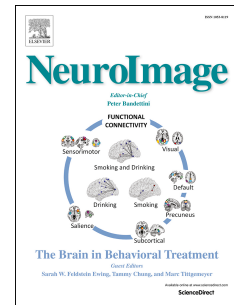


# Accepted Manuscript

What is changing when: Decoding visual information in movies from human intracranial recordings

Leyla Isik, Jedediah Singer, Joseph R. Madsen, Nancy Kanwisher, Gabriel Kreiman



PII: S1053-8119(17)30674-2

DOI: [10.1016/j.neuroimage.2017.08.027](https://doi.org/10.1016/j.neuroimage.2017.08.027)

Reference: YNIMG 14258

To appear in: *NeuroImage*

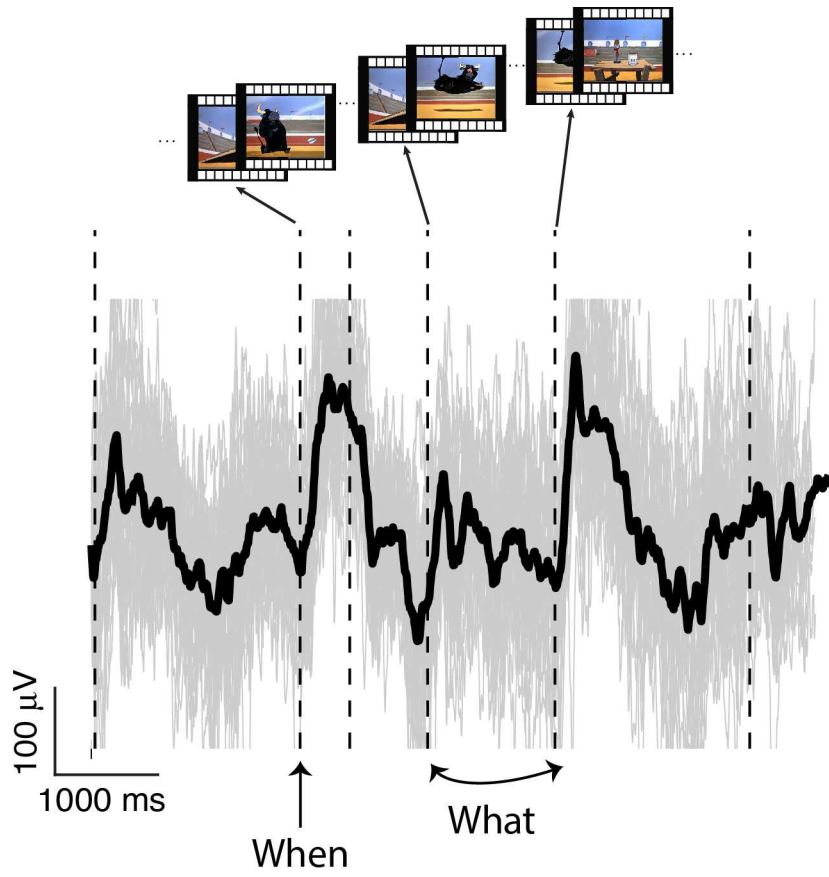
Received Date: 28 April 2017

Revised Date: 7 August 2017

Accepted Date: 8 August 2017

Please cite this article as: Isik, L., Singer, J., Madsen, J.R., Kanwisher, N., Kreiman, G., What is changing when: Decoding visual information in movies from human intracranial recordings, *NeuroImage* (2017), doi: [10.1016/j.neuroimage.2017.08.027](https://doi.org/10.1016/j.neuroimage.2017.08.027).

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



1 **What is changing when: Decoding visual information in movies from human intracranial**  
2 **recordings**

3 Leyla Isik<sup>a,b</sup>, Jedediah Singer<sup>b</sup>, Joseph R. Madsen<sup>c</sup>, Nancy Kanwisher<sup>b</sup>, and Gabriel Kreiman<sup>a</sup>

- 4  
5 a. Department of Ophthalmology, Boston Children's Hospital, Harvard Medical School,  
6 Boston, MA, United States  
7 b. Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology,  
8 Cambridge, MA, United States  
9 c. Department of Neurosurgery, Boston Children's Hospital, Harvard Medical School,  
10 Boston, MA, United States

11  
12 Correspondence should be addressed to Leyla Isik, lisik@mit.edu

13  
14 Number of figures 7

15 Number of supplementary figures 13

16 Number of supplementary tables 5

17 Number of words in abstract 156

18  
19 **Abstract**

20 The majority of visual recognition studies have focused on the neural responses to repeated  
21 presentations of static stimuli with abrupt and well-defined onset and offset times. In contrast,  
22 natural vision involves unique renderings of visual inputs that are continuously changing without  
23 explicitly defined temporal transitions. Here we considered commercial movies as a coarse  
24 proxy to natural vision. We recorded intracranial field potential signals from 1284 electrodes  
25 implanted in 15 patients with epilepsy while the subjects passively viewed commercial movies.  
26 We could rapidly detect large changes in the visual inputs within approximately 100 ms of their  
27 occurrence, using exclusively field potential signals from ventral visual cortical areas including  
28 the inferior temporal gyrus and inferior occipital gyrus. Furthermore, we could decode the  
29 content of those visual changes even in a single movie presentation, generalizing across the  
30 wide range of transformations present in a movie. These results present a methodological  
31 framework for studying cognition during dynamic and natural vision.

32  
33  
34 **Keywords:**

35 Electrocorticography (ECoG), Neural decoding, Object recognition, Ventral pathway, Natural  
36 Vision, Movies

## 37 **1. Introduction**

38 How does the brain interpret complex and dynamic inputs under natural viewing conditions?  
39 The majority of studies in visual recognition have simplified this question by examining neural  
40 responses to isolated shapes, presented in static images, with well-defined onset and offset  
41 times, and reporting averaged neural signals across multiple repetitions of identical stimuli. The  
42 extent to which the principles learned from these studies generalize to the complexities of  
43 temporally segmenting and interpreting the kind of rich, dynamic information present in real-  
44 world vision remains unclear (Felsen and Dan 2005; Rust and Movshon 2005).

45  
46 Studies of the neural responses to flashing static stimuli along the ventral visual stream have  
47 revealed a cascade of computational steps that show progressively increasing shape selectivity  
48 and invariance to stimulus transformations (for reviews, see (Logothetis and Sheinberg 1996;  
49 Riesenhuber and Poggio 1999; Connor et al. 2007; DiCarlo et al. 2012)). The starting point to  
50 analyze the responses to flashed stimuli involves aligning the neural signals to the stimulus  
51 onset, and showing raster plots and post-stimulus time histograms aligned to the transition from  
52 a blank screen to a screen containing the stimulus. Despite the trial-to-trial variability in the  
53 neural responses elicited by repeated presentation of the same stimulus, several studies have  
54 demonstrated that it is possible to read out information about image content in single trials by  
55 applying machine learning techniques (reviewed in (Kriegeskorte and Kreiman 2011)).  
56 Furthermore, it is also possible to identify the time at which the stimulus onset happens purely  
57 from the neural responses (Hung et al. 2005).

58  
59 In stark contrast to experiments that present stimuli with well-defined onsets and offsets, natural  
60 viewing conditions require interpreting the visual world from a continuous stream of visual input.  
61 These conditions present a series of important challenges: (i) there is no obvious “stimulus  
62 onset” to align responses to; (ii) the visual system is continuously bombarded with rapidly  
63 changing input; and (iii) natural images are significantly more complex and cluttered than those  
64 used in many studies with single shapes on a uniform background. In an attempt to begin to  
65 examine how the visual system responds under more naturalistic and dynamic conditions, there  
66 has been growing interest in using movies as stimuli in neurophysiological studies (e.g. (Vinje  
67 and Gallant 2000; Lewen et al. 2001; Fiser et al. 2004; Lei et al. 2004; Montemurro et al. 2008;  
68 Honey et al. 2012; McMahon et al. 2015; Podvalny et al. 2016)) and also in non-invasive studies  
69 (e.g. (Hasson et al. 2004; Bartels and Zeki 2005; Whittingstall et al. 2010; Nishimoto et al. 2011;  
70 Huth et al. 2012; Conroy et al. 2013; Russ and Leopold 2015)).

71  
72 These studies have demonstrated that general principles of visual processing derived from  
73 flashing static stimuli are maintained when considering dynamic stimuli but they have also  
74 highlighted important differences. For example, in primary visual cortex, investigators have  
75 reported that models built from responses to flashed gratings fail to capture all the variance in  
76 the neural responses to movies (Vinje and Gallant 2000; Carandini et al. 2005). In higher visual  
77 areas, the responses to complex shapes such as faces are strongly modulated by the dynamic  
78 aspects of movie stimuli (McMahon et al. 2015; Russ and Leopold 2015).

79

80 Here we describe a methodology to examine neural responses obtained from intracranial field  
81 potentials (IFP) in human epilepsy patients while they passively watched commercial movies.  
82 We tackled the central questions defined above by directly using the neural signals to: (i)  
83 evaluate *when* there are large visual changes in the continuous visual inputs, and hence how to  
84 align neural signals in response to movies, and (ii) identify *what* is the content in the changing  
85 movie frames, despite the complex and heterogeneous variations in the movies. In a first  
86 experiment, we presented multiple repetitions of short movie clips. We showed that we could  
87 decode intracranial field potentials to determine when a visual change happened and identify  
88 what changed in those visual events, generalizing across the transformations present in movie  
89 clips. In a second experiment, we extended this methodology to the analysis of neural  
90 responses to single presentations of a full-length movie.

91

## 92 **2. Material and Methods**

93 Raw data and code for this manuscript are available at  
94 [http://klab.tch.harvard.edu/resources/lsiketal\\_whatchangeswhen.html](http://klab.tch.harvard.edu/resources/lsiketal_whatchangeswhen.html)

95

### 96 **2.1 - Physiology subjects**

97 Subjects were 15 patients (ages 4-36, 8 males, 2 left handed) with pharmacologically intractable  
98 epilepsy treated at Children's Hospital Boston (CHB) or Brigham and Women's Hospital (BWH).  
99 They were implanted with intracranial electrodes to localize seizure foci for potential surgical  
100 resection (Ojemann 1997; Liu et al. 2009). All studies described here were approved by each  
101 hospital's institutional review board and were carried out with the subjects' informed consent.  
102 Electrode types, numbers and locations were driven solely by clinical considerations.

103

### 104 **2.2 - Recordings and data preprocessing**

105 Subjects were implanted with 2 mm diameter intracranial subdural electrodes (Ad-Tech, Racine,  
106 WI, USA) that were arranged into grids or strips with 1 cm separation. Each subject had  
107 between 26 and 144 electrodes ( $86 \pm 26$ , mean  $\pm$  SD). We conducted two experiments (described  
108 below). We studied a total of 1284 electrodes (954 in Experiment I, and 330 in Experiment II,  
109 Supplemental Table 1 and Supplemental Table 2). All data were collected during periods  
110 without seizures or immediately following a seizure. Data were recorded using XLTEK (Oakville,  
111 ON, Canada) and BioLogic (Knoxville, TN, USA) with sampling rates of 256 Hz, 500 Hz, 1000  
112 Hz or 2000 Hz.

113

114 For each electrode, a notch filter was applied at 60 Hz and harmonics, and the common  
115 average reference computed from all channels was subtracted. We focused on the broadband  
116 voltage signals in the 0.1-100 Hz range (referred to as broadband signals throughout the  
117 manuscript). In the Supplementary Material, we also considered the power in the intracranial  
118 field potential signals filtered in the following broadband frequency ranges: alpha (8-15 Hz,  
119 alpha broadband), low gamma (25-70 Hz, low gamma broadband), and high gamma (70-120  
120 Hz, high gamma broadband). All of these are broadband frequency ranges and not single  
121 frequency oscillatory signals. After notch filtering, signals were band passed filtered in each of  
122 those frequency bands. Power in each frequency band was extracted using a moving window

123 multi-taper Fourier transform (Chronux Toolbox, Mitral and Bokil, 2008) with a time-bandwidth  
124 product of five tapers. The window size was 200 ms with 10 ms increments.

