

1 **TDM: a temporal decomposition method for**
2 **removing venous effects from task-based fMRI**

3
4 Kendrick Kay^{*,1}, Keith W. Jamison^{1,#}, Ruyuan Zhang¹, Kamil Ugurbil¹

5
6 ¹Center for Magnetic Resonance Research (CMRR), Department of Radiology, University of Minnesota

7 [#]*Current address:* Department of Radiology, Weill Cornell Medical College

8 ^{*}Corresponding author (kay@umn.edu)

9
10 *Running title:* Temporal decomposition method for task-based fMRI

11
12 *Keywords:* high-resolution fMRI, veins, vasculature, hemodynamic timecourse, hemodynamic response
13 function, primary visual cortex, eccentricity, 7T

1 **Abstract**

2

3 Most functional magnetic resonance imaging (fMRI) is conducted with gradient-echo pulse sequences.
4 Although this yields high sensitivity to blood oxygenation level dependent (BOLD) signals, gradient-echo
5 acquisitions are heavily influenced by venous effects which limit the ultimate spatial resolution and spatial
6 accuracy of fMRI. While alternative acquisition methods such as spin-echo can be used to mitigate
7 venous effects, these methods lead to serious reductions in signal-to-noise ratio and spatial coverage,
8 and are difficult to implement without leakage of undesirable non-spin-echo effects into the data.
9 Moreover, analysis heuristics such as masking veins or sampling inner cortical depths using high-
10 resolution fMRI may be helpful, but sacrifice information from many parts of the brain. Here, we describe
11 a new analysis method that is compatible with conventional gradient-echo acquisition and provides
12 venous-free response estimates throughout the entire imaged volume. The method involves fitting a low-
13 dimensional manifold characterizing variation in response timecourses observed in a given dataset, and
14 then using identified early and late timecourses as basis functions for decomposing responses into
15 components related to the microvasculature (capillaries and small venules) and the macrovasculature
16 (veins), respectively. We show that this *Temporal Decomposition through Manifold Fitting (TDM)* method
17 is robust, consistently deriving meaningful timecourses in individual fMRI scan sessions. Moreover, we
18 show that by removing late components, TDM substantially reduces the superficial cortical depth bias
19 present in gradient-echo BOLD responses and eliminates artifacts in cortical activity maps. TDM is
20 general: it can be applied to any task-based fMRI experiment, can be used with standard- or high-
21 resolution fMRI acquisitions, and can even be used to remove residual venous effects from specialized
22 acquisition methods like spin-echo. We suggest that TDM is a powerful method that improves the spatial
23 accuracy of fMRI and provides insight into the origins of the BOLD signal.

24

1. Introduction

Among the handful of noninvasive techniques that permit the study of human brain activity, functional magnetic resonance imaging (fMRI) based on blood oxygenation level dependent (BOLD) contrast has emerged as by far the most commonly used approach and has proven to be an extremely useful tool in cognitive neuroscience. A primary advantage of fMRI over other measurement techniques that report on neural activity (e.g., EEG, MEG, PET) is its spatial resolution. Whole-brain volumes sensitized to the BOLD effect are routinely acquired at a resolution of 2–3 mm, and even at sub-second imaging times as exemplified by data from the Human Connectome Project (Ugurbil et al., 2013). However, efforts to increase the spatial resolution of fMRI further—especially to reach the sub-millimeter scale of mesoscopic brain organization—face a major challenge imposed by the indirect nature of fMRI signals. The BOLD response is mediated by neurovascular coupling (Attwell and Iadecola, 2002) that links neuronal activity to secondary hemodynamic and metabolic alterations. An undesirable consequence of this mechanism is the “draining vein” confound that affects the most commonly employed gradient-echo BOLD fMRI technique, first noted early in the history of fMRI (Menon et al., 1993) and investigated in numerous subsequent studies. Deoxyhemoglobin changes originating at the site of neuronal activity subsequently propagate through venules and veins; these effects may appear as activation displaced from the original site of neural activity by as much as 4 mm (Turner, 2002) and may reflect neural activity pooled over large spatial scales, thus degrading spatial specificity with respect to underlying neural activity patterns (Bianciardi et al., 2011; Kay et al., 2019; Olman et al., 2007; Shmuel et al., 2007). In response to this problem, the field has long sought to measure BOLD responses from the microvasculature (e.g. capillaries and venules) while avoiding BOLD responses from the macrovasculature (e.g. veins) (Cheng, 2018; Ugurbil, 2016; Yacoub and Wald, 2018). The problem of the macrovasculature is especially critical to resolve given the growing popularity in the neuroscience community of using fMRI to probe laminar-specific responses (see De Martino et al., 2018; Dumoulin et al., 2018; Lawrence et al., 2017 for reviews).

In order to avoid the specificity loss caused by veins, the field has traditionally turned to the use of spin-echo acquisition (e.g. Parkes et al., 2005; Yacoub et al., 2008) instead of conventional gradient-echo acquisition. However, spin-echo involves increased energy deposition, longer volume acquisition times, and lower signal-to-noise ratio. Thus, in order to maintain measurement sensitivity, the experimenter is generally forced to reduce spatial coverage and/or substantially increase the amount of data collected per experimental condition. These unappealing options are often dealbreakers for neuroscientists, given that measuring multiple brain regions is often critical, the sensitivity of fMRI is already relatively low to start with, and increasing the duration of an experiment beyond more than a factor of two or so is often impractical.

Here, we introduce an analysis method, called *Temporal Decomposition through Manifold Fitting (TDM)*, that identifies and removes venous-related signals from task-based fMRI data. The TDM method is simple, principled, and is compatible with a variety of experimental protocols including those based on gradient-echo acquisitions. We demonstrate TDM on visual experiments conducted on human subjects, and show that TDM consistently removes unwanted venous effects while maintaining a reasonable level of sensitivity. We end with a discussion of the strengths and limitations of the TDM method and a comparison of TDM to other approaches for reducing venous effects in fMRI. Code implementing TDM is freely available at <https://osf.io/j2wsc/>.

