in dendritic integration. Alternatively, the computations in equation (9) can be

820    done in two steps where, in the first step, reconstruction errors are computed (i.e. $e_0 = u - A^T x$, $e_1 = x-$

821    $B^T y$) and explicitly represented in a separate error coding population. These error signals can then be

822    integrated by their downstream target population to generate the requisite update to the state signal of

823    neuron population $x$.

824

33

825 $$\dot{\mathbf{x}} = A\mathbf{e}_0 - k_{td}\mathbf{e}_1 - \mathbf{x}/\tau \qquad\qquad \dot{\mathbf{y}} = B\mathbf{e}_1 - y/\tau \qquad\qquad (10)$$

826

827 An advantage of this strategy is that the a state unit now directly receives errors as inputs, and those

828 inputs allow implementation of an efficient Hebbian rule for learning weight matrices (Rao & Ballard,

829 1999) -- the gradient rule for learning is simply a product of the state activation and the input error

830 activation (weight updates obtained by differentiating equation (6) with respect to weight matrices *A* and

831 *B*: $\Delta A = \mathbf{x} \cdot \mathbf{e_0}^T$, $\Delta A^T = \mathbf{e_0} \cdot \mathbf{x}^T$, $\Delta B = y \cdot \mathbf{e_1}^T$, and $\Delta B^T = \mathbf{e_1} \cdot y$). Thus, the reconstruction errors serve as

832 computational intermediates for both the gradients of online inference mediated by dendritic integration

833 (dynamics in state space, equation (10)) and gradients for offline learning mediated by Hebbian

834 plasticity (dynamics in weight space).

835

836 In order for the reconstruction errors at each layer to be scaled appropriately in the feedback

837 model, we invoke an additional downstream variable *z* to predict activity at the top stage such that,

838 instead of $\mathbf{e_2} = y$ which scales as a state variable, we have $\mathbf{e_2} = y - C^T z$ (**Figure 5A**). This overall model

839 reflects a state and error coding model as opposed to a state only model.

840

841 *Feedback (three-stage)*. For the simulations in **Figures 6 and 8**, three-stage versions of the above

842 equations were used. These deeper networks were also wider such that they began with four input

843 units (***u***) instead of only two inputs in the two-stage models. These inputs converged through

844 successive processing stages (***w,x,y***) to one unit at the top node (***z***) (**Figure 5B**).

845

846 *Feedback (nonlinear reconstruction)*. To test the generality of our findings beyond a linear

847 reconstruction cost, we simulated feedback-based models which optimized different candidate cost

848 functions proposed for the ventral stream (**Figure 7C**). In nonlinear hierarchical inference,

849   reconstruction is performed using a monotonic nonlinearity with a threshold (*th*) and bias (*bi*):

850

$$C = \frac{1}{2}\left(\mathbf{u} - f(A^T\mathbf{x})\right)^2 + \frac{1}{2}\left(\mathbf{x} - f(B^T\mathbf{y})\right)^2 , \quad where\ f(x) = tanh(x - th) + bi \quad (11)$$

852

$$\dot{\mathbf{x}} = A(\mathbf{u} - f(A^T\mathbf{x}))(1 - tanh(A^T\mathbf{x} - th)^2) - k_{td}(\mathbf{x} - f(B^T y)) - \mathbf{x}/\tau$$

854                                                                                                  (12)

$$\dot{\mathbf{y}} = B(\mathbf{x} - f(B^T\mathbf{y}))(1 - tanh(B^T\mathbf{y} - th)^2) - \mathbf{y}/\tau$$

856

857   *Feedback (linear construction).* Instead of a reconstruction cost where responses match the input (i.e.

858   generative model) as in unsupervised learning, we additionally simulated the states and errors in a

859   feedback network minimizing a linear construction cost where the network is producing responses to

860   match a given output (i.e. discriminative model) similar to supervised learning:

861

$$C = \frac{1}{2}\left(A\mathbf{u} - \mathbf{x}\right)^2 + \frac{1}{2}\left(B\mathbf{x} - \mathbf{y}\right)^2 \qquad (13)$$

863

$$\dot{\mathbf{x}} = (A\mathbf{u} - \mathbf{x}) - k_{td}B^T(B\mathbf{x} - \mathbf{y}) - \mathbf{x}/\tau \qquad \dot{\mathbf{y}} = (B\mathbf{x} - \mathbf{y}) - \mathbf{y}/\tau \qquad (14)$$

865

866   **Model simulation.** To simulate the dynamical systems in equations (2)-(14), a step input *u* was

867   applied. This input was smoothed using a Gaussian kernel to approximate the lowpass nature of signal

868   propagation in the series of processing stages from the retina to pIT:

869

$$\mathbf{u}(t) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(t-t_0)^2}{2\sigma^2}} * \mathbf{h}(t) \Rightarrow \dot{\mathbf{u}} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(t-t_0)^2}{2\sigma^2}} \cdot \mathbf{h}$$

870                                                                                                  (15)

871

872    where the elements of **h** are scaled Heaviside step functions. The input is thus a sigmoidal ramp whose

873    latency to half height is set by $t_0$ and rise time is set by $\sigma$. For simulation of two-stage models, there

874    were ten basic parameters: latency of the input $t_0$, standard deviation of the Gaussian ramp $\sigma$, system

875    time constant $\tau$, input connection strength $a$, feedforward connection strength $b$, the four input values

876    across two stimulus conditions (i.e. $h_{11}$, $h_{12}$, $h_{21}$, $h_{22}$), and a factor $sc$ for scaling the final output to the

877    neural activity. In the deeper three-stage network, there were a total of fifteen parameters which

878    included an additional feedforward connection strength $c$ and additional input values since the three-

879    stage model had four inputs instead of two. The lateral inhibition model class required one additional

880    parameter $k_l$ as did the normalization model family $k_s$, and for feedback model simulations, there was

881    an additional feedback weight $k_{td}$ to scale the relative contribution of the top-down errors in driving

882    online inference. For the error coding variants of the feedback model, gain parameters $c$ (two-stage)

883    and $d$ (three-stage) were included to scale the overall magnitude of the top level reconstruction error

884    (also see **Figure 5** for locations of parameters in network diagrams).

885

886    **Model parameter fits to neural data**. In fitting the models to the observed neural dynamics, we

887    mapped the summed activity in the hidden stage (**w**) to population averaged activity in pIT, and we

888    mapped the summed activity in the output stage (**y**) to population averaged signals measured in aIT. To

889    simulate error coding, we mapped the reconstruction errors $e_1 = w\text{-}B^T x$ and $e_3 = y\text{-}C^T z$ to activity in pIT

890    and aIT, respectively. We applied a squaring nonlinearity to the model outputs as an approximation to

891    rectification since recorded extracellular firing rates are non-negative (and linear rectification is not

892    continuously differentiable). Analytically solving this system of dynamical equations (2)-(14) for a step

893    input is precluded because of the higher order interaction terms (the roots of the determinant and hence

894    the eigenvalues/eigenvectors of a 3x3 or larger matrix are not analytically determined, except for the

895    purely feedforward model which only has first-order interactions), and in the case of the normalization

896    models, there is an additional nonlinear dependence on the shunt term. Thus, we relied on

897    computational methods (constrained nonlinear optimization) to fit the parameters of the dynamical

898    systems to the neural data with a quadratic (sum of squares) loss function.