125

### 126 **2.3 - Electrode localization**

127 Electrodes were localized by coregistering the preoperative MRI with the postoperative  
128 computerized tomography (CT) (Liu et al. 2009; Destrieux et al. 2010; Tang et al. 2014). For  
129 each subject, the surface of the brain was reconstructed from the MRI and then assigned to one  
130 of 75 anatomically defined regions by Freesurfer. Each surface was also co-registered to a  
131 common brain (Freesurfer fsaverage template) for display purposes only, all analyses  
132 separating electrodes by brain region were based on localization in individual subject's own  
133 anatomical images. We emphasize that all electrode locations are strictly dictated by clinical  
134 criteria. In this type of study, comparisons across subjects are complicated because not all  
135 subjects have electrodes in the same anatomically defined brain region and there are also  
136 differences in electrode locations within each such region across subjects. The locations of the  
137 electrodes in Experiment I are shown in **Figure 4A**, and the locations of the electrodes in  
138 Experiment II are shown in **Figure 7A**. **Tables S4-S5** report the number of subjects contributing  
139 to each anatomically defined brain region in experiment I and II, respectively.

140

### 141 **2.4 – Neurophysiology experiments**

142

#### 143 **2.4.1 - Experiment I**

144 In the first experiment, 11 subjects viewed three 12 s cartoon clips from two separate movies  
145 (example frames for one of these movies are shown in **Figure 1A**). Each clip was repeated  
146 multiple times, between 10-68 repetitions (see Supplemental Table 1), depending on subject  
147 fatigue and clinical constraints. Clips were presented in a random order with a 1 second interval  
148 between clips. Subjects passively viewed the clips. Clips were presented at approximately 4x3  
149 degrees of visual angle. Clips were shown in color and had no sound.

150

#### 151 **2.4.2 - Experiment II**

152 In the second experiment, 4 different subjects viewed a full-length commercial movie: Home  
153 Alone (subject 12, see example frames in **Figure 1B**), Charlie and the Chocolate Factory  
154 (subjects 13-14) or In the Shadow of the Moon (subject 15). Movies were presented with sound  
155 and color at ~18x12 degrees of visual angle. Subjects passively viewed the movies once  
156 through. The movies were interleaved with static images presented for a separate experiment.  
157 The movie was played for 25s, followed by 20 static images from different categories, and then  
158 again by the next 25s of movie.

159

### 160 **2.5 - Eye tracking experiment**

161 Even though the stimulus size was relatively small to prevent large eye movements, we  
162 performed a post-hoc experiment to evaluate whether subjects generate consistent saccades  
163 under these viewing conditions (consistency within a subject across repetitions of the same clip  
164 and consistency across subjects). A post-hoc eye tracking experiment was conducted on 7 in-  
165 lab subjects to examine eye movements. Each subject viewed each 12s clip in Experiment 1,  
166 presented five times in a random order. The viewing conditions were the same as in the

167 physiology experiments. Eye position was recorded with an infrared camera eye tracker  
168 (EyeLink D1000, SR Research). The median eye position across subjects and repetitions is  
169 shown in **Figure S1**.

170

## 171 **2.5 - Data analyses**

172

### 173 **2.5.1 - Cut detection**

174 Movies were segmented based on sharp visual transitions between scenes referred to as movie  
175 cuts throughout (see examples in **Figure 1**). In Experiment I, the cuts within the 12s clips were  
176 manually labeled. In Experiment II, the cuts in the full-length movies were first detected  
177 automatically using an algorithm calculating and thresholding pixel differences between  
178 consecutive movie frames. The automatically detected movie cuts were then checked and  
179 refined manually. We refer to a “shot” as the time period in between two adjacent cuts and we  
180 refer to an “event” as a single occurrence of a shot.

181

### 182 **2.5.2 - Movie labeling**

183 We manually labeled shots in the movies by assigning one label to an entire segment between  
184 movie cuts (shots ranged in length from 0.4s to 3.73 s, with an average length of 1.67s). The  
185 objects and background within a given shot are generally different than those in the previous  
186 shot and are approximately constant throughout a shot.

187

188 In Experiment I, we labeled the presence or absence of the main characters (humanized  
189 versions of cartoon animals) in each 12s clip. This allowed us to test visual selectivity for each  
190 repeated event (e.g. appearance of a particular shape) across the course of the movie. In  
191 particular, in Experiment I, both 12s clips contained shots with a single animal, and shots with  
192 no animal. Two pairs of animal/no-animal scenes were selected in each 12s clip, one pair  
193 occurring at the beginning of the clip and one pair at the end. In the decoding analyses  
194 described below, pairs that were close in time were selected as foils (e.g. each animal shot was  
195 closer in time to its no-animal foil than to the other animal shot) so that the decoding algorithm  
196 could not simply exploit correlations in the physiological data that occur due to temporal  
197 proximity.

198

199 In Experiment II, we labeled in each movie shots with a single face and shots with no faces or  
200 bodies. Faces were selected as a target for visual decoding because they are a consistent,  
201 repeating visual element in all movies shown.

202

### 203 **2.5.3 - Correlation analyses**

204 In Experiment I, we evaluated how consistent the neural signals were across the repeated  
205 presentation of the same 12 s clip for all the cut-responsive electrodes. We correlated the time  
206 courses across repetitions of the same 12 s clip. For each of the n=954 electrodes, we  
207 calculated the Pearson correlation coefficient between each pair of repetitions in every 50 ms  
208 overlapping bin (step size of 1 sample) in each of the three 12 s clips (correlations for an  
209 example electrode during one movie clip are shown in **Figure 2D**). The choice of a 50 ms  
210 window was dictated by the attempt to make the window as small as possible while keeping a

211 sufficient number of sampled voltage values to compute a correlation. To quantify the statistical  
212 significance of the correlation coefficients thus obtained, we defined a null distribution by  
213 computing the correlation coefficients between each 50 ms bin in the movie and random  
214 temporally offset segments. We defined a segment as showing a significant consistency across  
215 repetitions when the correlations between repetitions were significantly above chance in at least  
216 20 consecutive 50 ms bins with  $p < 0.01$  with respect to the null distribution (e.g. horizontal marks  
217 in **Figure 2D**). To examine how the timing of consistent responses across repetitions revealed  
218 by the inter-repetition correlations related to movie cuts, we calculated the latency between the  
219 onset of significantly above chance consistency segments and the previous movie cut (**Figure**  
220 **3**).

221  
222 We repeated the above correlation analyses using a binning window of 400 ms in **Figure S13B,**  
223 **E, H**. This longer time window implies more time points in the calculation of each correlation  
224 coefficient. To ensure that this increase in the number of time points would not bias the results,  
225 we repeated the analyses with a bin size of 400 ms and a smoothing factor of 8 in **Figures**  
226 **S13C, F and I** to match the number of time points in **Figure 3**. Given the larger time window in  
227 the analyses in **Figure S13**, we explicitly removed windows that intersected a camera cut (to  
228 avoid, for example, a window from -200 to +200 ms with respect to a movie cut to be assigned  
229 to -200 ms and erroneously suggest windows of high correlation before movie cuts).

230

#### 231 **2.5.4 - Cut responsiveness**

232 To examine whether the physiological signals showed a significant evoked response to cuts  
233 (e.g. **Figure 2B**), we compared the IFP response, defined as the range (max-min) of the  
234 broadband signals or the total power in each frequency band in the 50 to 400 ms post-cut  
235 window to the corresponding values in the -400 to -50 ms pre-cut window. We defined cut  
236 responsive electrodes as those that showed a  $p < 0.01$  difference in the post-cut versus the pre-  
237 cut windows when considering all repetitions of the  $n=20$  cuts (all cuts, excluding the first cut –  
238 i.e. movie onset – in each movie) based on a permutation test where the pre-cut and post-cut  
239 windows were randomly shuffled 1000 times to define a null distribution. Channels that yielded a  
240 greater IFP response than 99% of the null distribution were defined as significant with  $p < 0.01$ .  
241 All of the electrodes that met this significance criterion are reported in Supplemental Tables 2  
242 through 5 and in Section 3.1.

243

#### 244 **2.5.5 - Decoding methods**

245 Several analyses in the manuscript describe the accuracy in discriminating between visual  
246 events during the movie using statistical classifiers. We describe next the methods for those  
247 analyses.

248

249 *Classifier features* – In each decoding analysis, we considered the average voltage in 50 ms  
250 non-overlapping time bins for each electrode as input to the classifiers described below. In the  
251 Supplementary Material we repeated these analyses examining the average power in the alpha  
252 (8-15 Hz), low gamma (25-70 Hz), or high gamma (70-120 Hz) frequency ranges. Depending on  
253 the specific analyses, we used either single electrodes, pseudo-populations composed of a



254 fixed number of electrodes per region or a population from multiple electrodes selected across  
255 subjects, as described below. The entire decoding procedure was repeated in each 50 ms bin.  
256

257 In Experiment I, because subjects viewed multiple repetitions of identical stimuli, electrodes  
258 were pooled across all subjects into pseudo-populations for specific locations. We first  
259 examined the decoding performance in each brain region by pooling electrodes within a given  
260 anatomical parcel from the Freesurfer Destrieux atlas (Section 2.3). For this analysis, we  
261 considered all anatomical parcels with at least 8 electrodes, and performed decoding with the  
262 pattern of activity across the top 8 electrodes (as measured by the electrode selection  
263 procedure described below) in each of these regions (**Figure 4B-C, Figure 5B-C**). Next, we  
264 also evaluated performance by combining electrodes across separate brain regions and  
265 subjects (**Figure 5D**, (Tang et al. 2014)).  
266

267 In Experiment II, because subjects did not all view the same movie, decoding was performed  
268 separately for each electrode and subject. We calculated decoding performance per brain  
269 region with at least 5 electrodes by averaging the single electrode decoding results for all  
270 electrodes in each anatomical region (**Figure 7B-C**). We also pooled all electrodes per subject  
271 and movie to perform population level decoding, and then again averaged the decoding results  
272 post-hoc across subjects (**Figure 7D**).  
273

274 *Feature pre-processing* - The data from each electrode (feature) was z-scored normalized  
275 based on the mean and standard deviation in the training data. In addition, an ANOVA was  
276 performed on each input feature using only the training data. The ANOVA selects electrodes  
277 that show a larger variance between “categories” than within a “category” as described next. In  
278 **Figures 4B** and **7B**, the ANOVA analysis was used to select those electrodes that showed a  
279 larger variance between cuts and non-cuts compared to the variance within repetition of cuts. In  
280 **Figure 5B**, the ANOVA was used to select electrodes that showed a larger variance between  
281 different movie shots compared to the variance within the same movie shots. In **Figures 5C-D**,  
282 the ANOVA was used to select electrodes that showed a larger variance between shots with an  
283 animal and shots without an animal compared to the variance within shots with an animal and  
284 within shots without an animal. This method has been shown empirically to improve the signal to  
285 noise ratio of decoding with human MEG and monkey LFP time series data (Meyers et al. 2008;  
286 Isik et al. 2014).  
287

288 *Classifier* - Decoding analyses were performed using a maximum correlation coefficient  
289 classifier. This classifier computes the correlation between each test data point and the mean of  
290 all training data points from each class. Each test point is assigned the label of the class of the  
291 training data with which it is maximally correlated (**Figure S12A**).  
292

293 *Cross-validation* - For each decoding run, the data were divided into 10 cross-validation splits.  
294 Feature pre-processing (z-scoring and ANOVA) was performed on 9 out of 10 of the splits, and  
295 testing was performed on the 10th held out split.  
296

297 The decoding at each time bin was repeated for 20 times, each with a different train/test data  
298 split. The average performance of the 20 decoding runs is displayed as “classification accuracy”  
299 as a function of time from cut onset in **Figures 5D** and **7D**. In other cases, we summarized  
300 classification accuracy by reporting the average value from 50 to 400 ms post-cut onset  
301 (**Figures 4B, 5B-C, 7B-C**).