2. Materials and Methods

2.1. Subjects

Eleven subjects (five males, six females; one subject, S1, was an author (K.K.)) participated in the experiments described in this study. All subjects had normal or corrected-to-normal visual acuity. Informed written consent was obtained from all subjects, and the experimental protocol was approved by the University of Minnesota Institutional Review Board.

We conducted four experiments. Experiment E1 measured responses to eccentricity stimuli using a high-resolution (7T, 0.8 mm) gradient-echo protocol. Experiment E2 measured responses to category stimuli also using the high-resolution gradient-echo protocol; data from this experiment are the same as described in a previous publication (Kay et al., 2019). Experiment E3 measured responses to the same eccentricity stimuli in E1 but used a spin-echo protocol (7T, 1.05 mm). Experiment E4 measured responses to the same eccentricity stimuli in E1 but used a low-resolution (3T, 2.4 mm) gradient-echo protocol.

A total of sixteen datasets (scan sessions) were collected: five corresponding to Experiment E1; seven corresponding to Experiment E2; two corresponding to Experiment E3; and two corresponding to Experiment E4. To facilitate direct comparison, Experiments E3 and E4 were conducted in subjects who also participated in Experiment E1. A full breakdown of subjects, experiments, and datasets is provided in **Figure 1**.

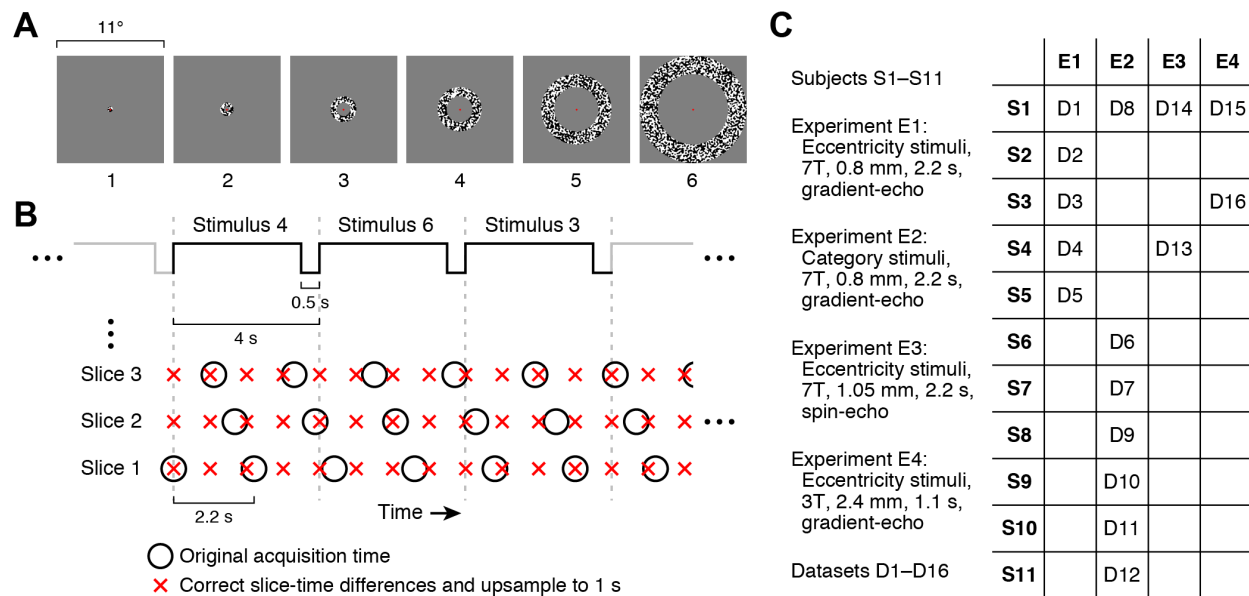


Figure 1. Schematic of experiment. A, Eccentricity stimuli. Stimuli consist of six rings varying in eccentricity. A small central dot serves as a fixation point. B, Timing of trials and fMRI acquisition. Each stimulus lasts 3.5 s and is followed by a 0.5-s gap before the next trial. Slices are acquired at different times within each 2.2-s TR (black circles). In pre-processing, cubic interpolation is used to resample each voxel's time-series data to a rate of 1 s such that the same time points are obtained for all voxels (red crosses). Because the trial duration (4 s) is not evenly divisible by the TR (2.2 s), the experiment automatically incorporates jitter between trial onsets and slice acquisition times. C, Summary of datasets. Eleven subjects participated in four experiments. Each dataset corresponds to one scan session, and sixteen datasets were collected. (For details on the category stimuli used in Experiment E2, see Section 2.3.)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

2.2. Stimulus presentation

For the 7T datasets, stimuli were presented using a Cambridge Research Systems BOLDscreen 32 LCD monitor positioned at the head of the scanner bed (resolution 1920 × 1080 at 120 Hz; viewing distance 189.5 cm). For the 3T datasets, stimuli were presented using a NEC NP4100 DLP projector that was focused onto a backprojection screen positioned at the head of the scanner bore (resolution 1024 × 768 at 60 Hz; viewing distance 102 cm). Subjects viewed the monitor or backprojection screen via a mirror mounted on the RF coil. A Mac Pro (7T) or iMac (3T) computer controlled stimulus presentation using code based on Psychophysics Toolbox (Brainard, 1997; Pelli, 1997). Behavioral responses were recorded using a button box.

2.3. Experimental design

In the eccentricity experiment (Experiments E1, E3, E4), stimuli consisted of rings positioned at six different eccentricities (i.e. distances from the center of gaze), and were confined to a circular region with diameter 11°. Each ring was filled with a black-and-white contrast pattern that updated at 10 Hz. Ring size scaled with eccentricity, and rings were presented on a neutral gray background (**Figure 1A**). Stimuli were presented in 4-s trials. In a trial, one of the six rings was presented for 3.5 s (35 images presented sequentially, each with duration 0.1 s) and was followed by a brief gap of 0.5 s. Each run lasted 368.116 s and included 12 presentations of each of the 6 rings as well as blank trials (also of 4-s duration). Throughout stimulus presentation, a small semi-transparent dot (50% opacity) was present at the center of the stimulus. The color of the dot switched between red, white, and black every 1–5 s, and subjects were instructed to maintain fixation on the dot and to press a button whenever the color changed. A total of 9 runs were collected in each 7T scan session, and a total of 6 runs were collected in each 3T scan session. (The eccentricity experiment was time-locked to the refresh rate of the LCD monitor, which caused the additional 116 ms in the total run duration. To compensate for this slight offset, we pre-processed the fMRI data for the eccentricity experiment at a sampling rate of 1.000316 s, and then, for simplicity, treated the data in subsequent analyses as if the sampling rate was exactly 1.0 s.)