899

900        Parameter values were fit in a two-step procedure. In the first step, we fit only the difference in

901    response between image classes (differential mode which is the selectivity profile over time, see **Figure**

902    **6A**, right data panel), and in the second step, we refined fits to capture an equally weighted average of

903    the differential mode and the common mode (the common mode is the average across images of the

904    response time course of visual drive). This two-step procedure was used to ensure that each model

905    had the best chance of fitting the dynamics of selectivity (differential mode) as these selectivity profiles

906    were the main phenomena of interest but were smaller in size (20% of response) compared to overall

907    visual drive. In each step, fits were done using a large-scale algorithm (interior-point) to optimize

908    coarsely, and the resulting solution was used as the initial condition for a medium-scale algorithm

909    (sequential quadratic programming) for additional refinement. The lower and upper parameter bounds

910    tested were: $t_0$=[50 70], $\sigma$=[0.5 25], $\tau$ =[0.5 1000], $k_l, k_s, k_{td}$=[0 1], $a,b,c,d$=[0 2], $h$=[0 20], $sc$=[0 100], $th$=[-

911    20 20], and $bi$=[-1 1] which proved to be adequately liberal as parameter values converged to values

912    that did not generally approach these boundaries. To avoid local minima, the algorithm was initialized to

913    a number of randomly selected points (n = 50), and after fitting the differential mode, we took the top

914    fits (n = 25) for each model class and used these as initializations in subsequent steps. The single best

915    fitting instance of each model class is shown in the main figures.

916

917    **Model predictions.** For the predictions in **Figure 8**, all architectural parameters obtained by the fitting

918    procedure above were held fixed; only the pattern of inputs to the network was varied. For **Figure 8A**,

919    to test the input integration properties of a model, we used the top-performing model and compared the

920    response to all inputs presented simultaneously with the sum of the responses to each input alone.

921

922        For **Figure 8B**, we approximated novel versus familiar images by parametrically varying the

923    degree to which input patterns were random (inputs drawn from i.i.d. uniform distributions on the

924    interval [0,1]) versus structured in a way that matched the weight pattern of the network. Here, we used

925    a version of the model with two independent outputs reflecting detectors for two familiarized input

926    patterns (output 1 tuned to pattern *A*: $u_1, u_2, u_3, u_4$ active and output 2 tuned to pattern 2: $u_5, u_6, u_7, u_8$

927    active) (**Figure 8B**). Alternating between these two input patterns simulates alternation of two

928    familiarized (learned) images as compared to purely random patterns ($u_{1-8}$ independent and identically

929    distributed). The first-to-second layer weights were [1,1,1,1,0,0,0,0] for pattern A and [0,0,0,0,1,1,1,1]

930    for pattern B, so to parametrically vary the degree of correlation of inputs to this weight pattern from

931    random (correlation = 0) to deterministic (correlation = 1), we drew input values from a joint distribution

932    $P(u_1,u_2,u_3,u_4,u_5,u_6,u_7,u_8)$ where $u_{1-4}$ were drawn from a high-valued uniform distribution on the interval

933    $[1-\varepsilon,1]$ and $u_{5-8}$ were drawn from a low-valued uniform distribution $[0, \varepsilon]$ for stimulus pattern *A* and the

934    opposite for pattern *B* ($u_{5-8}$ high-valued and $u_{1-4}$ low-valued). The parameter $\varepsilon$ determines the range of

935    values that could be drawn from purely deterministic (0 or 1, $\varepsilon = 0$) to randomly uniformly distributed

936    (from 0 to 1, $\varepsilon = 1$). Thus, the correlation of the inputs correspondingly varies according to $\rho(u_i, u_j)$

937    $= \rho(u_k, u_l) = (\varepsilon-1)^2/((\varepsilon-1)^2+\varepsilon^2/3)$ where $1 \leq i,j \leq 4$, $i \neq j$ and $5 \leq k,l \leq 8$, $k \neq l$ approaching correlation equal to 0 for a

938    purely, random pattern ($\varepsilon = 1$) that had a low probability of matching the learned patterns *A* and *B*.

939

940    **Code availability.** All data analysis and computational modeling were done using custom scripts

941    written in Matlab. All code is available upon request.

942

943    **Statistics.** Error bars represent standard errors of the mean obtained by bootstrap resampling (n =

944    1000). All statistical comparisons including those of means or correlation values were obtained by

945    bootstrap resampling (n = 1000) producing p-values at a resolution of 0.001 so that the lowest p-value

946    that can be reported is p = 0.000 given the resolution of this statistical analysis. All statistical tests were

947    two-sided unless otherwise specified. Spearman's rank correlation coefficient was used.

948

953    **REFERENCES**
954
955    Ackley, D. H., Hinton, G. E., & Sejnowski, T. J. (1985). A learning algorithm for boltzmann machines.
956         *Cognitive Science*, *9*(1), 147–169. https://doi.org/10.1016/S0364-0213(85)80012-4
957
958    Afraz, S.-R., Kiani, R., & Esteky, H. (2006). Microstimulation of inferotemporal cortex influences face
959         categorization. *Nature*, *442*(7103), 692–695. https://doi.org/10.1038/nature04982
960
961    Aparicio, P. L., Issa, E. B., & DiCarlo, J. J. (2016). Neurophysiological Organization of the Middle Face
962         Patch in Macaque Inferior Temporal Cortex. *Journal of Neuroscience*, *36*(50), 12729–12745.
963         https://doi.org/10.1523/JNEUROSCI.0237-16.2016
964
965    Brincat, S. L., & Connor, C. E. (2006). Dynamic Shape Synthesis in Posterior Inferotemporal Cortex.
966         *Neuron*, *49*(1), 17–24. https://doi.org/10.1016/j.neuron.2005.11.026
967
968    Cadieu, C. F., Hong, H., Yamins, D. L. K., Pinto, N., Ardila, D., Solomon, E. A., … DiCarlo, J. J. (2014). Deep
969         Neural Networks Rival the Representation of Primate IT Cortex for Core Visual Object
970         Recognition. *PLoS Comput Biol*, *10*(12), e1003963.
971         https://doi.org/10.1371/journal.pcbi.1003963
972
973    Carandini, M., & Heeger, D. J. (2012). Normalization as a canonical neural computation. *Nature Reviews
974         Neuroscience*, *13*(1), 51–62. https://doi.org/10.1038/nrn3136
975
976    Carandini, M., Heeger, D. J., & Movshon, J. A. (1997). Linearity and Normalization in Simple Cells of the
977         Macaque Primary Visual Cortex. *The Journal of Neuroscience*, *17*(21), 8621–8644.
978
979    Chang, L., & Tsao, D. Y. (2017). The Code for Facial Identity in the Primate Brain. *Cell*, *169*(6), 1013-
980         1028.e14. https://doi.org/10.1016/j.cell.2017.05.011
981
982    Chen, M., Yan, Y., Gong, X., Gilbert, C. D., Liang, H., & Li, W. (2014). Incremental Integration of Global
983         Contours through Interplay between Visual Cortical Areas. *Neuron*, *82*(3), 682–694.
984         https://doi.org/10.1016/j.neuron.2014.03.023
985
986    Coen-Cagli, R., Kohn, A., & Schwartz, O. (2015). Flexible gating of contextual influences in natural
987         vision. *Nature Neuroscience*, *18*(11), 1648–1655. https://doi.org/10.1038/nn.4128
988
989    Cox, D. D., Papanastassiou, A. M., Oreper, D., Andken, B. B., & DiCarlo, J. J. (2008). High-Resolution
990         Three-Dimensional Microelectrode Brain Mapping Using Stereo Microfocal X-ray Imaging.
991         *Journal of Neurophysiology*, *100*(5), 2966–2976. https://doi.org/10.1152/jn.90672.2008
992
993    Egeth, H. E., & Yantis,  and S. (1997). VISUAL ATTENTION: Control, Representation, and Time Course.
994         *Annual Review of Psychology*, *48*(1), 269–297. https://doi.org/10.1146/annurev.psych.48.1.269
995