302

### 303 *Decoding analyses, Experiment I*

304 (i) We compared movie segments with a movie cut versus random segments falling at least 400  
305 ms away from a movie cut (**Figure 4B**).

306

307 (ii) We evaluated whether we could detect visual transitions in the entire 12 second clip. The  
308 procedure is illustrated in **Figure S12B**. We used the average vector representing “cut” and “no-  
309 cut” events as described in (i) and **Figure S12A**. For each 50 ms window from held-out  
310 repetitions, if the correlation with the “cut” vector was larger than the correlation with the “no-  
311 cut”, we assigned a label of +1, otherwise we assigned a label of -1. We defined hits as those  
312 50 ms windows which had a label of +1 and which were within the 0 to 400 ms after a cut.  
313 Similarly, we defined false alarms as those 50 ms windows which had a label of +1 and which  
314 did not occur within 0 to 400 ms after a cut. We calculated the d prime measure across all 50  
315 ms time bins in the 12s clip:  $d' = Z(\text{hit rate}) - Z(\text{false alarm rate})$ , where Z is the inverse  
316 cumulative distribution function (**Figure S12B, Figure 4C**). We defined a predicted visual  
317 transition as a set of 1 or more continuous 50 ms windows classified as +1. For each predicted  
318 visual transition, we defined the time of the transition as the first 50 ms window in the set. We  
319 calculated how far away those predicted visual transitions were from the nearest prior cut in  
320 **Figure 4D**.

321

322 (iii) We tested for visually selective signals by decoding the different camera shots from each  
323 other (**Figure 5B**). We included the 13 camera shots in the first two movies (all the movie cuts  
324 that were presented at least 20 times across subjects, see **Supplemental Table 1, Figure S6**,  
325 excluding the first and last shot).

326

327 (iv) We compared shots with an animal versus shots without an animal (**Figure 5A,C-D**). We  
328 performed this animal versus no animal decoding first across repetitions of the same movie clips  
329 (referred to as the “within shot” condition). Next, we decoded across shots in the same 12s clips  
330 (referred to as the “across shot” condition), and finally and across shots from different movie  
331 clips (referred to as the “across clip” condition).

332

### 333 *Decoding analyses, Experiment II*

334 (i) We compared movie segments with a movie cut versus random segments falling at least 400  
335 ms away from a movie cut, as in experiment I (**Figure 7B**).

336 (ii) We compared shots with a single face versus shots with no face (**Figure 7C-D**).

337

## 338 **3. Results**

339 We investigated the neurophysiological responses elicited by dynamic movie stimuli by  
340 recording intracranial field potential (IFP) signals from 1324 electrodes implanted in 15 patients

341 with epilepsy (**Tables S1-S3**). We conducted two experiments: (i) Experiment I consisted of  
342 repeated presentation of three 12s commercial cartoon movie clips (**Figure 1A**, 954 electrodes);  
343 (ii) Experiment II consisted of a single presentation of a full-length commercial movie (**Figure**  
344 **1B**, 370 electrodes).

### 346 **3.1 - Neurophysiological responses to time-varying stimuli (Experiment I)**

347 In multiple Visual Neuroscience experiments, stimuli are flashed with well-defined onset and  
348 offset times and responses are analyzed by aligning activity to the appearance of a stimulus.  
349 Movies, as a coarse proxy to natural vision, lack those stimulus onsets. We conjectured, with  
350 others (McMahon et al. 2015), that the drastic changes between consecutive frames that occur  
351 during movie cuts provide a strong temporal demarcation. **Figure 1** shows two examples of  
352 movie cuts (transition from frame 130 to 131 in **Figure 1A** and from frame 15869 to 15870 in  
353 **Figure 1B**) and the accompanying large changes in the visual field. We set out to investigate  
354 whether such movie cuts trigger the onset of physiological responses and can thus be used to  
355 demarcate visual events in movies.

356  
357 We started by aligning the IFP signals to movie cuts. **Figure 2B** shows the responses of an  
358 example electrode located in the right inferior occipital gyrus (**Figure 2A**) that demonstrated a  
359 vigorous modulation after one of the movie cuts. The changes in IFP were evident in almost  
360 every single repetition of the movie clip, showed a consistent latency of approximately 100 ms  
361 after the cut and were short-lived, with the voltage returning to baseline within approximately  
362 400 ms after the cut. This electrode showed an appreciable modulation elicited by most, but not  
363 all, the cuts in the 12s clips (**Figure 2C**). To further quantify the modulation in IFP, we computed  
364 the degree of consistency in the responses evaluated by the Pearson correlation coefficient  
365 between the voltage time series for every possible pair of repetitions, using a window of 50 ms  
366 (**Figure 2D**). The correlation coefficient largely hovered around zero, indicating that the IFP  
367 signals were inconsistent across repetitions, except for sharp spikes in correlation, which were  
368 typically evident right after a movie cut. For the example electrode in **Figure 2** and movie clip 1,  
369 there was a significant increase in consistency after 9 of the 10 movie cuts.

370  
371 We defined an electrode as visually responsive if the range (max-min) of the broadband IFP  
372 signals from 50 to 400 ms after a movie cut was significantly different from the range from -400  
373 ms to -50 ms before a movie cut, using all cuts across the 3 movie clips ( $p < 0.01$  permutation  
374 test, **Section 2.5.4**, similar criteria have been used in other work, e.g., (Liu et al. 2009)). In the  
375 Supplementary Material, we report the results obtained by evaluating modulation in the alpha  
376 (8-15 Hz), low gamma (25-70 Hz) and high gamma (70-120 Hz) bands of the IFP signals.

377  
378 Using these criteria, out of the total of 954 electrodes in Experiment I, we obtained 51  
379 electrodes, which were mostly located in the occipital pole, and inferior and middle occipital gyri  
380 and, to a lesser degree, in the fusiform gyrus, medial lingual gyrus, and inferior temporal gyrus  
381 (**Table S4**). In order to avoid potential physiological changes elicited by eye movements, we  
382 kept the stimuli relatively small (~4 x 3 degrees) and we restricted the analyses to the initial  
383 neurophysiological response between 50 and 400 ms. Furthermore, we conducted a separate  
384 post-hoc experiment in non-epilepsy subjects to measure eye movements under the same

385 stimulus presentation conditions and we did not observe any consistent eye movements elicited  
386 by the movie cuts (**Figure S1**). Therefore, it seems more likely that the modulatory changes in  
387 the physiological signals were triggered by the large changes in the visual stimulus rather than  
388 by large saccadic eye movements. Reliable responses triggered by movie cuts were also  
389 evident in other frequency bands, an example in the high gamma band is shown in **Figure S2**.

### 391 **3.2 – Responses that were reproducible between repetitions largely clustered shortly** 392 **after movie cuts**

393 We next sought to evaluate the degree of trial-to-trial reproducibility in the physiological  
394 responses across the entire 12s clip and the whole set of electrodes in our sample. We plotted  
395 the statistical significance of the correlation coefficient over the entire 12s clips in each electrode  
396 on the Freesurfer fsaverage template brain (**Figure S5A**). Multiple electrodes along the ventral  
397 stream showed reliable responses (**Table S4**, see **Figures S5B-D** for the results in other  
398 frequency bands). As illustrated for the example electrode in **Figure 2**, the increase in  
399 correlation between repetitions was largely present in the initial ~300 ms after cut onset. We  
400 followed the procedure in **Figure 2D** to detect segments with statistically significant correlation  
401 between repetitions. The majority of consistent responses fell within ~300 ms of a movie cut  
402 (**Figure 3A**). Throughout the entire population of electrodes, there was a small number of  
403 consistent responses occurring >500 ms away from movie cuts (**Figure 3A**). For example, there  
404 was a small but statistically significant peak before the 3rd cut and another small non-significant  
405 peak before the 5th cut in **Figure 2D**. However, the degree of consistency, as quantified by the  
406 correlation coefficient between repetitions, showed a small drop with the time from movie cut  
407 onset (**Figure 3B**). Moreover, the duration of those segments showing consistency between  
408 repetitions also showed a small decrease as a function of time from the previous cut (**Figure**  
409 **3C**). To further illustrate this point, we searched in the entire electrode sample for two example  
410 electrodes with the most reliable response segments that were more than 400 ms away from a  
411 movie cut (**Figure S4**). Even though the peaks in **Figure S4** represent the strongest examples,  
412 they are still weaker and shorter than those illustrated in **Figure 2D**. Similar conclusions were  
413 drawn when considering other frequency bands (**Figure S3**). The correlation coefficients in  
414 **Figure 3** were calculated using a window size of 50 ms; similar conclusions were reached when  
415 considering a window size of 400 ms (**Figure S13**). In sum, the abundance, strength and  
416 duration of consistent responses was largely linked to the occurrence of movie cuts.

### 417 418 **3.3 - Detecting the presence of movie cuts (Experiment I)**

419 Under natural viewing conditions, in the absence of a blank screen followed by a flashed  
420 stimulus, the brain needs to determine *when* there is a visual change and *what* that change  
421 consists of. The when and what computations need to take place in single events, without  
422 averaging. To evaluate whether the neural signals are able to discriminate the timing of changes  
423 in the visual world, we built machine learning classifiers to discriminate between movie  
424 segments (350 ms duration) containing a movie cut versus movie segments without a movie cut  
425 (**Figure 4B**). The control movie segments consisted of random time periods that were at least  
426 400 ms away from a cut. We built pseudopopulations of electrodes in different anatomically  
427 defined brain regions that contained at least 8 electrodes by pooling data across all patients  
428 (**Figure 4A**, Methods). In each region, we used the 8 most selective electrodes per region (as

429 determined by an ANOVA applied to the training data, see **Methods**). We report the  
430 classification accuracy, i.e., the proportion of repetitions where the machine learning classifier  
431 correctly determined the presence or absence of a movie cut (chance = 0.5). Of the 25 regions  
432 with at least 8 electrodes (**Table S4**), 5 regions showed significantly above chance classification  
433 accuracy: inferior occipital gyrus, fusiform gyrus, middle occipital gyrus, inferior temporal gyrus  
434 and occipital pole. The average classification accuracy across these 5 regions was  $0.62 \pm 0.04$   
435 (mean $\pm$ SD, across regions; see **Supplemental Figure 7A-C** for the corresponding classification  
436 results using IFP signals filtered in different frequency bands).

437  
438 Whereas the analysis in **Figure 4B** compares 350 ms segments with and without cuts, the brain  
439 needs to be able to detect those events in single events and during a continuous stream. Next,  
440 we developed a classifier to investigate whether it is possible to detect visual transitions in  
441 single events during the entire 12s clips (Methods). The procedure is schematically described in  
442 **Figure S12B**. This classifier continuously determines whether there is a visual transition, thus  
443 making correct detections (hits) as well as false ones (false alarms). We evaluated the accuracy  
444 of this continuous prediction by measuring the classifier's sensitivity using  $d'$  prime. We found  
445 that classifiers using data from five of the seven regions described in **Figure 4A** (excluding the  
446 middle temporal gyrus and the occipital pole) detected visual transitions with above chance  
447 precision with an average  $d'$  of  $0.48 \pm 0.18$  (mean $\pm$ SD across significant regions, **Figure 4C**). We  
448 estimated the latency of these visual transition predictions by measuring the time difference to  
449 the nearest prior cut; the mean latency was  $690 \pm 610$  ms (mean $\pm$ SD, across all significant  
450 regions, **Figure 4D**). The distribution of these time differences was significantly different from  
451 the one expected under the null hypothesis defined by 10,000 runs of randomly selecting the  
452 same number of time points per movie as predicted transitions (**Figure 4D**, black line,  $p < 10^{-10}$ ,  
453 Kolmogorov-Smirnov test; see **Supplemental Figure D-I** for the corresponding classification  
454 results using IFP signals filtered in different frequency bands). In sum, it is possible to detect  
455 when the image changes within a continuous stream from the neural responses along the  
456 ventral visual stream.