The category experiment (Experiment E2) was the same as the ‘functional localizer’ experiment conducted in a previous paper (Kay et al., 2019). This experiment (<http://vpnl.stanford.edu/fLoc/>) was developed by the Grill-Spector lab (Stigliani et al., 2015). Stimuli consisted of grayscale images of different semantically meaningful categories. There were 10 categories, grouped into 5 stimulus domains: characters (word, number), body parts (body, limb), faces (adult, child), places (corridor, house), and objects (car, instrument). Each stimulus was presented on a scrambled background and occupied a square region with dimensions 10° × 10°. Stimuli were presented in 4-s trials. In a trial, 8 images from a given category were presented sequentially, each with duration 0.5 s. Each run lasted 312.0 s and included 6 presentations of each of the 10 categories as well as blank trials (also of 4-s duration). Throughout stimulus presentation, a small red fixation dot was present at the center of the stimulus. Subjects were instructed to maintain fixation on the dot and to press a button whenever they noticed an image in which only the background was present (“oddball” task). A total of 10–12 runs were collected in each scan session.

2.4. MRI data acquisition and pre-processing

Acquisition and pre-processing procedures are the same as described in a previous paper (Kay et al., 2019), except for the addition of a spin-echo acquisition protocol. A summary of all procedures is provided below, and we refer the reader to the previous paper for details.

1 2.4.1. Acquisition

2

3 MRI data were collected at the Center for Magnetic Resonance Research at the University of Minnesota.
4 Some data were collected using a 7T Siemens Magnetom scanner equipped with SC72 body gradients
5 and a custom 4-channel-transmit, 32-channel-receive RF head coil. Other data were collected using a 3T
6 Siemens Prisma scanner and a standard Siemens 32-channel RF head coil. Head motion was mitigated
7 using standard foam padding.

8

9 Anatomical data were collected at 3T at 0.8-mm isotropic resolution. We used a whole-brain T1-weighted
10 MPRAGE sequence (TR 2400 ms, TE 2.22 ms, TI 1000 ms, flip angle 8°, bandwidth 220 Hz/pixel, no
11 partial Fourier, in-plane acceleration factor (iPAT) 2, TA 6.6 min/scan) and a whole-brain T2-weighted
12 SPACE sequence (TR 3200 ms, TE 563 ms, bandwidth 744 Hz/pixel, no partial Fourier, in-plane
13 acceleration factor (iPAT) 2, TA 6.0 min/scan). Several T1 and T2 scans were acquired for each subject
14 in order to increase signal-to-noise ratio.

15

16 Functional data for Experiments E1 and E2 were collected at 7T using gradient-echo EPI at 0.8-mm
17 isotropic resolution with partial-brain coverage (84 oblique slices covering occipitotemporal cortex, slice
18 thickness 0.8 mm, slice gap 0 mm, field-of-view 160 mm (FE) × 129.6 mm (PE), phase-encode direction
19 inferior-superior (F >> H in Siemens' notation), matrix size 200 × 162, TR 2.2 s, TE 22.4 ms, flip angle 80°,
20 echo spacing 1 ms, bandwidth 1136 Hz/pixel, partial Fourier 6/8, in-plane acceleration factor (iPAT) 3,
21 multiband slice acceleration factor 2). Gradient-echo fieldmaps were also acquired for post-hoc correction
22 of EPI spatial distortion (same slice slab as the EPI data, resolution 2 mm × 2 mm × 2.4 mm, TR 391 ms,
23 TE1 4.59 ms, TE2 5.61 ms, flip angle 40°, bandwidth 260 Hz/pixel, no partial Fourier, TA 1.3 min).
24 Fieldmaps were periodically acquired over the course of each scan session to track changes in the
25 magnetic field.

26

27 Functional data for Experiment E3 were collected at 7T using spin-echo EPI at 1.05-mm isotropic
28 resolution with partial-brain coverage (64 (or 48 for Dataset D14) slices, slice thickness 1.05 mm, slice
29 gap 0 mm, field-of-view 128 mm (FE) × 111.2 mm (PE), phase-encode direction inferior-superior (F >> H
30 in Siemens' notation; Dataset D13 was reversed H << F), matrix size 122 × 106, TR 2.2 s, TE 39 ms, flip
31 angle 90°, echo spacing 1 ms, bandwidth 1138 Hz/pixel, partial Fourier 6/8, in-plane acceleration factor
32 (iPAT) 2, multiband slice acceleration factor 2). Corresponding gradient-echo fieldmaps were also
33 acquired.

34

35 Functional data for Experiment E4 were collected at 3T using gradient-echo EPI at 2.4-mm isotropic
36 resolution with partial-brain coverage (30 slices, slice thickness 2.4 mm, slice gap 0 mm, field-of-view 192
37 mm (FE) × 192 mm (PE), phase-encode direction anterior-posterior (A >> P in Siemens' notation), matrix
38 size 80 × 80, TR 1.1 s, TE 30 ms, flip angle 62°, echo spacing 0.55 ms, bandwidth 2232 Hz/pixel, no
39 partial Fourier, no in-plane acceleration, multiband slice acceleration factor 2). Corresponding gradient-
40 echo fieldmaps were also acquired.