996   Epshtein, B., Lifshitz, I., & Ullman, S. (2008). Image interpretation by a single bottom-up top-down
997         cycle. *Proceedings of the National Academy of Sciences*, *105*(38), 14298–14303.
998         https://doi.org/10.1073/pnas.0800968105
999

1000  Freedman, D. J., Riesenhuber, M., Poggio, T., & Miller, E. K. (2006). Experience-Dependent Sharpening
1001        of Visual Shape Selectivity in Inferior Temporal Cortex. *Cerebral Cortex*, *16*(11), 1631–1644.
1002        https://doi.org/10.1093/cercor/bhj100
1003

1004  Freiwald, W. A., & Tsao, D. Y. (2010). Functional Compartmentalization and Viewpoint Generalization
1005        Within the Macaque Face-Processing System. *Science*, *330*(6005), 845–851.
1006        https://doi.org/10.1126/science.1194908
1007

1008  Freiwald, W. A., Tsao, D. Y., & Livingstone, M. S. (2009). A face feature space in the macaque temporal
1009        lobe. *Nature Neuroscience*, *12*(9), 1187–1196. https://doi.org/10.1038/nn.2363
1010

1011  Friston, K., & Kiebel, S. (2009). Cortical circuits for perceptual inference. *Neural Networks*, *22*(8), 1093–
1012        1104. https://doi.org/10.1016/j.neunet.2009.07.023
1013

1014  Grimaldi, P., Saleem, K. S., & Tsao, D. (2016). Anatomical Connections of the Functionally Defined "Face
1015        Patches" in the Macaque Monkey. *Neuron*, *90*(6), 1325–1342.
1016        https://doi.org/10.1016/j.neuron.2016.05.009
1017

1018  Hung, C. P., Kreiman, G., Poggio, T., & DiCarlo, J. J. (2005). Fast Readout of Object Identity from
1019        Macaque Inferior Temporal Cortex. *Science*, *310*(5749), 863–866.
1020        https://doi.org/10.1126/science.1117593
1021

1022  Hyvärinen, A., Hurri, J., & Hoyer, P. O. (2009). Lateral interactions and feedback. In *Natural Image
1023        Statistics: A Probabilistic Approach to Early Computational Vision.* (Vol. 39, pp. 307–318).
1024        Springer Science & Business Media.
1025

1026  Issa, E. B., & DiCarlo, J. J. (2012). Precedence of the Eye Region in Neural Processing of Faces. *The
1027        Journal of Neuroscience*, *32*(47), 16666–16682. https://doi.org/10.1523/JNEUROSCI.2391-
1028        12.2012
1029  Issa, E. B., Papanastassiou, A. M., Andken, B. B., & DiCarlo, J. J. (2010). Towards large-scale, high
1030        resolution maps of object selectivity in inferior temporal cortex. Presented at the
1031        Computational and Systems Neuroscience, Salt Lake City, UT: Front in Neuroscienc Abstract.
1032        https://doi.org/10.3389/conf.fnins.2010.03.00154
1033

1034  Issa, E. B., Papanastassiou, A. M., & DiCarlo, J. J. (2013). Large-Scale, High-Resolution
1035        Neurophysiological Maps Underlying fMRI of Macaque Temporal Lobe. *The Journal of
1036        Neuroscience*, *33*(38), 15207–15219. https://doi.org/10.1523/JNEUROSCI.1248-13.2013
1037

1038  Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional
1039        Neural Networks. In *Advances in Neural Information Processing Systems* (pp. 1106–1114).

1040

1041  Lee, T. S., & Mumford, D. (2003). Hierarchical Bayesian inference in the visual cortex. *Journal of the*
1042        *Optical Society of America A*, *20*(7), 1434. https://doi.org/10.1364/JOSAA.20.001434

1043

1044  Lee, T. S., Yang, C. F., Romero, R. D., & Mumford, D. (2002). Neural activity in early visual cortex reflects
1045        behavioral experience and higher-order perceptual saliency. *Nature Neuroscience*, *5*(6), 589–
1046        597. https://doi.org/10.1038/nn0602-860

1047

1048  Leopold, D. A., Bondar, I. V., & Giese, M. A. (2006). Norm-based face encoding by single neurons in the
1049        monkey inferotemporal cortex. *Nature*, *442*(7102), 572–575.
1050        https://doi.org/10.1038/nature04951

1051

1052  Majaj, N. J., Hong, H., Solomon, E. A., & DiCarlo, J. J. (2015). Simple Learned Weighted Sums of Inferior
1053        Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition
1054        Performance. *The Journal of Neuroscience*, *35*(39), 13402–13418.
1055        https://doi.org/10.1523/JNEUROSCI.5181-14.2015

1056

1057  Marblestone, A. H., Wayne, G., & Kording, K. P. (2016). Toward an Integration of Deep Learning and
1058        Neuroscience. *Frontiers in Computational Neuroscience*, *10*.
1059        https://doi.org/10.3389/fncom.2016.00094

1060

1061  Meyer, T., & Olson, C. R. (2011). Statistical learning of visual transitions in monkey inferotemporal
1062        cortex. *Proceedings of the National Academy of Sciences*, *108*(48), 19401–19406.
1063        https://doi.org/10.1073/pnas.1112895108

1064

1065  Meyer, T., Walker, C., Cho, R. Y., & Olson, C. R. (2014). Image familiarization sharpens response
1066        dynamics of neurons in inferotemporal cortex. *Nature Neuroscience*, *17*(10), 1388–1394.
1067        https://doi.org/10.1038/nn.3794

1068

1069  Meyers, E. M., Borzello, M., Freiwald, W. A., & Tsao, D. (2015). Intelligent Information Loss: The Coding
1070        of Facial Identity, Head Pose, and Non-Face Information in the Macaque Face Patch System. *The*
1071        *Journal of Neuroscience*, *35*(18), 7069–7081. https://doi.org/10.1523/JNEUROSCI.3086-14.2015

1072

1073  Moeller, S., Crapse, T., Chang, L., & Tsao, D. Y. (2017). The effect of face patch microstimulation on
1074        perception of faces and objects. *Nature Neuroscience*, *advance online publication*.
1075        https://doi.org/10.1038/nn.4527

1076

1077  Moeller, S., Freiwald, W. A., & Tsao, D. Y. (2008). Patches with Links: A Unified System for Processing
1078        Faces in the Macaque Temporal Lobe. *Science*, *320*(5881), 1355–1359.
1079        https://doi.org/10.1126/science.1157436

1080

1081  Murray, E. A., Bussey, T. J., & Saksida, L. M. (2007). Visual Perception and Memory: A New View of
1082        Medial Temporal Lobe Function in Primates and Rodents. *Annual Review of Neuroscience*,
1083        *30*(1), 99–122. https://doi.org/10.1146/annurev.neuro.29.051605.113046