### 457 458 **3.4 - Decoding visual events in movies (Experiment I)**

459 After detecting when there is a visual transition in the movie, we asked whether it is possible to  
460 selectively identify *what* visual event changes occur. To address this question, we assessed  
461 whether the neural signals could discriminate among the 350 ms windows (ranging from 50 to  
462 400 ms) after movie cuts. We selected 13 movie cuts that were presented at least 20 times  
463 (**Figure S6, Methods**). Of the 25 regions with at least 8 electrodes, we found 7 regions that  
464 showed above chance classification accuracy based on a  $p < 0.01$  permutation test (**Figure 5B**).  
465 These 7 regions included the 5 regions described in **Figure 4B** and also the medial lingual  
466 gyrus and middle temporal gyrus. The average classification accuracy across these 7 regions  
467 was  $0.27 \pm 0.07$  (chance =  $1/13 = 0.08$ ; see **Supplemental Figure 8A-C** for the corresponding  
468 classification results using IFP signals filtered in different frequency bands).

469  
470 The results in **Figure 5B** show classification accuracy averaged from 50 to 400 ms with respect  
471 to movie cuts. To summarize and visualize dynamic changes in classification accuracy as a  
472 function of time, we pooled electrodes across all subjects and selected those electrodes that

473 showed larger variation across the 13 movie cuts than within repetitions of the same movie cut  
474 using only training data (described under feature selection in **Methods**). We performed the  
475 same 13-way cut classification analysis described in **Figure 5B**. This analysis shows that  
476 classification accuracy started to increase at around 100 ms after a movie cut and peaked at  
477 around 400 ms (**Figure S8D**), consistent with the example electrode dynamics shown in **Figure**  
478 **2** and also with previous work decoding different objects with static images (Liu et al. 2009;  
479 Tang et al. 2014). **Figure S8D** shows that classification accuracy was also high at  $t=0$ , and even  
480 *before* the onset of the movie cut. Unlike experiments where static images are presented in  
481 random order and are preceded by a blank screen, in the movie presentation, the visual  
482 stimulus preceding a movie cut was always the same across different repetitions. Furthermore,  
483 several movie cuts were preceded by another movie cut within a few hundred ms (e.g. cut  
484 numbers 2 and 3 in movie clip 1, **Figure S6**), contributing to the significant classification  
485 accuracy before and at  $t=0$  in **Figure S8D**.

486

### 487 **3.5 - Invariant decoding of visual events in movies (Experiment I)**

488 A central challenge in visual recognition involves combining selectivity to different shapes with  
489 invariance to the myriad transformations in those shapes (Booth and Rolls 1998; Riesenhuber  
490 and Poggio 1999; Serre et al. 2007; DiCarlo et al. 2012). After identifying when visual transitions  
491 occur and what changes during each event, we asked whether these visual shape-selective  
492 signals generalize across transformations in the stimuli. To test the degree of invariance in the  
493 visual shape-selective responses, we labeled the content of each shot with the presence or  
494 absence of a cartoon humanized animal. We selected four animal/no-animal shot pairs (from  
495 movies 1 and 2, **Figure S6, Methods**), and used the same methodology described above to  
496 determine in each event whether an animal was present or not, with varying amounts of  
497 generalization described next (**Figure 5A**).

498

499 First, the classifier was trained on a subset of the repetitions and tested on the remaining  
500 repetitions of the same shots (“within shot”, **Figure 5C**, blue bars), requiring generalization  
501 across different repetitions of identical stimuli (similar to **Figure 5B**, here using a subset of the  
502 shots for comparison with the next set of analyses and specifically distinguishing shots  
503 containing an animal versus shots not containing an animal, chance = 0.5). As expected from  
504 the previous analyses, in **Figure 5C** we observed significant classification accuracy in 6 of the 7  
505 regions described in **Figure 4A** (the medial lingual gyrus did not reach statistical significance in  
506 this analysis). The mean within-shot classification accuracy in these 6 regions was  $0.74 \pm 0.06$   
507 (mean  $\pm$  SD across regions).

508

509 Next, we evaluated the degree of generalization across different shots containing an animal  
510 within the same movie (“across shots”, **Figure 5C**, red bars). To avoid conflating tolerance to  
511 different shots with correlated activity in time, each animal versus no animal pair was selected to  
512 be closer in time to each other than to the second animal versus no-animal pair (i.e., each shot  
513 containing an animal was closer in time to its no-animal foil shot than to the other animal  
514 containing shot, **Figure S6**). This analysis revealed significant classification accuracy in 4 of the  
515 7 regions described in **Figure 4A**: inferior occipital gyrus, fusiform gyrus, inferior temporal gyrus

516 and occipital pole. The mean within-shot classification accuracy in these 4 regions was  
517  $0.71\pm 0.05$ .

518  
519 Finally, we considered the most extreme case of visual generalization by asking whether we  
520 could train a classifier to discriminate shots containing an animal or not in one movie and test it  
521 on a different movie ("across clip", **Figure 5C**, green bars). Three brain regions, inferior occipital  
522 gyrus, fusiform gyrus and inferior temporal gyrus, yielded significant classification accuracy with  
523 an average performance of  $0.68\pm 0.07$ .

524  
525 To summarize and visualize the temporal dynamics in classification accuracy, we followed the  
526 procedure described in the previous section for **Figure S8D-G** and combined electrodes across  
527 all subjects in **Figure 5D**. The dynamics revealed an increase in classification accuracy  
528 commencing around 100ms post cut onset and peaking around 400ms post cut onset. As noted  
529 in **Figure S8D-G**, the within-shot condition (blue curve) also revealed strong classification  
530 accuracy at and before cut onset in **Figure 5D**. **Figure S9** presents corresponding results  
531 examining IFP signals filtered in different frequency bands. In sum, the results presented in the  
532 previous section and this section show that we can selectively extract information about what  
533 changes in the image in single events and with a considerable degree of invariance to the pixel-  
534 level transformations.

### 535 536 **3.6 - Detecting the presence of movie cuts in single presentation of movies (Experiment 537 II)**

538 The insights and analyses derived from Experiment I relied on multiple repeated presentations  
539 of the same identical movies. Under natural viewing conditions, the brain must rely strictly on  
540 unique presentations of single events. **Figure 5C-D** showed that it was possible to decode the  
541 presence of absence of an animal by generalizing across different shots and even different  
542 movie clips. However, all the classifiers in **Figure 5C-D** were still trained using multiple  
543 repetitions of identical stimuli. As a more stringent test of generalization across events, we  
544 conducted Experiment II where subjects passively viewed a single repetition of a full-length  
545 commercial movie (**Methods**). In lieu of identical stimulus repetitions, we leverage the repetition  
546 of similar visual events across the duration of a movie.

547  
548 We assessed whether it was possible to detect when large visual changes occurred in the full-  
549 length movies. As described in **Figure 2**, neural signals showed strong changes in voltage  
550 shortly after movie cuts in the full-length movie. **Figure 6** illustrates the responses of an  
551 example electrode located in the left occipital pole that showed consistent (but not identical)  
552 changes after almost every movie cut (see raster plot depicting every movie cut in **Figure 6B**),  
553 despite the fact that the cuts vary enormously in content and were only shown once (see also  
554 **Figure S10**). The voltage deflections commenced approximately 100 ms after a movie cut  
555 (**Figure 6C**). In total, we found 61 (out of 330 total) cut-responsive electrodes (Section 2.5.4),  
556 located primarily in the cuneus, medial lingual gyrus, fusiform gyrus, inferior occipital gyrus and  
557 occipital pole (**Table S5**).

558

559 Following the procedure used in **Figure 4B**, we evaluated whether we could distinguish a  
560 segment from 50 to 400 ms post cut onset from random time points in single events (**Figure**  
561 **7B**). Because of the smaller total number of electrodes in Experiment II, we considered regions  
562 with at least 5 electrodes (as opposed to the threshold of 8 electrodes used in **Figures 4** and **5**).  
563 Also, because subjects watched different full-length movies, we did not build pseudo-  
564 populations combining electrodes in the same labeled region across subjects. Instead, we used  
565 single electrodes and report average classification accuracy for single electrodes in **Figure 7B**  
566 (whereas **Figure 4B** is based on a pseudopopulation of 8 electrodes in each region). Of the 20  
567 regions with at least 5 electrodes, we observed a small but significant classification accuracy in  
568 4 regions: inferior occipital gyrus, cuneus, medial lingual gyrus and occipital pole (see **Figure**  
569 **S11A-C** for the corresponding analyses after filtering the IFP signals in different frequency  
570 bands).

571  
572 Not all the same regions were interrogated in the different subjects that participated in  
573 Experiment I and II (**Tables S4** and **S5** provide detailed information about electrode locations in  
574 the two experiments). All of the 7 regions described in **Figure 4A** had enough coverage to be  
575 considered in Experiment II. Three of these regions - inferior occipital, medial lingual gyrus and  
576 the occipital pole – showed significant classification accuracy to detect the presence of movie  
577 cuts in both experiments, while the other four regions did not reach significant classification  
578 accuracy in Experiment II. In addition, the cuneus showed significant classification accuracy to  
579 detect the presence of movie cuts in Experiment II but not in Experiment I.

### 580 **3.7 - Invariant decoding of visual events in single presentation of movies (Experiment II)**

581 Following the steps in Experiment I, we next asked whether we could decode *what* changed in  
582 the image at a given movie cut. We trained the classifier to distinguish those shots containing a  
583 face from shots that did not contain a face following the procedures in **Figure 5**, with two  
584 important differences. First, given the extensive preponderance of frames including human  
585 faces in the full-length movies in Experiment II, we labeled each shot as containing a face or not  
586 (as opposed to the animal faces in Experiment I, **Methods**). Second, as described above, we  
587 also considered single electrodes and report average classification accuracy in **Figure 7C**, as  
588 opposed to results based on pseudopopulations. Of the 5 regions described in **Figure 7B**, we  
589 could discriminate with small but significant classification accuracy shots containing a face from  
590 those with no face from single electrodes in the inferior occipital gyrus. Additionally, the fusiform  
591 gyrus also showed even smaller but still significant classification accuracy (see **Figures S11D-F**  
592 for the corresponding analyses considering IFP signals filtered in different frequency bands.  
593  
594

595 To summarize the temporal dynamics in classification accuracy, we followed the procedure  
596 described in **Figure S8D, 5D** for Experiment I and combined electrodes across all subjects in  
597 **Figures 7D**. Again, because subjects watched different full-length movies, we did not combine  
598 electrodes across subjects but instead built pseudo-populations using each subjects' electrodes  
599 and averaged the four subjects' classification accuracies post-hoc. There was an increase in the  
600 classification accuracy to detect the presence or absence of a face starting slightly before 200  
601 ms post cut onset and peaking around 300ms post cut onset (**Figure 7D**; see **Figures S11G-I**  
602 for the corresponding analyses considering IFP signals filtered in different frequency bands). In



603 sum, the previous section and this section demonstrate that the results obtained in Experiment I  
604 extrapolate to the conditions in Experiment II, whereby we can discriminate when there are  
605 visual changes and what those visual changes consist of in single presentations of a full-length  
606 movie.

607

#### 608 **4. Discussion**

609 Parsing a continuous stream of visual stimuli is a fundamental challenge for the visual system.  
610 Here we considered commercial movies as a coarse proxy for natural visual input and described  
611 a methodology to extract visual information from invasive physiological recordings from the  
612 human brain during a continuous movie. Intracranial field potentials recorded along the ventral  
613 visual stream showed strong modulation approximately 100 ms after movie cuts, defined as  
614 discontinuous changes from one frame to the next (**Figure 2**). Such vigorous physiological  
615 responses allowed us to detect *when* there are visual changes during the continuous stimulus  
616 (**Figure 4B-D**). By aligning the responses to those changes, we identified *what* visual  
617 information was present in each shot (e.g., shots with or without an animal), generalizing across  
618 different events within the same movie or even across different movies (**Figure 5**). We further  
619 demonstrated that these findings extend to detecting the timing of visual changes and decoding  
620 events in a single presentation of a full-length movie (**Figures 6-7**).