41

42 2.4.2. Pre-processing

43

44 T1- and T2-weighted anatomical volumes were corrected for gradient nonlinearities, co-registered, and
45 averaged (within modality). The averaged T1 volume (0.8-mm resolution) was processed using
46 FreeSurfer (Fischl, 2012) version 6 beta (build-stamp 20161007) with the *-hires* option. We generated 6
47 cortical surfaces spaced equally between 10% and 90% of the distance between the pial surface and the
48 boundary between gray and white matter, increased the density of surface vertices by bisecting each
49 edge, and truncated the surfaces retaining only posterior cortex in order to reduce memory requirements.
50 The resulting surfaces are termed 'Depth 1' through 'Depth 6' where 1 corresponds to the outermost

1 surface and 6 corresponds to the innermost surface. Cortical surface visualizations were generated using
2 nearest-neighbor interpolation of surface vertices onto image pixels.

3
4 Functional data were pre-processed by performing one temporal resampling and one spatial resampling.
5 The temporal resampling consisted of one cubic interpolation of each voxel's time-series data; this
6 interpolation corrected differences in slice acquisition times and also upsampled the data to 1.0 s (**Figure**
7 **1B**). Data were prepared such that the first time-series data point coincides with the acquisition time of
8 the first slice acquired in the first EPI volume. The motivation for upsampling is to exploit the intrinsic jitter
9 between the data acquisition and the experimental paradigm (**Figure 1B**). The spatial resampling
10 consisted of one cubic interpolation of each volume; this interpolation corrected head motion (rigid-body
11 transformation) and EPI distortion (determined by regularizing the fieldmaps and interpolating them over
12 time) and also mapped the functional volumes onto the cortical surface representations (affine
13 transformation between the EPI data and the averaged T2 volume). After pre-processing, the data
14 consisted of EPI time series sampled every 1.0 s at the vertices of the depth-dependent cortical surfaces
15 (Depth 1–6). Finally, for the purposes of identifying vertices affected by venous susceptibility effects, we
16 computed the mean of the EPI time-series data obtained for each vertex and divided the EPI intensities
17 by a fitted 3D polynomial (up to degree 4); this produced bias-corrected EPI intensities that can be
18 interpreted as percentages (e.g. 0.8 means 80% of the brightness of typical EPI intensities).

20 2.5. Data analysis

21
22 Code used for data analysis is available at <http://github.com/kendrickkay/>. The core functions that
23 implement the TDM method are referenced by name in the text below. Sample data and scripts
24 demonstrating the TDM method are available at <https://osf.io/j2wsc/>.

26 2.6. GLM analysis

27
28 We analyzed the pre-processed time-series data using three different GLM models (FIR, Standard, TDM).

29
30 The first GLM model, termed FIR (finite impulse response), is a GLM in which separate regressors are
31 used to model each time point in the response to each experimental condition (Dale, 1999). Results from
32 this model are used as inputs to the TDM method. The FIR model characterized the response from 0 s to
33 30 s after condition onset, yielding a total of 31 regressors for each condition. (Modeling the response to
34 30 s was sufficient to capture the majority of the hemodynamic responses; see **Figure 5** and
35 **Supplementary Figure 1**.) We divided the trials for each experimental condition into 2 groups using a
36 “condition-split” strategy (Kay et al., 2019), thereby producing two estimates for each response
37 timecourse. Fitting the FIR model produced BOLD response timecourses (timecourses of ‘betas’) with
38 dimensionality N vertices \times 6 depths \times M conditions \times 31 time points \times 2 condition-splits where N is the
39 number of surface vertices for a given subject and M is the number of conditions in the experiment.

40
41 The second GLM model, termed Standard, is a GLM in which a canonical hemodynamic response
42 function (*getcanonicalhrf.m*) is convolved with stimulus onsets to create a regressor for each experimental
43 condition. We used six condition-splits, thereby producing six response estimates for each condition.
44 Fitting the Standard model produced BOLD response amplitudes (‘betas’) with dimensionality N vertices \times
45 6 depths \times M conditions \times 6 condition-splits.

46
47 The third GLM model, termed TDM, is a GLM in which two hemodynamic timecourses (Early, Late) are
48 separately convolved with stimulus onsets to create two regressors for each experimental condition. The
49 exact nature of these timecourses is determined by the TDM method as described in Section 2.8. We
50 used six condition-splits, thereby producing six response estimates for each combination of condition and

1 timecourse. Fitting the TDM model produced BOLD response amplitudes ('betas') with dimensionality N
2 vertices \times 6 depths \times M conditions \times 2 timecourses \times 6 condition-splits.

3
4 GLMs were prepared and fit to the data using GLMdenoise (Charest et al., 2018; Kay et al., 2013). In
5 GLMdenoise, the GLM consists of experimental regressors (which may take on different forms, as
6 described above), polynomial regressors that characterize the baseline signal level in each run, and data-
7 derived nuisance regressors. In the case of the FIR model, experimental regressors consisted of binary
8 values (0s and 1s). In the case of the Standard and TDM models, experimental regressors consisted of
9 the convolution of stimulus onsets (1s) with hemodynamic timecourses that are normalized to peak at 1
10 (for example, see **Figure 2C**). After fitting the GLMs, estimated betas were converted from raw scanner
11 units to units of percent BOLD signal change by dividing by the mean signal intensity observed at each
12 vertex and multiplying by 100.

13
14 Betas were further analyzed using simple summary metrics. To quantify overall BOLD activity at a given
15 vertex, we calculated *mean absolute beta* (e.g. **Figure 6, left**) by averaging betas across condition-splits,
16 taking the absolute value of the results, and then averaging across conditions. To summarize responses
17 to the eccentricity stimuli, we calculated *peak eccentricity* (e.g. **Figure 6, middle**) by averaging betas
18 across condition-splits, performing positive half-wave rectification (i.e. setting negative values to zero),
19 and then calculating center-of-mass (Hansen et al., 2007). Specifically, center-of-mass was calculated as
20 the weighted average of the integers 1–6 (corresponding to the 6 ring eccentricities from fovea to
21 periphery) using the rectified betas as weights.