1084
1085   Nassi, J. J., Gómez-Laberge, C., Kreiman, G., & Born, R. T. (2014). Corticocortical feedback increases the
1086           spatial extent of normalization. *Frontiers in Systems Neuroscience*, *8*.
1087           https://doi.org/10.3389/fnsys.2014.00105
1088
1089   Ohayon, S., Freiwald, W. A., & Tsao, D. Y. (2012). What Makes a Cell Face Selective? The Importance of
1090           Contrast. *Neuron*, *74*(3), 567–581. https://doi.org/10.1016/j.neuron.2012.03.024
1091
1092   Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive field properties by learning a
1093           sparse code for natural images. *Nature*, *381*(6583), 607–609.
1094           https://doi.org/10.1038/381607a0
1095
1096   Patel, A. B., Nguyen, T., & Baraniuk, R. G. (2015). A Probabilistic Theory of Deep Learning.
1097           *ArXiv:1504.00641 [Cs, Stat]*. Retrieved from http://arxiv.org/abs/1504.00641
1098
1099   Rao, R. P. N., & Ballard, D. H. (1999). Predictive coding in the visual cortex: a functional interpretation
1100           of some extra-classical receptive-field effects. *Nature Neuroscience*, *2*(1), 79–87.
1101           https://doi.org/10.1038/4580
1102
1103   Rifai, S., Vincent, P., Muller, X., Glorot, X., & Bengio, Y. (2011). Contractive auto-encoders: Explicit
1104           invariance during feature extraction. In *Proceedings of the 28th international conference on
1105           machine learning (ICML-11)* (pp. 833–840).
1106
1107   Sadagopan, S., Zarco, W., & Freiwald, W. A. (2017). A causal relationship between face-patch activity
1108           and face-detection behavior. *ELife*, *6*, e18558. https://doi.org/10.7554/eLife.18558
1109
1110   Salakhutdinov, R., & Hinton, G. (2012). An Efficient Learning Procedure for Deep Boltzmann Machines.
1111           *Neural Computation*, *24*(8), 1967–2006. https://doi.org/10.1162/NECO_a_00311
1112
1113   Schwartz, O., & Simoncelli, E. P. (2001). Natural signal statistics and sensory gain control. *Nature
1114           Neuroscience*, *4*(8), 819–825. https://doi.org/10.1038/90526
1115
1116   Schwiedrzik, C. M., & Freiwald, W. A. (2017). High-Level Prediction Signals in a Low-Level Area of the
1117           Macaque Face-Processing Hierarchy. *Neuron*, *96*(1), 89-97.e4.
1118           https://doi.org/10.1016/j.neuron.2017.09.007
1119
1120   Seung, H. S. (1997). Pattern analysis and synthesis in attractor neural networks. *Theoretical Aspects of
1121           Neural Computation: A Multidisciplinary Perspective, Singapore*.
1122
1123   Sugase, Y., Yamane, S., Ueno, S., & Kawano, K. (1999). Global and fine information coded by single
1124           neurons in the temporal visual cortex. *Nature*, *400*(6747), 869–873.
1125           https://doi.org/10.1038/23703
1126

1127  Tsao, D. Y., Freiwald, W. A., Tootell, R. B. H., & Livingstone, M. S. (2006). A Cortical Region Consisting
1128        Entirely of Face-Selective Cells. *Science*, *311*(5761), 670–674.
1129        https://doi.org/10.1126/science.1119983
1130
1131  Tsao, D. Y., Moeller, S., & Freiwald, W. A. (2008). Comparing face patch systems in macaques and
1132        humans. *Proceedings of the National Academy of Sciences*, *105*(49), 19514–19519.
1133        https://doi.org/10.1073/pnas.0809662105
1134
1135  Williams, D. R. G. H. R., & Hinton, G. E. (1986). Learning representations by back-propagating errors.
1136        *Nature*, *323*, 533–536.
1137
1138  Woloszyn, L., & Sheinberg, D. L. (2012). Effects of Long-Term Visual Experience on Responses of
1139        Distinct Classes of Single Units in Inferior Temporal Cortex. *Neuron*, *74*(1), 193–205.
1140        https://doi.org/10.1016/j.neuron.2012.01.032
1141
1142  Yamins, D. L. K., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., & DiCarlo, J. J. (2014). Performance-
1143        optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of
1144        the National Academy of Sciences*, *111*(23), 8619–8624.
1145        https://doi.org/10.1073/pnas.1403112111
1146
1147  Zeiler, M. D., & Fergus, R. (2013). Stochastic Pooling for Regularization of Deep Convolutional Neural
1148        Networks. *ArXiv:1301.3557 [Cs, Stat]*. Retrieved from http://arxiv.org/abs/1301.3557
1149
1150  Zhang, N. R., & Heydt, R. von der. (2010). Analysis of the Context Integration Mechanisms Underlying
1151        Figure–Ground Organization in the Visual Cortex. *The Journal of Neuroscience*, *30*(19), 6482–
1152        6496. https://doi.org/10.1523/JNEUROSCI.5168-09.2010
1153