621

622 We separately considered broadband signals from 0.1 to 100 Hz and broadband, band-limited  
623 signals in the alpha (8-15 Hz), low gamma (25-70 Hz) and high gamma (70-120 Hz) bands. We  
624 observed fewer and weaker visual responses in the alpha band, consistent with previous  
625 studies (e.g. Bansal et al 2012). The qualitative and conceptual conclusions derived from  
626 examining the low and high gamma band were consistent with those based on the broadband  
627 signals. Yet, there were quantitative differences in terms of the numbers of responsive  
628 electrodes, classification performance and, in some cases, the specific areas that showed  
629 significant decoding accuracy. These differences are discussed in further detail in the  
630 Supplementary Material. These qualitative similarities and quantitative differences between  
631 broadband and gamma band responses were noted in several previous studies (e.g., (Vidal et  
632 al. 2010; Privman et al. 2011; Bansal et al. 2012; Miller et al. 2014)).

633

634 Commercial movies such as the ones used here and in other studies clearly constitute artificial  
635 stimuli that are different from natural viewing conditions. Movies are commercial forms of art  
636 specifically and carefully designed to evoke strong emotional experiences, producing  
637 memorable audiovisual scenes in a compressed time frame beyond the occurrences of  
638 everyday life. Movie cuts are introduced in videos by the director to manipulate spatial  
639 coordinates, context, attention, and interactions (Dudai 2012; Smith et al. 2012). These cuts  
640 only constitute a first order approximation to the type of discontinuities that arise under natural  
641 viewing conditions as a result of sudden changes in moving objects, occlusion, lighting and  
642 internally dictated changes such as eye movements. Despite these caveats, movies provide a  
643 rich stimulus for probing neural responses in situations where the brain is continuously subject  
644 to incoming inputs, as opposed to a blank screen followed by the onset of a picture. Indeed,  
645 several previous studies have demonstrated that sharp transitions between frames in movies  
646 can trigger a strong neural response all along ventral visual cortex from early visual areas (Vinje

647 and Gallant 2000; Montemurro et al. 2008) to the highest visual areas (Privman et al. 2007;  
648 Honey et al. 2012; McMahon et al. 2015).

649

650 Critically, the brain must be able to capture these dynamic transitions in single events without  
651 averaging responses over multiple repetitions. Even with the type of coarse signals and limited  
652 spatial sampling considered here, it is possible to detect visual changes in a movie within  
653 approximately 100 ms of those changes (Figures 2, 4, and 6). These latencies are close to  
654 those reported in monkey and human ventral visual cortex in response to static images  
655 (Richmond et al. 1990; Rolls and Tovee 1995; Keysers et al. 2001; Hung et al. 2005; Liu et al.  
656 2009). Thus, our intuitions about the initial dynamics of neural responses triggered by flashing  
657 static pictures seem to extrapolate to dynamic and continuous viewing conditions.

658

659 The rapid field potential changes were elicited by most movie cuts and were consistent  
660 throughout tens of repetitions. Intriguingly, we observed few consistent physiological responses  
661 across repetitions outside of movie cuts (**Figure 2D**, **Figure 3**). In other words, we largely failed  
662 to note consistent responses from one repetition of the movie clip to another except within a few  
663 hundred milliseconds after a movie cut. There are several non-exclusive possibilities for this  
664 observation. First, our sampling of brain locations was far from exhaustive. The electrode  
665 locations were strictly dictated by clinical criteria. Although we interrogated a relatively large  
666 number of brain regions for this type of study (almost 1,000 different electrodes distributed over  
667 46 brain regions, **Table S4**), there could well be many other brain loci that show consistent  
668 responses to other aspects of the movies unrelated to the movie cuts. Second, we studied  
669 coarse field potential signals recorded from low-impedance electrodes that capture neural  
670 activity over vast numbers of neurons (Buzsáki et al. 2012). It is quite possible that there are  
671 strong neuronal responses to other aspects of the movies that are not captured by field potential  
672 signals. Third, it is conceivable that other aspects of cognition beyond visual processing are  
673 modulated or even governed by different mechanisms that do not lead to the type of sharp and  
674 consistent responses illustrated in **Figure 2**. In particular, other aspects of cognition beyond  
675 visual processing during a movie may not have a well-defined temporal onset (e.g. when exactly  
676 emotions are triggered during a scene), or they may show rapid adaptation (e.g. the first viewing  
677 of a movie scene might trigger stronger emotions than the tenth viewing), both of which would  
678 reduce the reproducibility of these signals across multiple trials. In sum, while we argue here  
679 that we can rapidly decode visual transitions in single events during a movie, there remain  
680 important questions about how to study the neural basis of higher cognitive functions under  
681 natural conditions.

682

683 In the absence of fixed image onset times or movie cuts, the brain must segment continuous  
684 information into discrete visual events. How are visually evoked signals aligned under natural  
685 viewing conditions? Several sources in the brain could in principle provide an internal alignment  
686 signal to the ventral visual stream, including a copy of a motor efferent from eye movements, or  
687 external object movement onset information conveyed by the dorsal stream. While this study  
688 does not explain the mechanistic origin for the physiological changes triggered by movie cuts,  
689 the results presented here show that it is possible to align and interpret signals directly from the  
690 field potentials recorded from electrodes in the ventral visual stream. During natural viewing

691 conditions, we speculate that signals along the ventral visual stream may be sufficient to  
692 interpret what changes when without the need for additional sources of information.

693  
694 The main regions along the ventral visual stream that contributed to decoding when and what  
695 information included the inferior occipital gyrus, the fusiform gyrus, the inferior temporal gyrus  
696 and the occipital pole (**Figures 4 and 7**). All of these regions have also revealed selective visual  
697 responses in previous invasive human neurophysiology studies (e.g. (Privman et al. 2007; Liu et  
698 al. 2009; Vidal et al. 2010)). These areas are also consistent with locations highlighted in non-  
699 invasive human fMRI studies (e.g., (Grill-Spector and Malach 2004)) and with putative  
700 homologous regions in the macaque brain (e.g., (Logothetis and Sheinberg 1996; Tanaka 1996;  
701 Connor et al. 2007)).

702  
703 Once the onset of visual changes is detected, approximately the same ventral visual regions  
704 provide a rich representation that contains selective information about the nature of those  
705 changes (**Figures 5, 7C-D**). Selective visual information arose within the first 200 ms of a movie  
706 cut, and was relatively robust to the many highly varied transformations that took place in these  
707 commercial movies. Specifically, in Experiment I, classifiers trained to detect the presence  
708 versus absence of a humanized animal, using electrodes in the inferior occipital gyrus, inferior  
709 temporal gyrus or fusiform gyrus, showed a significant degree of extrapolation to independent  
710 test data from a completely different movie clip (**Figure 5C**, green bars). In Experiment II,  
711 classifiers trained to discriminate the presence versus absence of human faces from the field  
712 potential responses from single electrodes in the inferior occipital gyrus or fusiform gyrus  
713 showed a weak but significant degree of extrapolation to independent test data during single  
714 repetitions of other parts of the movie (**Figure 7C-D**). The results in **Figure 5** should *not* be  
715 interpreted to imply that those electrodes were selective to “humanized animals” or that the  
716 corresponding analyses in **Figure 7** imply selectivity for “human faces”. This study used  
717 commercial movies and no attempt was made to circumscribe the visual changes to the  
718 appearance of animals or faces. The appearance of animals and faces was correlated and  
719 accompanied by changes in motion, contrast and many other visual properties. It seems likely  
720 that the main drivers of the strong visually evoked transitions, such as the ones illustrated in  
721 **Figure 2**, are the sharp contrast changes and motion energy changes triggered by movie cuts.  
722 Further studies directly comparing the responses to dynamic stimuli versus stimulus flashes will  
723 be needed to further dissect the specific features that dictate selectivity to movie events  
724 revealed here. The current results demonstrate that it is possible to distill reliable, selective and  
725 invariant information, even in single events during a continuous stream of frames.

726  
727 Moving from repeated presentations of identical, static stimuli with fixed onsets and offsets to  
728 movie stimuli constitutes an important step to bridge the gap between laboratory studies and  
729 understanding vision in the real world. Furthermore, movies present rich visual and social input.  
730 The initial methodological steps suggested here open the doors to interpreting neural responses  
731 to complex cognitive events during single presentations of movies.

732  
733 **Acknowledgements**

734 We thank all of the patients for participating in this study. This paper is based upon work  
735 supported by NIH and the Center for Brains Minds and Machines (CBMM), funded by NSF STC  
736 award CCF-1231216.