22 23 2.7. Timecourse quantification and metrics

24
25 In some analyses (see **Figure 2A–B, Supplementary Figure 1**), we summarize the typical timecourse
26 shape observed in a set of timecourses. This was accomplished using a PCA-based procedure
27 (*derivehrf.m*). In the procedure, we first subtract off the mean of each timecourse. We then perform
28 principal components analysis (PCA) on the entire set of timecourses and extract the first PC (this is the
29 timecourse vector along which variance in the total set of timecourses is maximized). Next, we add a
30 constant offset to the first PC such that the first time point equals 0, and flip the sign of the PC such that
31 the mean over the range 0–10 s is positive. Finally, we calculate the weight that minimizes squared
32 reconstruction error for each timecourse, compute the absolute value of these weights, and then scale the
33 PC by the average weight. The motivation for de-meaning the timecourses prior to PCA is to suppress
34 low-frequency noise present in weak BOLD responses. In general, the use of PCA for summarizing
35 timecourse shape (Kay et al., 2008a) has advantages over simply computing the mean timecourse: PCA
36 elegantly handles negative BOLD timecourses, and PCA allows timecourses with larger BOLD responses
37 to have greater influence on the resulting timecourse shape (thereby producing more robust results).

38
39 Several timecourse metrics were computed (see **Figure 2, Figure 5, Supplementary Figure 1**). Given a
40 timecourse, we upsampled the timecourse to a sampling rate of 0.01 s using sinc interpolation. We then
41 identified the maximum of the resulting timecourse (*peak amplitude*) and its associated time (*time-to-*
42 *peak*). We used linear interpolation to calculate the time at which the timecourse rises to half of the
43 maximum value (*rise time*) and the time at which the timecourse falls to half of the maximum value (*fall*
44 *time*). Finally, we computed the time elapsed between the rise time and the fall time (*full-width-at-half-max*
45 *or FWHM*).

46 47 2.8. TDM method

48 49 2.8.1. Theory

1 TDM is a data-driven technique that identifies a principal axis of timecourse variation present in a set of
2 experimentally measured timecourses. It does this by examining timecourses projected into a low-
3 dimensional space (3 dimensions) and extracting a 1-dimensional manifold (a line) that captures the
4 variation of interest. The procedure can be viewed as a powerful method for summarizing and extracting
5 the signal present in timecourses which considered individually (i.e. one response timecourse at a time)
6 would likely be insufficiently reliable. In our fMRI measurements of responses to 4-s visual stimuli, we
7 consistently find that one endpoint of the line corresponds to an early timecourse peaking at around 5–7 s
8 and the other endpoint of the line corresponds to a late timecourse peaking at around 6–9 s (see **Figure**
9 **2D, Figure 5A, Supplementary Figure 1**). These timecourses are interpreted as reflecting hemodynamic
10 responses from the microvasculature (e.g. capillaries and venules) and hemodynamic responses from the
11 macrovasculature (e.g. veins), respectively. TDM then uses the identified timecourses in a regression
12 model in order to decompose observed hemodynamic responses into early and late components. The
13 researcher can choose to analyze further the early component, the late component, or both.

14
15 There are three main quantities involved in the TDM technique: *density*, referring to the timecourse
16 shapes that tend to be present in the data; *vector length*, referring to the amplitudes of the timecourses in
17 the data; and *EPI intensity*, referring to the bias-corrected EPI intensity of the vertex (or voxel) to which
18 each timecourse belongs. TDM combines density and vector length and fits an oriented 2D Gaussian to
19 the result in order to identify the 1-dimensional manifold. The motivation for incorporating vector length
20 beyond density alone is to ensure that veins—which generate BOLD responses with large amplitudes but
21 constitute only a fraction of the total set of responses—have sufficient influence on the determination of
22 the manifold. Note that EPI intensity does not directly participate in the determination of the manifold, and
23 can therefore provide useful validation of the manifold results (see **Figure 4A, Supplementary Figures**
24 **2–3**).

25
26 There are a few important conceptual points regarding the nature of the TDM method. For any given
27 voxel (or vertex), the BOLD response to an experimental event is expected to reflect a mixture of early
28 (microvasculature) and late (macrovasculature) timecourses. The specific proportions of these
29 timecourses is expected to vary from voxel to voxel simply due to heterogeneity in the spatial structure of
30 the vasculature (e.g., one voxel might be centered on a large vein, whereas another voxel may only
31 partially overlap the vein). Different proportions of the timecourses manifest in TDM as different points
32 along the 1-dimensional manifold, and these points collectively trace out an arc on the unit sphere (see
33 **Figure 3**).

34
35 Another important point is that TDM is not equivalent to estimating a different hemodynamic response
36 function for each voxel (e.g., Kay et al., 2008b, 2008a; Pedregosa et al., 2015). Using a single
37 hemodynamic timecourse for different experimental conditions (time-condition separability) yields at most
38 one amplitude estimate (beta) for each condition. In contrast, TDM allows experimental conditions to have
39 different loadings on the early and late timecourses, and yields two amplitude estimates (betas) for each
40 condition (see **Figure 6, inset boxes**). This is the critical feature that gives TDM the power to remove
41 unwanted venous effects from a given voxel's response.

42
43 Finally, note that the GLM analyses performed in TDM rely on the assumption of linear summation of
44 BOLD responses over time. Our experiments, like many used in cognitive neuroscience, involve a large
45 number of trials that are presented fairly rapidly in order to maximize statistical power (e.g. less than 10 s
46 of rest in between trials). In such experiments, it is a practical necessity to assume temporal linearity.
47 Moreover, nonlinear effects are likely to average out when using randomized experimental designs.
48 Nevertheless, the accuracy with which TDM identifies timecourses may be limited if there exist nonlinear
49 effects (Friston et al., 1998b) and especially if these nonlinearities vary for different types of vasculature
50 (Goodyear and Menon, 2001; Thompson et al., 2014; Zhang et al., 2009).