**Figure 1**. **Neural recordings and experimental design in face-selective subregions of the ventral visual stream. (A)** Neurons were recorded along the lateral convexity of the inferior temporal lobe spanning the posterior to anterior extent of IT (+0 to +20 mm AP, Horsely-Clarke coordinates) in two monkeys (data from monkey 1 are shown). Based on prior work, face-selective sites (red) were operationally defined as those with a response preference for images of frontal faces versus images of non-face objects (see **Methods**). While these neurons were found throughout IT, they tended to be found in clusters that mapped to previously identified subdivisions of IT (posterior, central, and anterior IT) and corresponded to face-selective areas identified under fMRI in the same subjects (Issa & DiCarlo, 2012)(Issa et al., 2013) (STS = superior temporal sulcus, IOS = inferior occipital sulcus, OTS = occipitotemporal sulcus). **(B)** (top diagram) The three visual processing stages in IT lie downstream of early visual areas V1, V2, and V4 in the ventral visual stream. (left) We designed our stimuli to focus on the intermediate stage pIT by seeking images of faces and images of non-faces that would, on average, drive equally strong initial responses in pIT. Novel images were generated from an exemplar monkey face by positioning the face parts in different positions within the face outline. This procedure generated both frontal face and non-face arrangements of the face parts, and we identified 21 images (red and black boxes) that drove the mean, early (60-100 ms) pIT population response to ≥90% of its response to the intact face (first image in red box is synthesized whole face; compare to the second image which is the original whole face), and of these 21 images, 13 images contained non-face arrangements of the face parts. For example, images with an eye centered in the outline (black box, 3[rd] and 4[th] rows) as opposed to the lateralized position of the eye in a frontal face (red box) have a global interpretation ("cyclops") that is not consistent with a frontal face but still evoked strong pIT responses. Selectivity of neural sites (see **Figs. 3 & 4**) for face versus non-face images was quantified using a d' measure. (middle) Computational hypotheses of cortical dynamics make differing predictions about how neural selectivity in pIT may evolve following an initial, weak preference signal for images of frontal faces. (right) Predictions of how aIT would behave as an output stage building selectivity for images of frontal faces through multiple stages of processing. **(C)** A population decoder, trained on average firing rates (60-200 ms post image onset, linear SVM classifier) for frontal face versus non-face arrangements of the face parts in this image subset, performed poorly in pIT on held-out trials of the same images (trial splits used so that the same images were shown in classifier training and testing). However, the particular class (face vs non-face) could be determined at above chance levels when reading the cIT and aIT population responses.
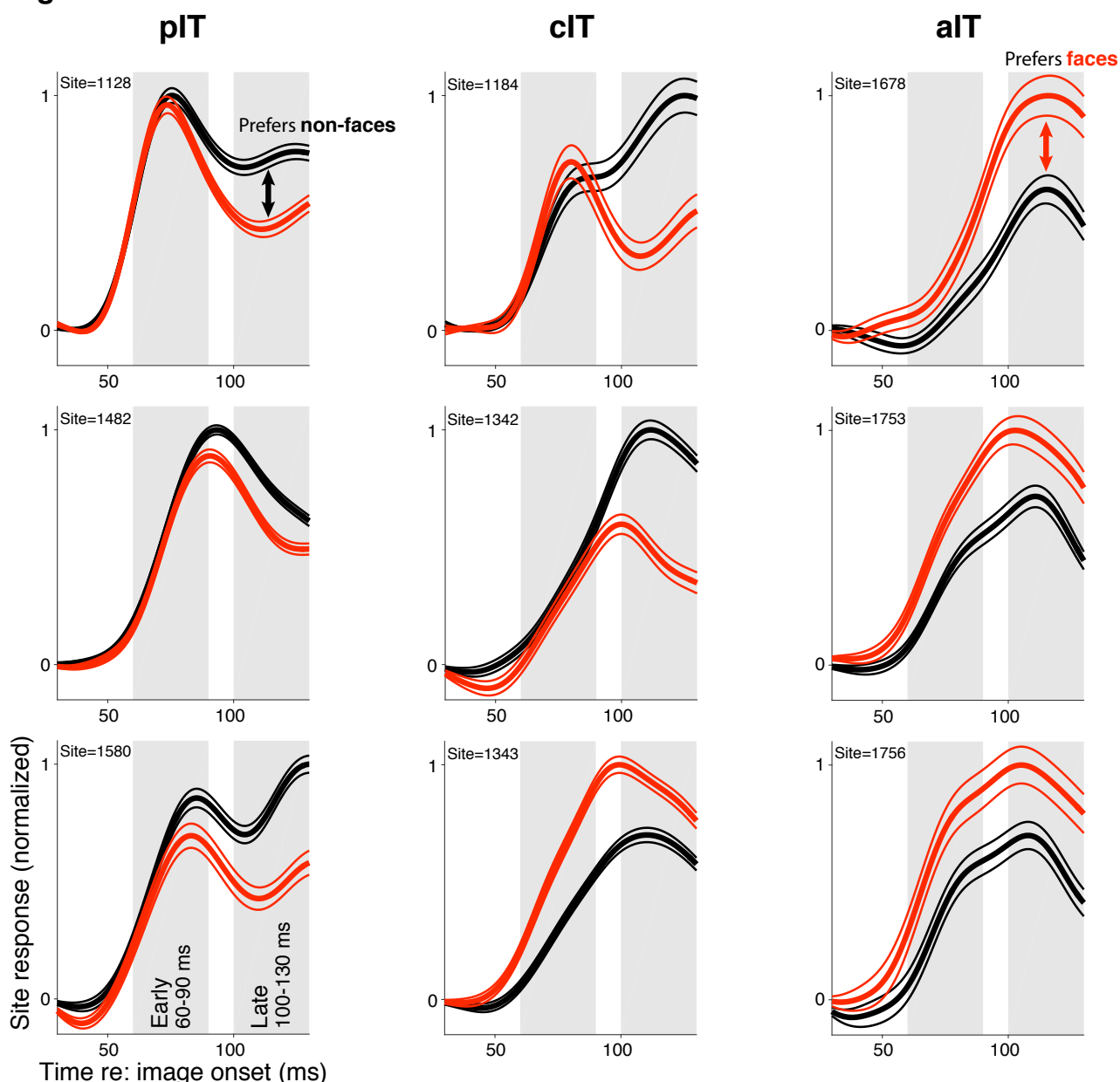
45

## Figure 2



1188
1189
1190  **Figure 2. Responses in example sites to face-like images with typical and atypical face part**
1191  **arrangements.** The three sites with the highest selectivity in the late response phase in each region
1192  are shown (pIT, cIT, and aIT; left, middle, and right columns, respectively) (d' selectivity measured in a
1193  100-130 ms window, gray shaded region shown in bottom, left panel). While the three aIT sites (right
1194  column) demonstrated late phase selectivity for face images, the three pIT sites evolved the opposite
1195  preference in their late phase (100-130 ms) responses (red line = mean response of 8 images shown in
1196  **Figure 1B** red box, and black line = mean response of 13 images shown in **Figure 1B** black box).
1197

**Figure 3**



1198

**Figure 3. Time course of neural response preferences in pIT, cIT, and aIT for images with face versus non-face arrangements of the parts. (A)** Preferences for frontal face versus non-face part arrangements for each site are plotted in both early (60-90 ms post image onset) and late (100-130 ms) time windows. Sites are grouped based on region (pIT, cIT, aIT) and whether they showed a significant change in selectivity from early to late time windows (light gray = increased preference, black = decreased preference, and dark gray = no change in preference for face versus non-face images, significance tested at $p < 0.01$ level; example sites from **Figure 2** are plotted using thicker, darker lines). Many sites in pIT and cIT showed a decreasing preference for frontal face versus non-face images over time (black lines, middle row, left and center panels). In contrast, no sites in aIT had this dynamic (middle row, right panel). **(B)** The fraction of sites whose responses showed a preference for images of typical, face-like arrangements of the face parts in pIT (blue), cIT (green) and aIT (red) in the early (60-90 ms) and late (100-130 ms) phase of the response. Note that, in the late phase of the response, most pIT neurons paradoxically showed a preference for non-face arrangements of face parts. **(C)** Selectivity measured for images driving similar responses within a site. This procedure ensured matched initial responses on a site-by-site basis rather than using a fixed set of images based on the overall population response (i.e. the fixed image set of **Figure 1B**; here, the initial d' for 60-90 ms is close to zero when images are selected site by site). Although initial response differences were near zero when using site based image selection, a late phase preference for non-face images still emerged in pIT and cIT but not in aIT similar to the decreasing selectivity profile observed when using a fixed image set for all sites.
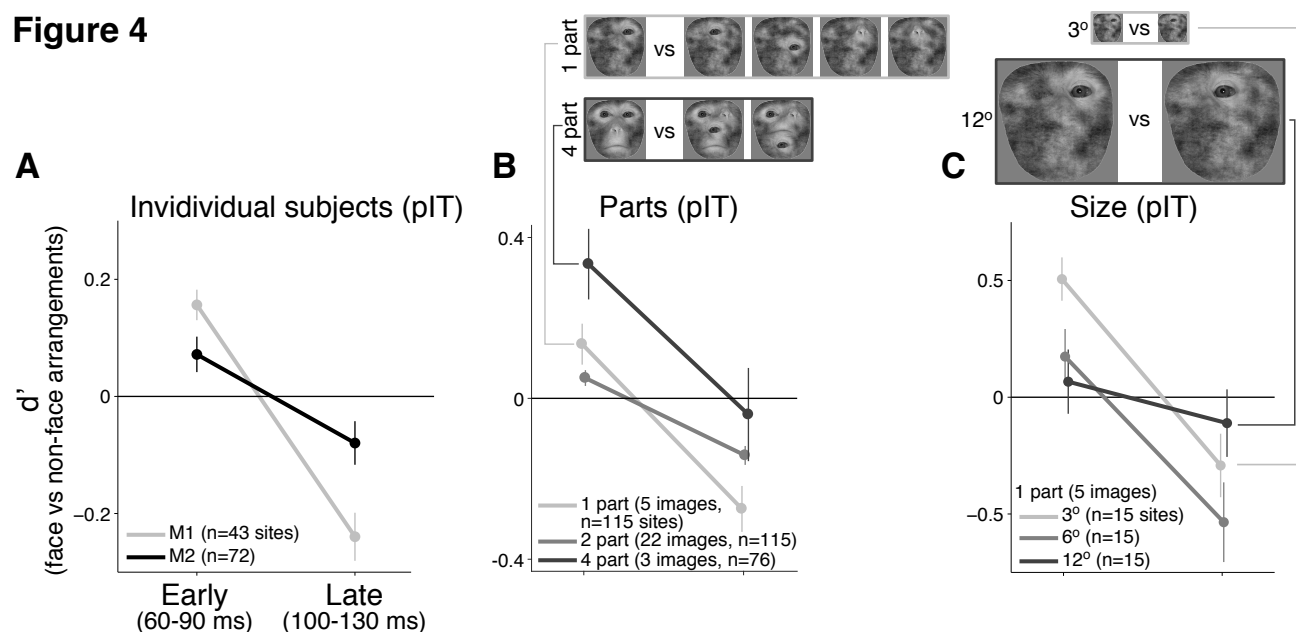
**Figure 4**