737

738

739

ACCEPTED MANUSCRIPT

740 **References**

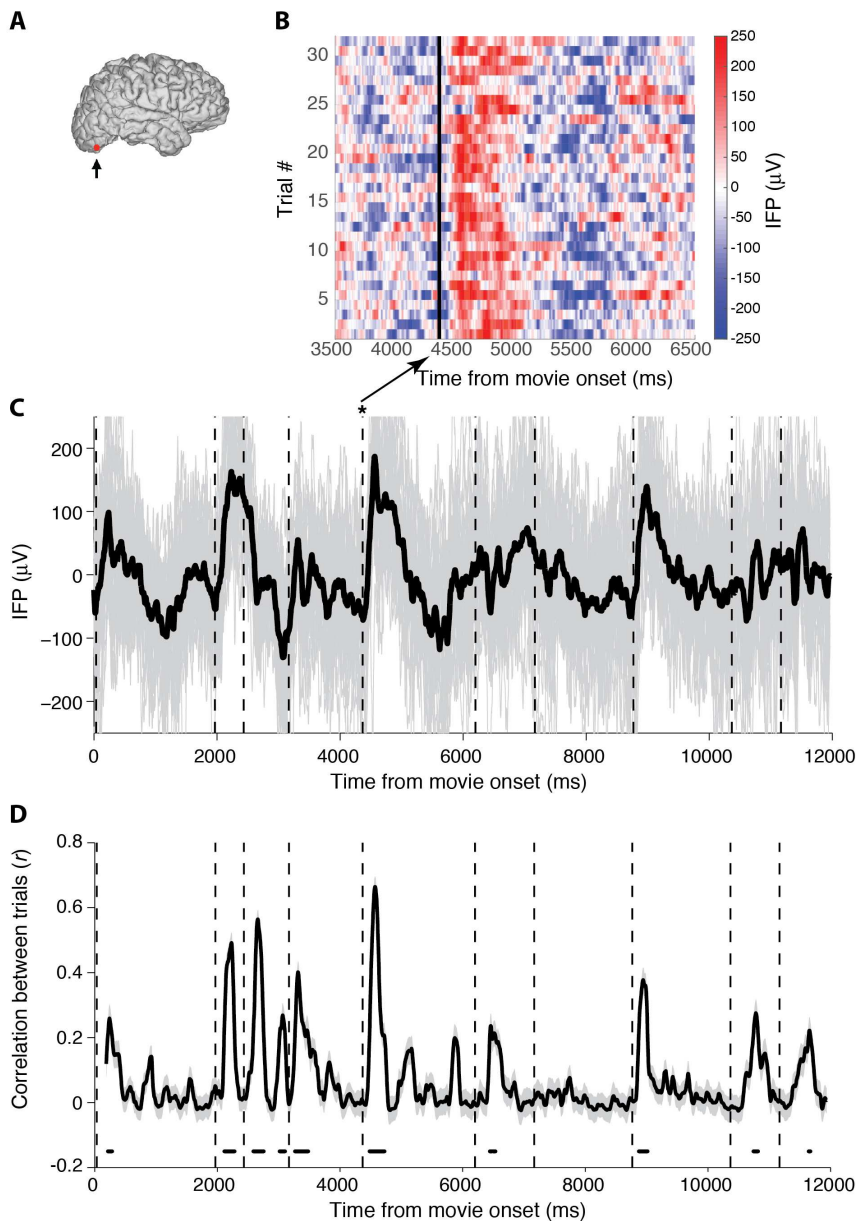
- 741 Bansal AK, Singer JM, Anderson WS, Golby A, Madsen JR, Kreiman G. 2012. Temporal  
742 stability of visually selective responses in intracranial field potentials recorded from human  
743 occipital and temporal lobes. *J Neurophysiol.* 108:3073–3086.
- 744 Bartels A, Zeki S. 2005. Brain dynamics during natural viewing conditions—A new guide for  
745 mapping connectivity in vivo. *Neuroimage.* 24:339–349.
- 746 Booth MC, Rolls ET. 1998. View-invariant representations of familiar objects by neurons in the  
747 inferior temporal visual cortex. *Cereb Cortex.* 8:510–523.
- 748 Buzsáki G, Anastassiou CA, Koch C. 2012. The origin of extracellular fields and currents —  
749 EEG, ECoG, LFP and spikes. *Nat Rev Neurosci.* 13:407–420.
- 750 Carandini M, Demb JB, Mante V, Tolhurst DJ, Dan Y, Olshausen BA, Gallant JL, Rust NC.  
751 2005. Do We Know What the Early Visual System Does? *J Neurosci.* 25.
- 752 Connor CE, Brincat SL, Pasupathy A. 2007. Transformation of shape information in the ventral  
753 pathway. *Curr Opin Neurobiol.* 17:140–147.
- 754 Conroy BR, Singer BD, Guntupalli JS, Ramadge PJ, Haxby J V. 2013. Inter-subject alignment of  
755 human cortical anatomy using functional connectivity. *Neuroimage.* 81:400–411.
- 756 Destrieux C, Fischl B, Dale A, Halgren E. 2010. Automatic parcellation of human cortical gyri  
757 and sulci using standard anatomical nomenclature. *Neuroimage.* 53:1–15.
- 758 DiCarlo JJ, Zoccolan D, Rust NC. 2012. How does the brain solve visual object recognition?  
759 *Neuron.* 73:415–434.
- 760 Dudai Y. 2012. The cinema-cognition dialogue: a match made in brain. *Front Hum Neurosci.*  
761 6:248.
- 762 Felsen G, Dan Y. 2005. A natural approach to studying vision. *Nat Neurosci.* 8:1643–1646.
- 763 Fiser J, Chiu C, Weliky M. 2004. Small modulation of ongoing cortical dynamics by sensory  
764 input during natural vision. *Nature.* 431:573–578.
- 765 Grill-Spector K, Malach R. 2004. The human visual cortex. *Annu Rev Neurosci.* 27:649–677.
- 766 Hasson U, Nir Y, Levy I, Fuhrmann G, Malach R. 2004. Intersubject synchronization of cortical  
767 activity during natural vision. *Science.* 303:1634–1640.
- 768 Honey CJ, Thesen T, Donner TH, Silbert LJ, Carlson CE, Devinsky O, Doyle WK, Rubin N,  
769 Heeger DJ, Hasson U. 2012. Slow cortical dynamics and the accumulation of information  
770 over long timescales. *Neuron.* 76:423–434.
- 771 Hung CP, Kreiman G, Poggio T, DiCarlo JJ. 2005. Fast readout of object identity from macaque  
772 inferior temporal cortex. *Science.* 310:863–866.
- 773 Huth AG, Nishimoto S, Vu AT, Gallant JL. 2012. A continuous semantic space describes the  
774 representation of thousands of object and action categories across the human brain.  
775 *Neuron.* 76:1210–1224.
- 776 Isik L, Meyers EM, Leibo JZ, Poggio T. 2014. The dynamics of invariant object recognition in the  
777 human visual system. *J Neurophysiol.* 111:91–102.
- 778 Keysers C, Xiao D-K, Földiák P, Perrett DI. 2001. The Speed of Sight. *J Cogn Neurosci.* 13:90–  
779 101.
- 780 Kriegeskorte N, Kreiman G. 2011. *Visual Population Codes: Toward a Common Multivariate  
781 Framework for Cell ...* - Google Books.
- 782 Lei Y, Sun N, Wilson FAW, Wang X, Chen N, Yang J, Peng Y, Wang J, Tian S, Wang M, Miao  
783 Y, Zhu W, Qi H, Ma Y. 2004. Telemetric recordings of single neuron activity and visual  
784 scenes in monkeys walking in an open field. *J Neurosci Methods.* 135:35–41.
- 785 Lewen GD, Bialek W, Steveninck RR d. R v. 2001. Neural coding of naturalistic motion stimuli.  
786 *Netw Comput Neural Syst.* 12:317–329.
- 787 Liu H, Agam Y, Madsen JR, Kreiman G. 2009. Timing, timing, timing: fast decoding of object  
788 information from intracranial field potentials in human visual cortex. *Neuron.* 62:281–290.
- 789 Logothetis NK, Sheinberg DL. 1996. Visual object recognition. *Annu Rev Neurosci.* 19:577–621.
- 790 McMahon DBT, Russ BE, Elnaïem HD, Kurnikova AI, Leopold DA. 2015. Single-Unit Activity

- 791 during Natural Vision: Diversity, Consistency, and Spatial Sensitivity among AF Face Patch  
792 Neurons. *J Neurosci.* 35.
- 793 Meyers EM, Freedman DJ, Kreiman G, Miller EK, Poggio T. 2008. Dynamic population coding  
794 of category information in inferior temporal and prefrontal cortex. *J Neurophysiol.*  
795 100:1407–1419.
- 796 Miller KJ, Honey CJ, Hermes D, Rao RP, denNijs M, Ojemann JG. 2014. Broadband changes in  
797 the cortical surface potential track activation of functionally diverse neuronal populations.  
798 *Neuroimage.* 85:711–720.
- 799 Montemurro MA, Rasch MJ, Murayama Y, Logothetis NK, Panzeri S. 2008. Phase-of-Firing  
800 Coding of Natural Visual Stimuli in Primary Visual Cortex, *Current Biology.*
- 801 Nishimoto S, Vu AT, Naselaris T, Benjamini Y, Yu B, Gallant JL. 2011. Reconstructing visual  
802 experiences from brain activity evoked by natural movies. *Curr Biol.* 21:1641–1646.
- 803 Ojemann GA. 1997. Treatment of Temporal Lobe Epilepsy. *Annu Rev Med.* 48:317–328.
- 804 Podvalny E, Yeagle E, Mégevand P, Sarid N, Harel M, Chechik G, Mehta AD, Malach R. 2016.  
805 Invariant Temporal Dynamics Underlie Perceptual Stability in Human Visual Cortex. *Curr*  
806 *Biol.*
- 807 Privman E, Fisch L, Neufeld MY, Kramer U, Kipervasser S, Andelman F, Yeshurun Y, Fried I,  
808 Malach R. 2011. Antagonistic Relationship between Gamma Power and Visual Evoked  
809 Potentials Revealed in Human Visual Cortex. *Cereb Cortex.* 21:616–624.
- 810 Privman E, Nir Y, Kramer U, Kipervasser S, Andelman F, Neufeld MY, Mukamel R, Yeshurun Y,  
811 Fried I, Malach R. 2007. Enhanced Category Tuning Revealed by Intracranial  
812 Electroencephalograms in High-Order Human Visual Areas. *J Neurosci.* 27.
- 813 Richmond BJ, Optican LM, Spitzer H. 1990. Temporal encoding of two-dimensional patterns by  
814 single units in primate primary visual cortex. I. Stimulus-response relations. *J Neurophysiol.*  
815 64.
- 816 Riesenhuber M, Poggio T. 1999. Hierarchical models of object recognition in cortex. *Nat*  
817 *Neurosci.* 2:1019–1025.
- 818 Rolls ET, Tovee MJ. 1995. Sparseness of the neuronal representation of stimuli in the primate  
819 temporal visual cortex. *J Neurophysiol.* 73.
- 820 Russ BE, Leopold DA. 2015. Functional MRI mapping of dynamic visual features during natural  
821 viewing in the macaque. *Neuroimage.* 109:84–94.
- 822 Rust NC, Movshon JA. 2005. In praise of artifice. *Nat Neurosci.* 8:1647–1650.
- 823 Serre T, Oliva A, Poggio T. 2007. A feedforward architecture accounts for rapid categorization.  
824 *Proc Natl Acad Sci U S A.* 104:6424–6429.
- 825 Smith TJ, Levin D, Cutting JE. 2012. A Window on Reality. *Curr Dir Psychol Sci.* 21:107–113.
- 826 Tanaka K. 1996. Inferotemporal cortex and object vision. *Annu Rev Neurosci.* 19:109–139.
- 827 Tang H, Buia C, Madhavan R, Crone NE, Madsen JR, Anderson WS, Kreiman G. 2014.  
828 Spatiotemporal dynamics underlying object completion in human ventral visual cortex.  
829 *Neuron.* 83:736–748.
- 830 Vidal JR, Ossandón T, Jerbi K, Dalal SS, Minotti L, Ryvlin P, Kahane P, Lachaux J-P. 2010.  
831 Category-Specific Visual Responses: An Intracranial Study Comparing Gamma, Beta,  
832 Alpha, and ERP Response Selectivity. *Front Hum Neurosci.* 4:195.
- 833 Vinje WE, Gallant JL. 2000. Sparse Coding and Decorrelation in Primary Visual Cortex During  
834 Natural Vision. *Science (80- ).* 287.
- 835 Whittingstall K, Bartels A, Singh V, Kwon S, Logothetis NK. 2010. Integration of EEG source  
836 imaging and fMRI during continuous viewing of natural movies. *Magn Reson Imaging.*  
837 28:1135–1142.
- 838

**Figures****Figure 1 - Experimental paradigm and movie cuts**

**A.** Experiment I - Three 12s clips from commercial cartoon movies were presented multiple times without sound, at 30 frames per second, and subtending  $\sim 4 \times 3$  degrees of visual angle (see Supplemental table 1). The first frame, three middle frames (demonstrating a movie cut between frames 130-131), and the final frame from clip 1 are shown (Methods). Subjects passively viewed the 12s movie clips.

**B.** Experiment II – A full-length movie was shown once through with sound, at 24 frames per second, and subtending  $\sim 18 \times 12$  degrees of visual angle. We considered data from patients watching one of three movies in this study (Methods). Example frames from one of the movies, Home Alone 2, including the first frame, three middle frames (demonstrating a movie cut between frames 15869-15870), and the final frame are shown. Subjects passively watched these full-length movies.



**Figure 2 – Example electrode showing consistent physiological responses to movie cuts (Experiment I)**

**A.** Electrode location. The electrode was located in the right inferior occipital gyrus (Talairach coordinates = [35.9, -82.8, -14.5]).