51

1 2.8.2. Algorithm

2

3 The TDM algorithm starts with a set of timecourses and determines a pair of timecourses that
4 characterize the overall variation in the timecourses. The following are the steps in the TDM algorithm
5 (*extracthrfmanifold.m*):

- 6 1. *Perform PCA on the timecourses.* We collect timecourses into a 2D matrix of dimensionality L
7 timecourses $\times T$ time points, and then perform singular value decomposition. This produces a
8 matrix with dimensionality T time points $\times T$ eigenvectors where columns correspond to principal
9 component (PC) timecourses in decreasing order of variance explained. We flip the sign of the
10 first PC if necessary to ensure that the mean of the timecourse is positive over the range 0–10 s.
11 Note that previous studies have applied PCA to fMRI response timecourses but in different
12 contexts (d'Avossa et al., 2003; Woolrich et al., 2004).
- 13 2. *Use PC1–PC3 to define a 3-dimensional space for further analysis.* Our convention for
14 visualization is that PC1 points out of the page (positive z-axis), PC2 points to the right (positive
15 x-axis), and PC3 points to the top (positive y-axis).
- 16 3. *Map timecourses onto the unit sphere.* We project each timecourse onto PC1, PC2, and PC3,
17 and mirror coordinates across the origin if necessary to ensure that the loading on PC1 is
18 positive. (Thus, negative BOLD timecourses are flipped and treated in the same way as positive
19 BOLD timecourses.) Each timecourse is represented by a set of coordinates (loadings) and can
20 be interpreted as a 3-dimensional vector. We normalize each vector to unit length (thereby
21 placing the vector on the unit sphere), and also save the original vector length for later use.
- 22 4. *Calculate a 2D image that represents density.* We orthographically project timecourses onto the
23 xy plane, and then calculate a 2D histogram. This produces a 2D image where pixel values
24 represent frequency counts. This image indicates typical timecourse shapes found in the data.
- 25 5. *Calculate 2D images that represent vector length and EPI intensity.* Using the same orthographic
26 projection and binning scheme of Step 4, we calculate the median vector length of the
27 timecourses found in each bin. We also calculate the median bias-corrected EPI intensity of the
28 voxels (vertices) associated with the timecourses in each bin. This produces 2D images where
29 pixel values represent vector lengths (indicating timecourse amplitudes) and EPI intensities
30 (indicating static susceptibility effects caused by veins), respectively.
- 31 6. *Regularize the density image by subtracting a DC bias.* We distribute a collection of particles on
32 the unit sphere (S2 Sampling Toolbox, <https://www.github.com/AntonSemechko/S2-Sampling-Toolbox>),
33 assign timecourses to their nearest particles, and count the number of timecourses
34 associated with each particle. We then calculate a histogram of these counts and determine the
35 bin B with the highest frequency. Finally, we stochastically subsample the timecourses such that
36 the number of timecourses associated with each particle is reduced by the middle value of bin B .
37 The resulting subsampled timecourses are used to generate a new density image.
- 38 7. *Scale the density image.* The regularized density image from Step 6 is scaled such that 0 maps to
39 0 and the maximum value maps to 1. Values are then truncated to the range [0, 1].
- 40 8. *Regularize the vector-length image by subtracting a DC bias.* This is accomplished in a similar
41 manner as Step 6: we calculate a histogram of the values in the vector-length image, determine
42 the bin B with the highest frequency, and subtract the middle value of bin B from all image pixels.
- 43 9. *Scale the vector-length image.* The regularized vector-length image from Step 8 is scaled such
44 that 0 maps to 0 and the maximum value maps to 1. Values are then truncated to the range [0, 1].
- 45 10. *Average the density and vector-length images.* Although the default behavior is to simply average
46 the density and vector-length images, if the user desires a different weighting (e.g. giving more
47 weight to the vector-length image), a flag can be used (*opt.vlengthweight*).
- 48 11. *Fit 2D Gaussian.* The image resulting from Step 10 is fit with a 2D Gaussian. The Gaussian is
49 controlled by two parameters specifying the center, two parameters specifying the spreads along
50 the major and minor axes, a rotation parameter, a gain parameter, and an offset parameter. In
51 model fitting, the error metric is set up such that the image is interpreted as a probability

1 distribution (pixels with larger values reflect higher density and thus contribute more heavily to the
2 error metric).

3 12. *Extract two points along the major axis of the Gaussian.* We determine points corresponding to
4 the mean \pm one standard deviation along the major axis of the Gaussian. The choice of one
5 standard deviation is somewhat arbitrary but appears to produce satisfactory results.

6 13. *Reconstruct timecourses corresponding to the identified points.* We place the points determined
7 in step 12 on the unit sphere, and use their associated coordinates to weight and sum the PC1,
8 PC2, and PC3 timecourses. This yields two reconstructed timecourses. Based on time-to-peak
9 (see Section 2.7), we label one timecourse as ‘Early’ and the other timecourse as ‘Late’.

10

11 2.8.3. Application of algorithm

12

13 We obtained timecourses by fitting an FIR model to the fMRI time-series data (as described in Section
14 2.6). We then calculated the amount of variance explained (R^2) by the FIR model. In order to focus the
15 TDM algorithm on cortical locations with BOLD responses, we selected all vertices within a region of
16 interest (see Section 2.9) that exceeded an automatically determined threshold (specifically, the value at
17 which the posterior probability switches between two Gaussian distributions fitted as a mixture model to
18 the data; see *findtailthreshold.m*). This produced a set of timecourses with dimensionality P vertices \times M
19 conditions \times 31 time points \times 2 condition-splits. We applied the TDM algorithm to the timecourses
20 averaged across the condition-splits (thus, reflecting the entire dataset) and also to the timecourses from
21 each condition-split separately in order to assess reliability. In both cases, the number of timecourses
22 given to the TDM algorithm is $L = P \times M$ and the number of time points is $T = 31$. After completion of the
23 TDM algorithm, the identified timecourses (Early, Late) were incorporated into a GLM model to
24 decompose the fMRI time-series data into early and late components (see Section 2.6).