**Figure 4. Individual monkey comparison and image controls for the decreasing selectivity profile in pIT. (A)** Preference for images with frontal face part arrangements analyzed separately for each monkey. Median d' of pIT sites in both early and late time windows is shown. **(B)** Preference for images with face versus non-face arrangements of the parts was re-computed using image subsets containing the same number of parts in the outline (the five 1-part and the three 4-part image subsets shown at top; the larger 2-part subset contained 30 images and is not shown). **(C)** The 1-part image subset was further tested at three different sizes ($3^o$, $6^o$, and $12^o$). In all cases, pIT responses showed a decreasing preference over time for typically-arranged face parts leading to a preference for atypically arranged face parts in the later time window (100-130 ms).
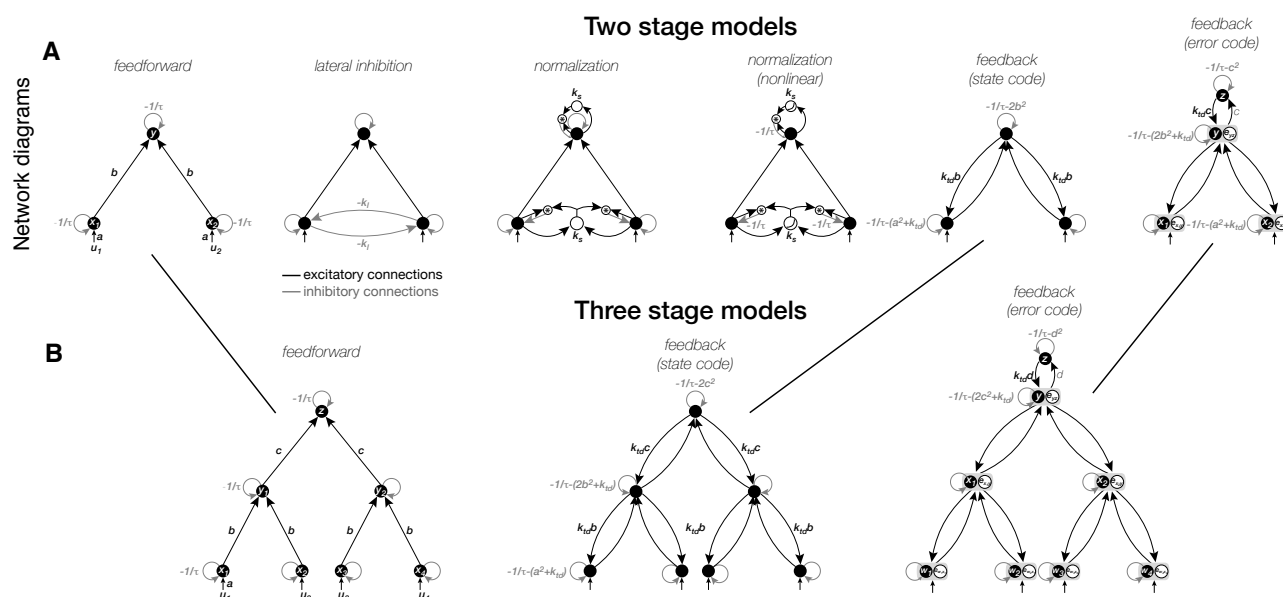
**Figure 5**



1232
1233
**Figure 5. Model diagrams. (A)** Network diagrams of two-stage models used in **Figure 7A,B**. The base feedforward architecture and parameters are shown in the first column while recurrent models with added connections and parameters are shown in the remaining columns. All models have two inputs ($u_1$, $u_2$), at least two hidden units ($x_1$, $x_2$) (a.k.a. layer 1 of the network), and at least one output unit ($y$). For simulations, the inputs $u_1$, $u_2$ are set independently to simulate each hidden node receiving different amounts of external drive, depending on the choice of the applied image relative to the unit's preferred image. The connection weights B = [$b_1$, $b_2$] transforming the hidden stage activations to the output unit are modeled as the same ($b_1 = b_2$). All units have self-connections that determine the degree of leak current set by the time constant $\tau$. In the normalization models, the leak term is additionally controlled (linearly or nonlinearly) by the total activity in each stage (third and fourth columns). In the feedback-based model, the feedback connections are symmetric to the feedforward connections with weights $B^T$ = [$b_1$, $b_2$]$^T$, a column vector (fifth and sixth columns). The error coding feedback model (sixth column) has an additional stage that contributes to computation of error in the second stage (see **Methods** for details). **(B)** Extensions of the two-stage model architectures to three stages are shown only for the feedforward and feedback models (compare to two-stage diagrams in **(A)**). The three-stage model has an additional hidden processing stage compared to the two-stage model. An extra node is introduced at the top of the hierarchy to produce an error signal in the third stage of the error coding model.
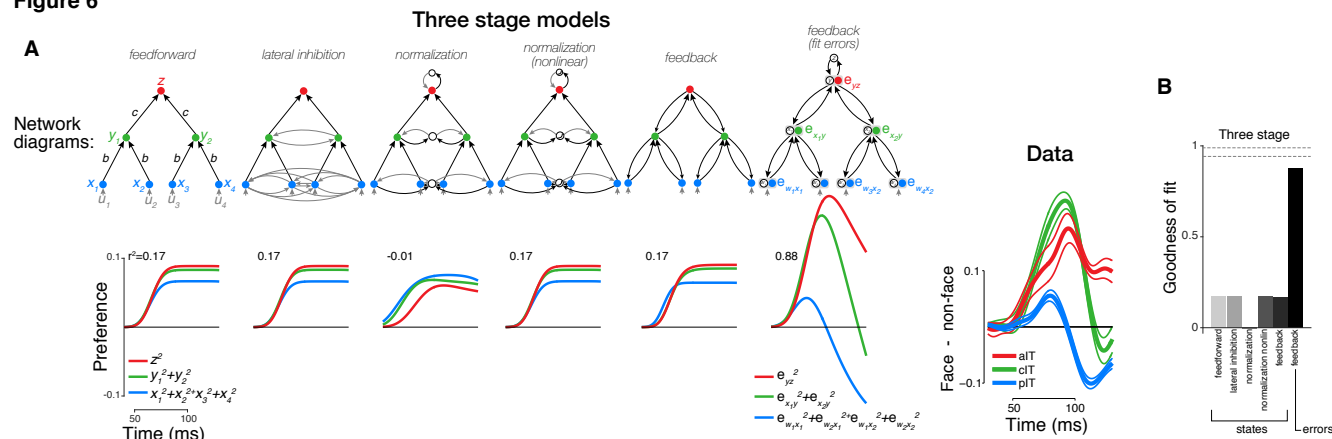1251