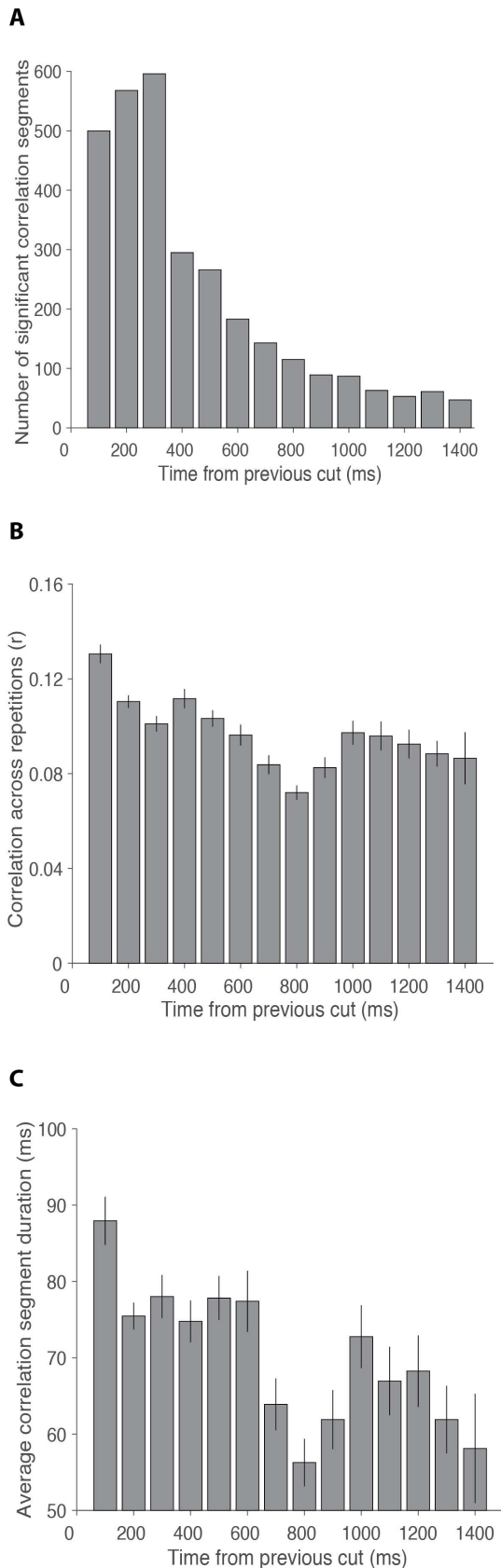
**B.** Raster plot showing the intracranial field potential (IFP) surrounding the cut transition shown in **Figure 1A** (frame 130-131 in movie clip 1). Each row denotes a repetition of the movie ( $n=32$  repetitions). The color indicates the IFP at each time point (bin size = 3.9 ms, see color scale on right). The movie cut triggered a large change in voltage in almost every repetition.

**C.** Electrode's broadband voltage time course over the entire 12s movie clip 1. Mean activity is shown with a thick black line, and 32 individual repetitions are shown with gray traces. Dashed vertical lines indicate movie cuts, and the cut shown in **B** is indicated with an asterisk.

Several, but not all, of the cut transitions elicited large voltage changes that can be observed even in individual repetitions. The y-axis is cut at  $-250$  and  $250 \mu\text{V}$  but some individual traces extend beyond these limits.

**D.** Average pairwise correlation (Pearson coefficient,  $r$ ,  $\text{mean} \pm \text{SEM}$ ) across the 32 choose 2 (496) pairwise comparisons between repetitions calculated in 50 ms non-overlapping bins. Horizontal black lines at the bottom of the plot indicate time periods when the average pairwise correlation across repetitions was significantly above chance based on a  $p < 0.01$  permutation test (Methods). See **Figure S2** for a similar example in the high gamma frequency band.





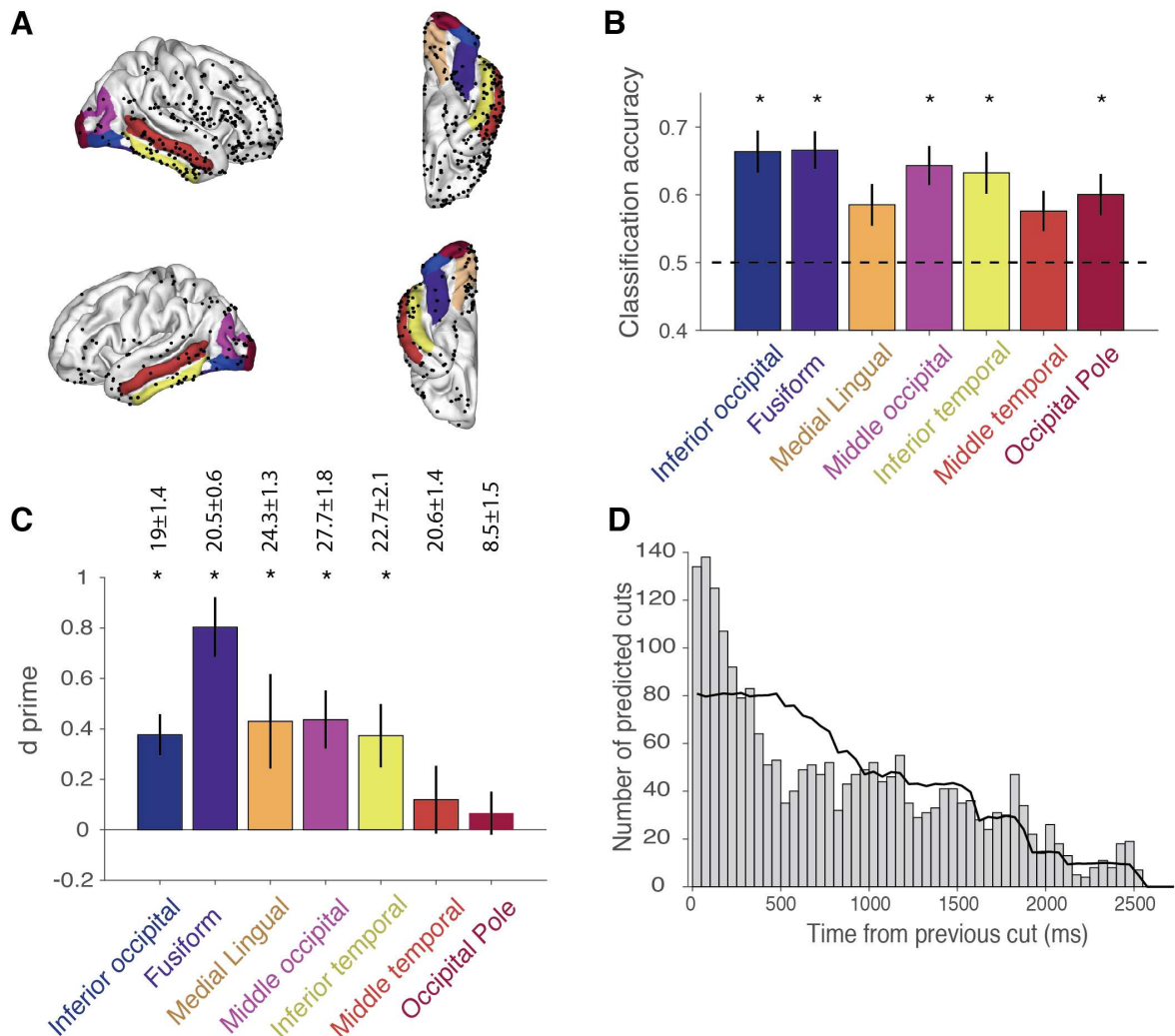
**Figure 3. Properties of neural responses that were consistent across trials**

**A.** Distribution of the onset of segments with statistically significant correlation across repetitions in all electrodes ( $n=954$ ), calculated with a sliding window of 50 ms duration, as a function of time from the previous cut. Bin size = 100 ms (**Methods**). These segments of consistent correlation across repetitions begin mostly within the 300 ms following a cut.

**B.** Average correlation coefficient between repetitions in each time bin for all the segments with statistically significant correlation between repetitions in **A** (mean $\pm$ SEM).

**C.** Average duration between the beginning of the first and last time points for all the consecutive segments with statistically significant correlation between repetitions in **A** (mean $\pm$ SEM).

See **Figure S3** for corresponding analyses in different frequency bands and **Figure S13** for the same analyses using different window sizes.



#### Figure 4 – Movie cuts and shots can be decoded from ventral visual cortex regions

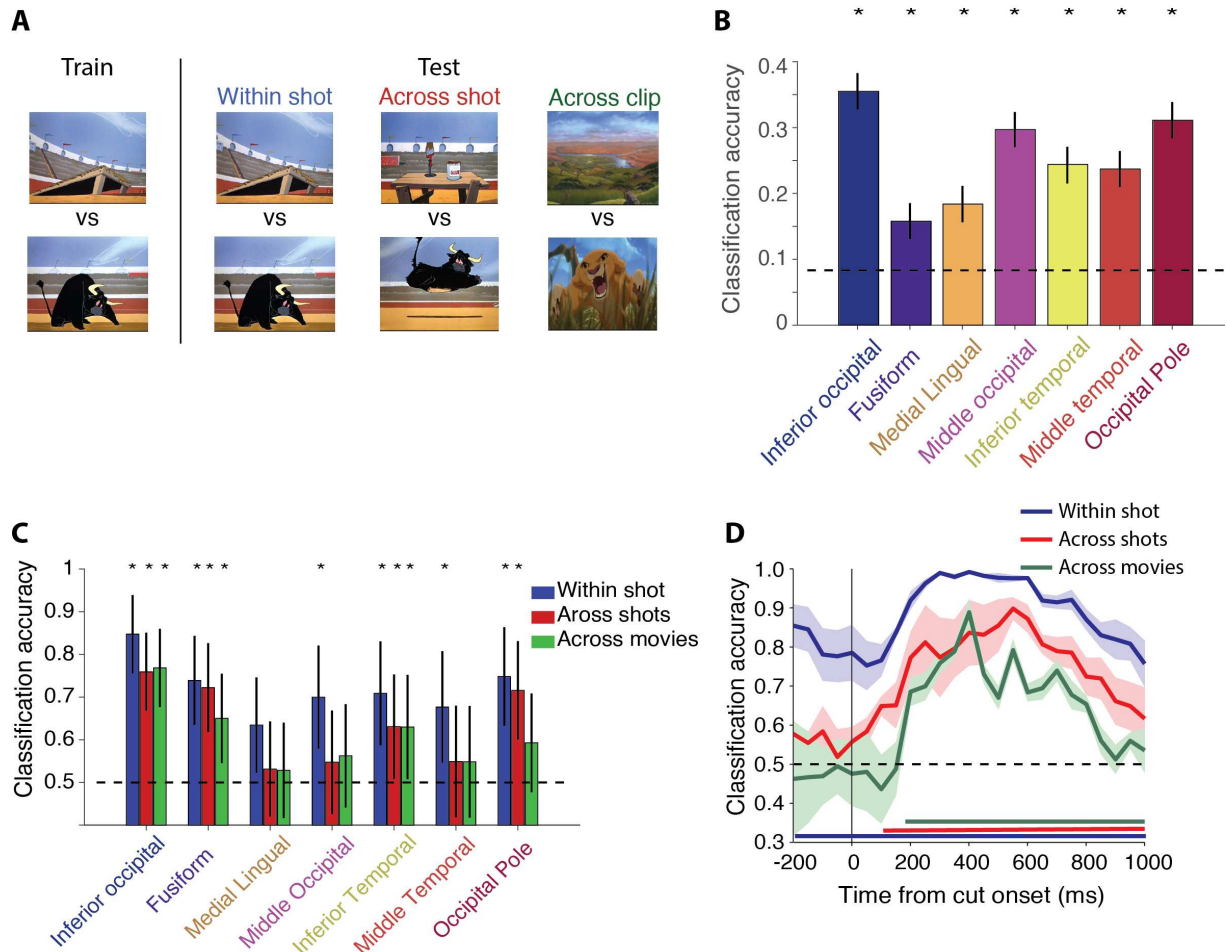
**A.** Location of all electrodes in Experiment I projected onto a common reference brain (Freesurfer fsaverage brain) shown at lateral and ventral views. Each dot corresponds to one electrode (total = 994 electrodes, Supplemental Table 1). Seven anatomical regions (out of 25 regions with at least eight electrodes) with significantly above chance decoding performance in any of the decoding tasks in **B** or **C** are highlighted.

**B.** Classification accuracy from  $n=8$  electrodes in each region, between movie segments with a cut versus segments without a cut in the seven regions highlighted in **Fig. 4A** (mean  $\pm$  SD across 20 decoding runs, Methods). Chance = 0.5. The classification accuracy is reported as the average from 50-400 ms post cut onset. Asterisks indicate significant decoding based on a  $p<0.01$  permutation test (Methods, **Supplemental Figure 12A**).

**C.** Sensitivity ( $d'$ ) to detect visual transitions during the entire 12s clip time course for held out repetitions of movie clips 1 and 2 (mean  $\pm$  SD across 20 decoding runs, **Methods**, **Supplemental Figure 12B**). Number at the top of each bar plot indicates the number of predicted visual transitions per region (the actual number of all cuts in movie clips 1 and 2 was 17).