25

26 2.8.4. Alternative ICA-based procedure

27

28 Independent components analysis (ICA) is a potential alternative method for determining latent
29 timecourses. We designed an ICA-based procedure (*icadecomposehrf.m*) that serves as a drop-in
30 replacement for the TDM algorithm described above. In the procedure, we take the two condition-split
31 versions of the FIR-derived timecourses and perform ICA on each set of timecourses (FastICA Toolbox,
32 <https://research.ics.aalto.fi/ica/fastica/>, default nonlinearity). This produces 31 independent component
33 (IC) timecourses for each split of the data. Because ICA does not provide a natural ordering or grouping
34 of the ICs, the challenge is to determine which specific pair of ICs to use as the latent timecourses in
35 TDM.

36

37 To determine the pair of ICs to use, we devised the following heuristic procedure. We first greedily order
38 the ICs within each split to maximize variance explained. That is, we choose the IC that maximizes
39 variance explained in the timecourses, choose a second IC that, when combined with the first, maximizes
40 variance explained, and so on. For normalization purposes, we flip each IC if necessary to ensure that it
41 is positive over the range 0–10 s and normalize it to unit length. Next, we perform greedy matching in
42 order to match the ICs in the second split of the data to the ICs in the first split. Specifically, we choose
43 the IC in the second split that is most similar in a squared-error sense to the first IC in the first split,
44 choose the remaining IC in the second split that is most similar to the second IC in the first split, and so
45 on. To reduce the number of ICs under consideration, we select ICs that both exhibit high consistency
46 across the two splits of the data ($R^2 > 50\%$) and are within the top 50% of the ICs in the first split (with
47 respect to the ordering based on variance explained). Finally, from the ICs that meet these selection
48 criteria, we perform an exhaustive search to determine the unique pair of ICs that explain the most
49 variance in the original timecourses. Based on time-to-peak (see Section 2.7), we label one IC as ‘Early’
50 and one IC as ‘Late’.

51

1 2.9. Region-of-interest (ROI) definition

2

3 We used two regions-of-interest (ROIs). For the eccentricity experiment, we used the union of visual
4 areas V1, V2, and V3 from a publicly available atlas of visual topography (Wang et al., 2015). For the
5 category experiment, we used a manually defined region in occipital, parietal, and temporal cortex that
6 covers visually responsive vertices (same region used in Kay et al., 2019). Both ROIs were defined in
7 FreeSurfer's *fsaverage* space and backprojected to individual subjects for the purposes of vertex
8 selection.

9

3. Results

We collected BOLD fMRI measurements in human visual cortex while subjects viewed briefly presented visual stimuli (3.5-s or 4-s duration). The main datasets are D1–D5 in which a 7T gradient-echo 0.8-mm 2.2-s protocol was used to measure responses to rings at different eccentricities (**Figure 1A**). Additional datasets D6–D12 used the same acquisition protocol but involved images of different semantic categories (e.g. faces). The data were pre-processed by performing one temporal interpolation that corrected slice-time differences and prepared the data at a 1-s sampling rate (**Figure 1B**) and one spatial interpolation that corrected head motion and EPI distortion and sampled the fMRI volumes at the locations of cortical surface vertices. Given the high resolution of the fMRI data, we prepared multiple cortical surface representations at different depths through the gray matter.

3.1. Systematic variation in timecourse amplitude, delay, and width

To investigate BOLD timecourse characteristics, we analyzed the time-series data in each dataset using a finite impulse response model (0–30 s, 31 time points). This produced, for each surface vertex, an estimate of the BOLD response timecourse to each experimental condition. We then binned these timecourses either with respect to cortical depth or with respect to EPI intensity. For each bin, we summarized the timecourses found in that bin using a PCA-based procedure (see Methods for details).

Results for a representative dataset are shown in **Figures 2A–B**. Proceeding from inner cortical depths (Depth 6) to outer cortical depths (Depth 1), we observe an increase in timecourse amplitude and an increase in timecourse delay (**Figure 2A**). Proceeding from darker EPI intensities (0–0.5) to lighter EPI intensities (>1), we again observe increases in amplitude and delay but also a small increase in timecourse width (**Figure 2B**). Summarizing results across different datasets (reflecting different subjects and experiments), we see that these effects are consistently observed (**Figure 2D**). Notice that although different datasets show similar patterns in relative timing (e.g., time-to-peak is longer for low EPI intensity than for high EPI intensity), there is large variance in absolute timing across datasets (e.g., time-to-peak is approximately 8 s in one dataset but 6 s in another dataset). This is consistent with well-established observations of variability of hemodynamic response functions across subjects (Handwerker et al., 2012).

Notice that the variations in amplitude (**Figure 2D**, top row) and the variations in timecourse delays (**Figure 2D**, middle rows) are more pronounced when binning by EPI intensity (**Figure 2D**, middle column) than when binning by cortical depth (**Figure 2D**, left column). We suggest that the underlying source of these effects consists in the high-amplitude, delayed BOLD responses carried by macroscopic veins. Binning by EPI intensity provides a relatively direct proxy for where these venous effects occur (Kay et al., 2019; Menon et al., 1993; Ogawa et al., 1990). In contrast, binning by cortical depth provides a less direct proxy for these effects (e.g., pial veins affect outer cortical depths more than inner cortical depths). Thus, binning by EPI intensity will tend to accentuate and highlight the amplitude and delay effects.

The fact that veins carry delayed BOLD responses has been previously shown (de Zwart et al., 2005; Kim and Ress, 2017; Lee et al., 1995; Siero et al., 2011). Prior studies have also demonstrated increased temporal delays at superficial cortical depths (Kim and Ress, 2017; Siero et al., 2011). Thus, the observations we make here are not novel, but provide reassurance that these effects can be reproduced in our dataset and establish a starting point for the development of the TDM method.