**Figure 6**



**Figure 6. Computational modeling of neural dynamics in IT. (A)** Three stage neural networks with recurrent dynamics (network diagrams in top row; see **Methods** and **Figure 5B**) were constructed to model neural signals measured in pIT, cIT, and aIT corresponding to the first (blue), second (green), and third (red) model processing stages. All models received four inputs (gray) into four hidden stage units (blue) which sent feedforward projections that converged onto two units in the next layer (green). Besides this feedforward architecture, additional excitatory and inhibitory connections between units were used to implement recurrent dynamics (self-connections reflecting leak currents are not shown here for clarity; see **Figure 5B** for detailed diagrams). In the five models on the left, the responses of the simulated neurons are assumed to code the current estimates of some set of features in the world (a.k.a states), as is standard in most such networks. The best fit to the population averaged neural data (far right) of the states of each model class are shown (first five columns). These state coding models generally showed increasing selectivity over time from hidden to output layers and did not demonstrate the strong decrease of stimulus preference in their hidden processing stage as observed in the pIT and cIT neural population (blue and green lines). However, the neurons coding errors in a feedback-based hierarchical model did show a strong decrease of stimulus preference in the hidden processing stage (sixth column; reconstruction errors instead of the states were fit directly to the data). This model which codes the error signals (filled circles) also codes the states (open circles) (network diagram in sixth column of top row). **Far right,** population averaged neural selectivity profile for difference between frontal face versus non-face arrangements (normalized by the mean population response to the whole face) used in model fitting. **(B)** Goodness of fit of all three stage models tested to population averaged selectivity profiles (dashed lines represent mean and standard error of reliability of neural data as estimated by bootstrap resampling).
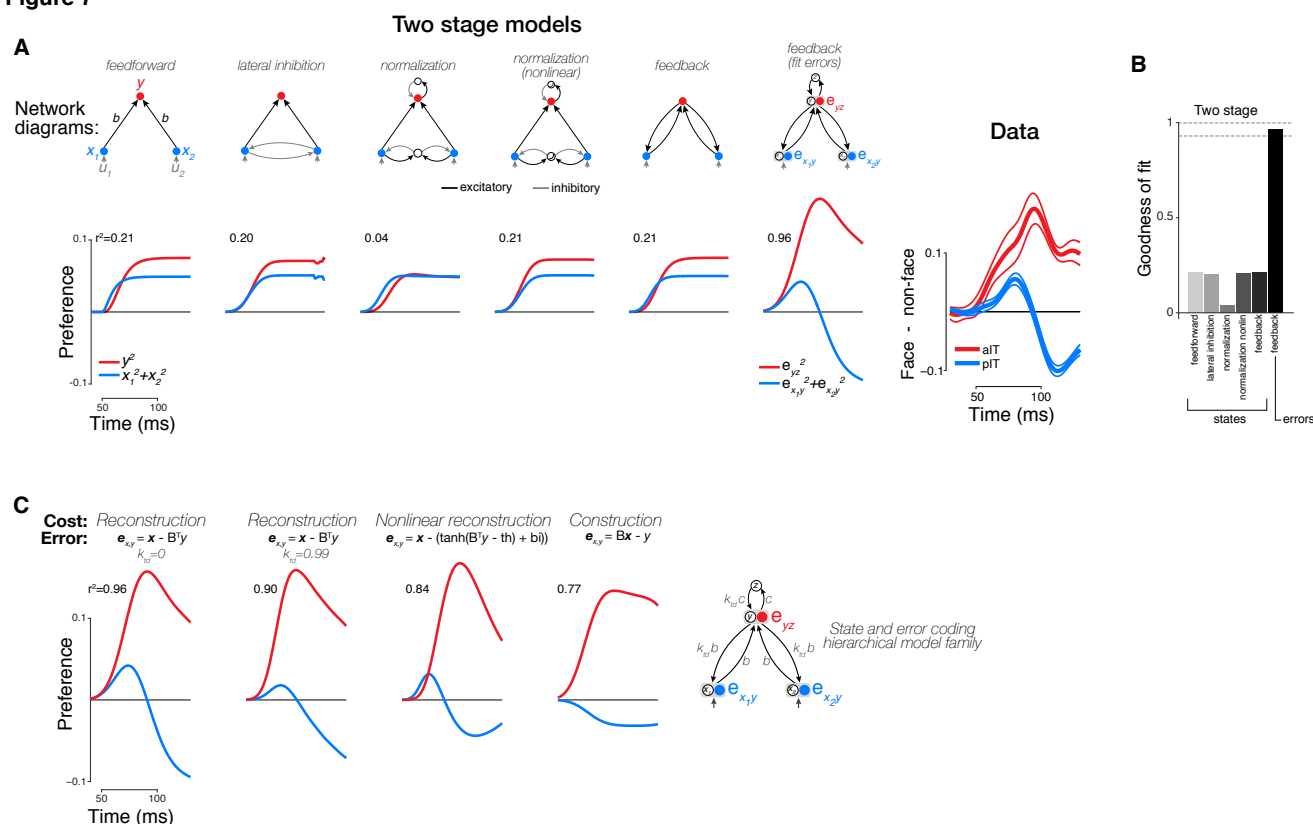
**Figure 7**



1277
1278
1279 **Figure 7**. **Two stage model fits to neural dynamics and comparison of variants of error coding**
1280 **hierarchical models that use different algorithms for online inference. (A)** We tested how well two
1281 stage models could reproduce decreasing selectivity observed in the hidden layer (same format as
1282 **Figure 6A**). The two stage error coding feedback model displayed similar dynamical phenomena as the
1283 data suggesting that three stages of processing were not necessary to obtain error signal dynamics. **(B)**
1284 Goodness of fit to population neural data (right plot in **(A)**) for all two stage models. **(C)** Existing forms
1285 of error computing networks can be distinguished by the type of online inference algorithm that they
1286 use. In one case, inference does not utilize top-down information between stages (classic error
1287 backpropagation; between-stage feedback connections shown are not used in these networks during
1288 runtime). On the other hand, between-stage feedback can be used such as in more general forms of
1289 error backpropagation and predictive coding. We approximated these two extremes by including a
1290 parameter ($k_{td}$, see **Methods**) controlling the relative weighting of bottom-up (feedforward) and top-
1291 down (feedback) evidence during online inference (first and second panels). We found that top-down
1292 inference between stages was not necessary to produce the appropriate error signal dynamics, and $k_{td}$
1293 was equal to zero in our best fitting two-layer and three-layers models (first panel is same model as
1294 two-layer error coding model in **(A)**) although models with $k_{td} \sim 1$ also performed well (second panel).
1295 Models can also differ in their goal (cost function) which directly impacts the error signals required
1296 (equations in top row). Under a nonlinear reconstruction goal (emulating the nonlinear nature of spiking
1297 output), the resulting error signals were still consistent with our data (third column). A simple sigmoidal
1298 nonlinearity, however, did lead to additional details present in our neural data such as a rapid return of
1299 stimulus preference to zero in the hidden layer. When we tested a discriminative, construction goal
1300 more consistent with a supervised learning setting (e.g. classic error backpropagation) where bottom-
1301 up responses simply have to match a downstream target signal in classification tasks, we found that the