**D.** The bars show the latency difference between the time of the predicted visual transitions (first time point in visual transition predicted periods, see **Figures S12B**) in **C** and the time of the

previous true cut for the five regions with significantly above chance  $d'$  values in **C**. Bin size = 50 ms. The line shows the average distribution obtained from randomly selecting the same number of times as predicted visual transitions. The distribution of selected transition times is significantly different from the random distribution ( $p < 10^{-10}$ , Kolmogorov-Smirnov test). See **Figure S7** for corresponding analyses in different frequency bands.



**Figure 5 – Visual information generalizes across movies in 12s clips**

**A.** We decoded shots with an animal versus shots without an animal, first from repetitions of the exact same shot pairs (“within shot”, blue), next with generalization across different shots in the same movie clip (“across shot”, red), and finally across movie clips (“across clips”, green). One example pair of frames (first frame in shot) depicting the different conditions is shown. Decoding was repeated for four pairs of clips (from two of the three 12s clips that represent the two unique movies, Fig. S4, Methods).

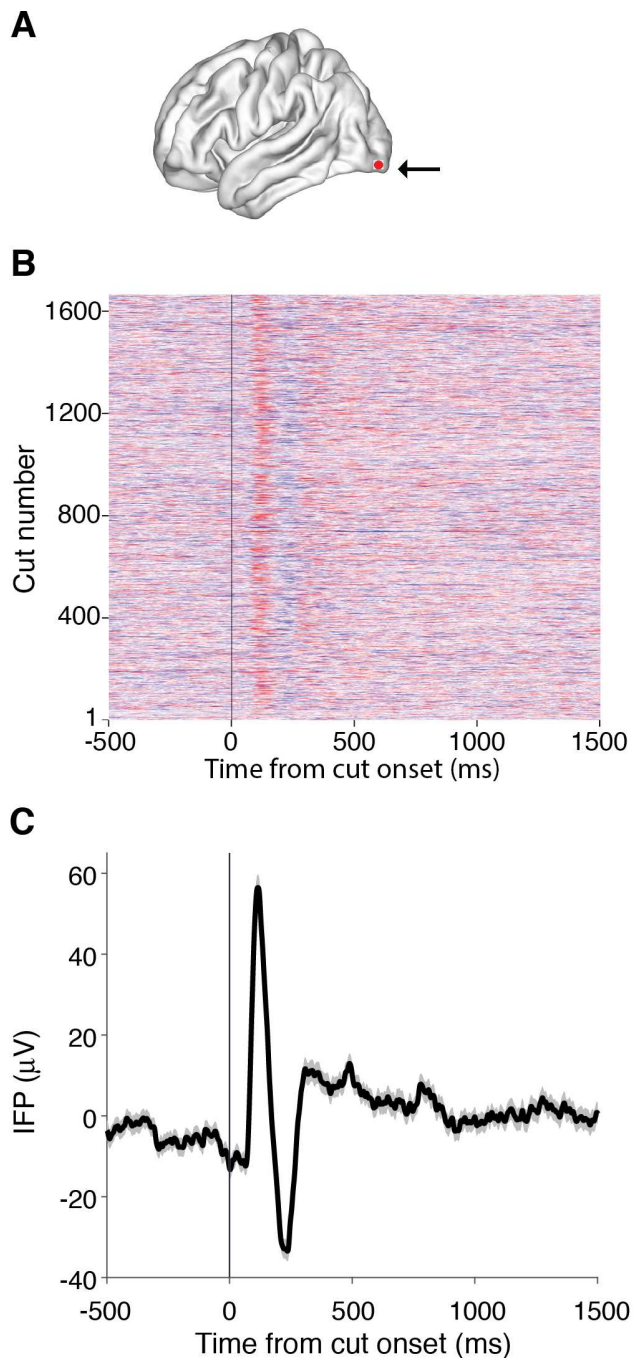
**B.** Classification accuracy to label each of the 13 cuts from clips 1 and 2 (excluding the first and last cut from each movie, **Figure S6**) using  $n=8$  electrodes in each of the seven regions highlighted in **A** (mean  $\pm$  SD across 20 decoding runs, Methods). Chance =  $1/13$ . The classification accuracy is reported as the average from 50-400 ms post cut onset. Asterisks indicate significant decoding based on a  $p<0.01$  permutation test (Methods).

See **Figure S7** for corresponding analyses in different frequency bands.

**C.** Classification accuracy from  $n=8$  electrodes in each region for shots with versus without an animal (mean  $\pm$  SD across 20 decoding runs, chance = 0.5) in the seven highlighted regions described in **Fig. 4A**. We considered 3 conditions corresponding to different levels of extrapolation: within shot (blue), across shots (red), and across movies (green). The classification accuracy is reported as the average from 50-400 ms post cut onset. Region labels are color coded following the conventions in **Fig. 4A**. Asterisks indicate significant decoding for each of the three decoding conditions based on  $p<0.01$  permutation test (Methods).

**D.** Visualization of dynamic classification accuracy for shots with an animal versus without an animal across time relative to cut onset from a pseudo population based on feature selection

from all electrodes across all subjects (mean  $\pm$  SD across 20 decoding runs, Methods). Feature selection was applied at each time point to choose selective electrodes in the training data to be used in the classifier (Methods). Horizontal line indicates chance classification. Note that the 'within shot' classification accuracy was significantly above chance even before the cut onset, because the visual stimulus pre-cut was identical in the training and test sets (see discussion in text). While the 'Within shot' classification accuracy was significantly above chance for the entire time course, the 'Across shot' and 'Across clip' classification accuracies were significantly above chance from 100-1000 ms and 200-1000 ms post-cut onset, respectively. See **Figure S9** for corresponding analyses in different frequency bands.



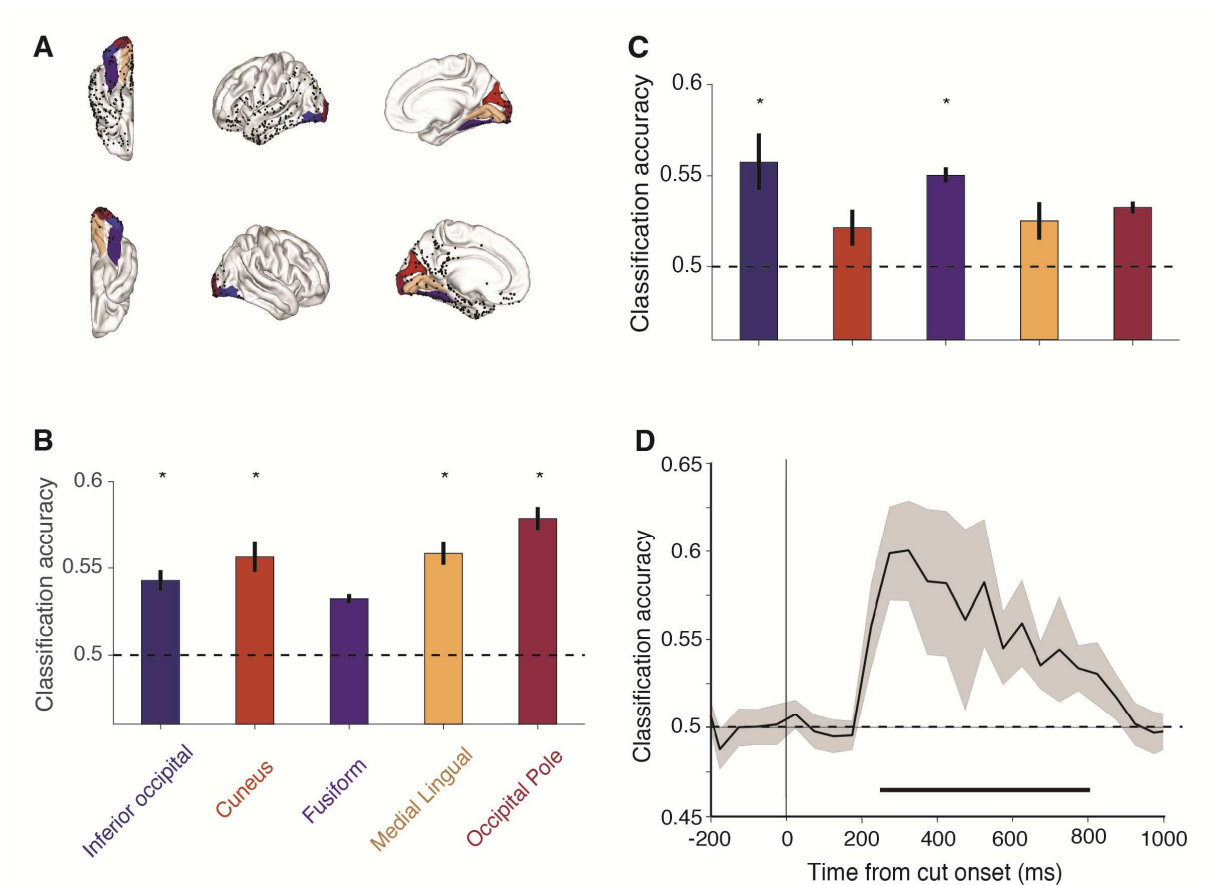
**Figure 6 – Example electrode showing a consistent physiological response to movie cuts in full-length movies**

**A.** Location of one example cut-responsive electrode (Experiment II) in the left occipital pole (Talairach coordinates = [-2.2, -92.4, -4.3])

**B.** Raster plot showing the intracranial field potential (IFP) surrounding all cut transitions in the full-length movie (Home Alone 2). Each row denotes a different cut ( $n=1630$  cuts). The color indicates the IFP at each time point (bin size = 0.5 ms, see color scale on right).

**C.** Average IFP time course (mean  $\pm$  SEM) over all movie cuts.

See Figure S11 for a similar example in the high gamma frequency band.



**Figure 7 - Movie cuts can be decoded from a single presentation of a full-length movie**

**A.** Location of all electrodes in Experiment II projected onto a common reference brain (Freesurfer fsaverage brain) shown at lateral and ventral views. Each dot corresponds to one electrode (total = 330 electrodes, **Supplemental Table 1**). The five anatomical regions (out of 20 regions with at least five electrodes) with significantly above chance classification accuracy in any of the decoding tasks in **Fig. 7** or **Fig. 8** are highlighted.

**B.** Average single electrode classification accuracy between movie segments with a cut versus those without a cut for the 5 regions highlighted in **Fig. 7A** (mean  $\pm$  SEM across all electrodes in each region). Chance = 0.5 (horizontal dashed line). The classification accuracy is calculated as the average from 50-400 ms post cut onset. The number of electrodes averaged is: inferior occipital, n=7; cuneus, n=10; fusiform, n=26, medial lingual, n=11, occipital pole, n=13. Asterisks indicate regions with significantly above chance average classification accuracy based on a  $p < 0.01$  permutation test (Methods).

**C.** Average single electrode classification accuracy between movie shots with a face versus those without a face for those regions highlighted in **Fig. 7A** (mean  $\pm$  SEM across all electrodes in each region). Chance = 0.5, horizontal dashed line. The classification accuracy is calculated as the average from 50-400 ms post cut onset. The number of electrodes averaged is the same as in **Fig. 7B**. Asterisks indicate regions with significantly above chance average classification accuracy based on a  $p < 0.01$  permutation test.

**D.** Visualization of dynamic classification accuracy for shots with a face versus those without a face versus time relative to cut onset using feature selection from all subjects and all electrodes (mean  $\pm$  SEM across four subjects, Methods). Feature selection across all electrodes based on the training data only was applied at each time point to choose selective electrodes to be used.

in the classifier (Methods). Since the subjects viewed different movies, decoding results were then averaged post-hoc. Horizontal line indicates chance classification. *The decoding was significantly above chance from 250-850 ms post-cut onset based on a  $p < 0.01$  permutation test.*

See **Figure S11** for corresponding analyses in different frequency bands.

ACCEPTED MANUSCRIPT