1302  errors of construction did not match the data as well as reconstruction errors (compare fourth column to
1303  first three columns) although both types of error outperformed all state-based models and were overall
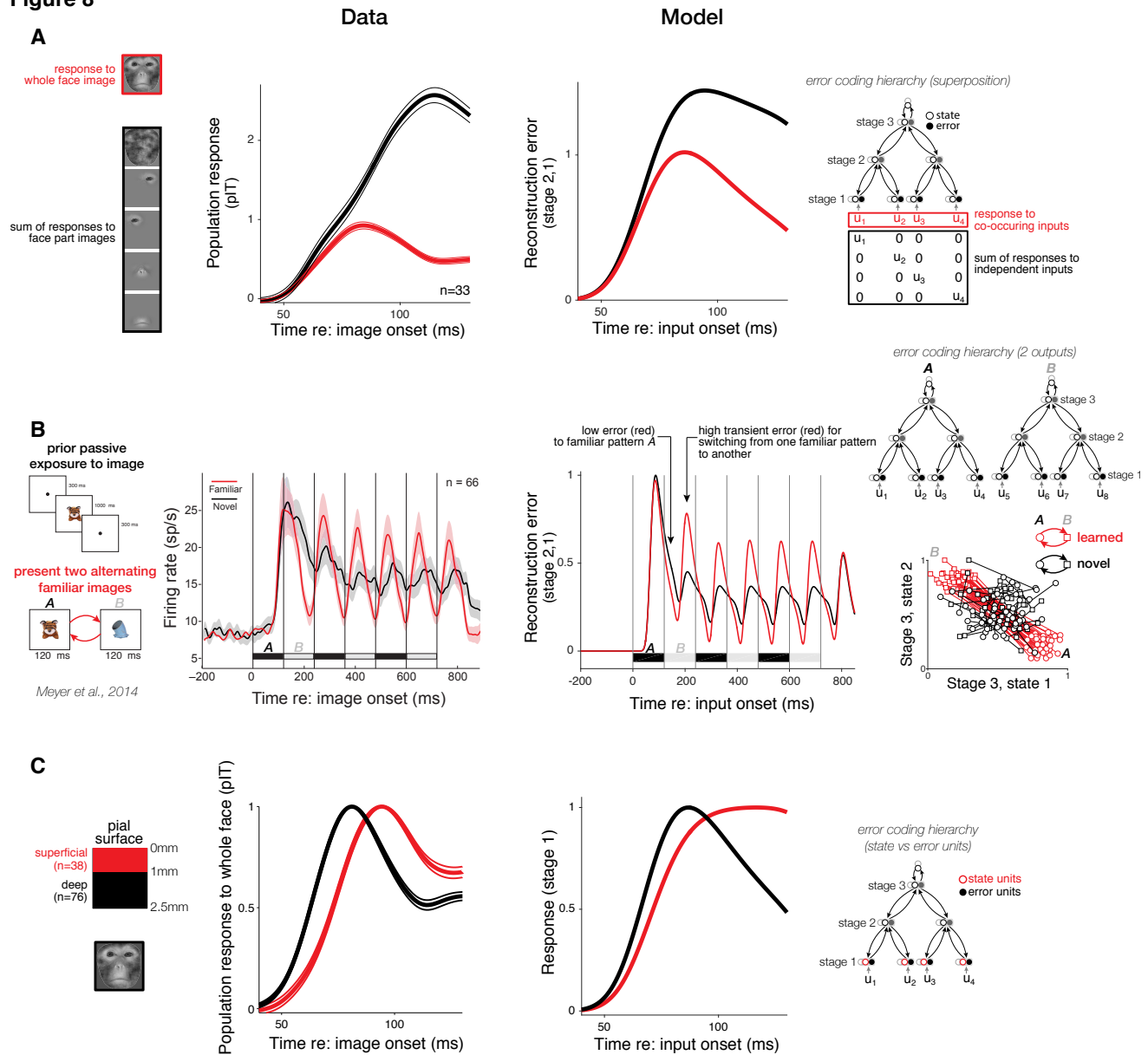1304  very similar (compare to **(A)**).
1305

**Figure 8**



**Figure 8. Comparison of neuronal phenomena to predictions of a feedback-based error coding model.** To generate model predictions, the architectural parameters of the model (i.e. connection weights and time constants) were held fixed, and only the input patterns *u* were varied. **(A)** Neuronal (left panel): in pIT, we found that the sum of the neuronal response to the face parts presented individually (black) exceeded the response to the same parts presented simultaneously (i.e. a whole face, red). Each line is the mean response of 33 pIT sites normalized by the peak response to the whole face. Model (right panel): The magnitude of errors between stage 1 and stage 2 of the model showed a similar degree of sublinear integration (responses were normalized by peak response to the simultaneous input condition). **(B)** We extended the model to include two units in the third, output stage that responded to two learned input patterns (see **Methods**) with increased separation of patterns *A* and *B* in this high-level feature space (red markers in far right panel; 50 draws were made from

1319   distributions for *A* and *B* and were compared to pseudo-randomly drawn inputs, black markers). When
1320   alternating the two learned or familiar patterns *A* and *B*, activations in the top layer of the networks
1321   experienced greater changes than when two randomly selected patterns were alternated (compare
1322   distances traversed in high-level feature space by red lines versus black lines, far right panel). Because
1323   of this large change in high-level activations, transient errors were generated back through the network
1324   for learned patterns. Strong oscillations could be observed in the error signals between stage 1 and
1325   stage 2 of the model for alternated familiar inputs *A* and *B* (120 ms period) (red curve, middle right
1326   panel). In contrast, alternating novel inputs with similar amplitude but with patterns not matching the
1327   learned weights led to small amplitude oscillations in upstream error signals (black curve). Note that the
1328   average response strength of state signals to novel and familiar inputs was matched at all network
1329   levels by construction (mean-matched inputs). The differing response dynamics of error signals under
1330   familiar versus novel patterns are qualitatively consistent with the IT findings for novel versus familiar
1331   images (left, reproduced with permission from *Figs. 1 & 2b* of Meyer et al., 2014). **(C)** The average
1332   response to the whole face for pIT sites recorded 0 to 1 mm below the pial surface (left panel;
1333   superficial recordings, red curve) and for sites 1 to 2.5 mm beneath the pial surface (black curve). (right
1334   panel) Average response of state units (red) and error units (black) in stage 1 of the model. Note the
1335   lagged response in state units which is similar to the lagged response of units in superficial recordings
1336   (red curve, left panel).