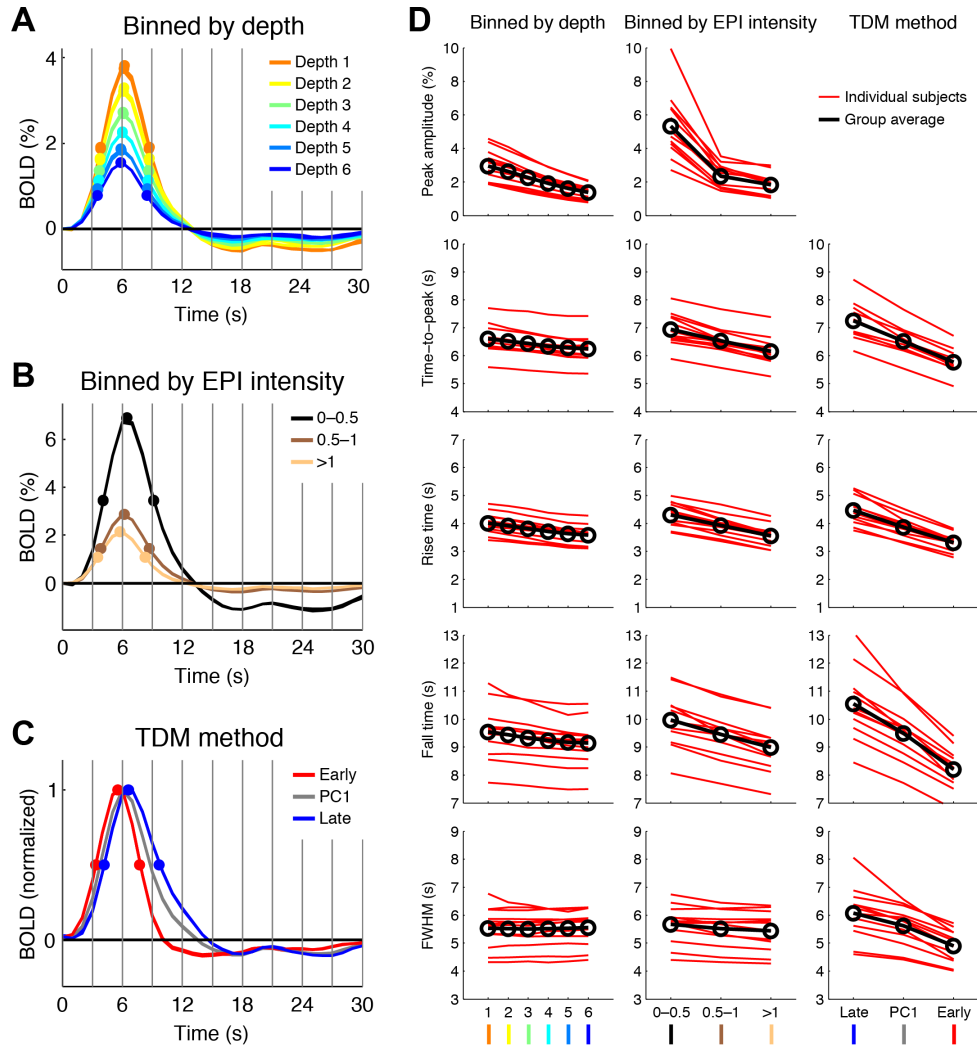
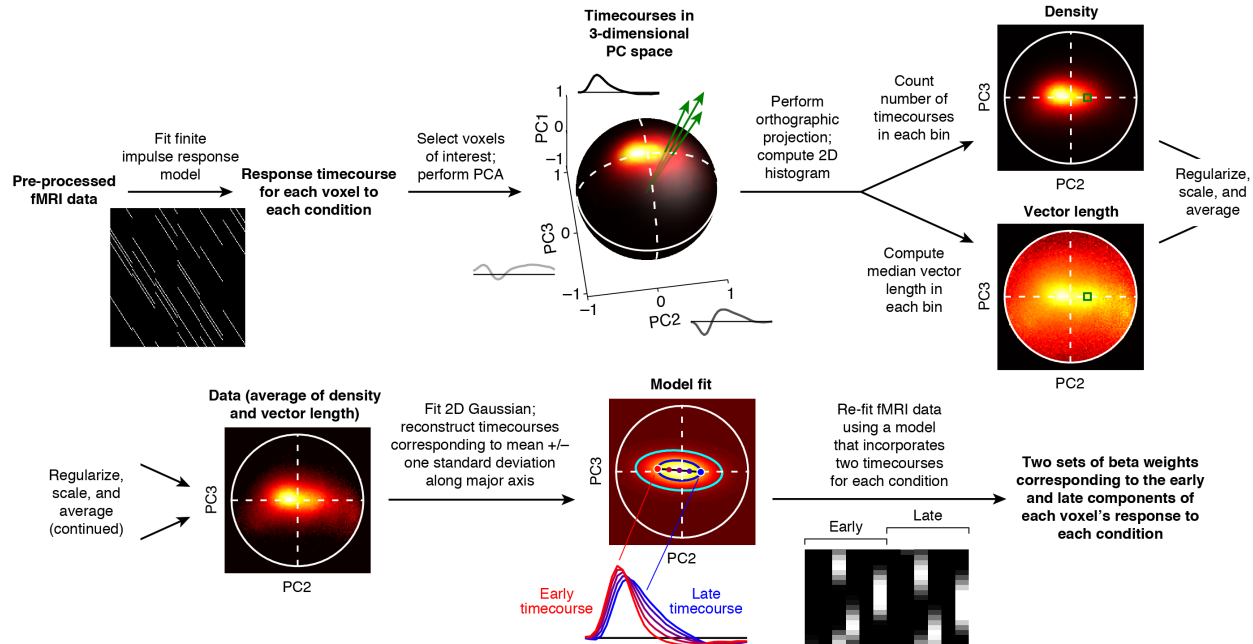


Figure 2. BOLD timecourses exhibit systematic variation in amplitude, delay, and width.

Panels A–C show detailed results for Dataset D10 (detailed results for all datasets are shown in **Supplementary Figure 1**). *A*, Timecourses binned by cortical depth (Depth 1 is superficial; Depth 6 is deep). Timecourses from split-halves of the data are shown for each bin (the two sets of traces are nearly identical, indicating high reliability). Vertical gray lines mark 3-s intervals, a convention used throughout this paper. Solid dots mark timecourse peak, rise time (time at which the signal rises to half of the peak value), and fall time (time at which the signal falls to half of the peak value). Full-width-half-max (FWHM) is calculated as fall time minus rise time. *B*, Timecourses binned by EPI intensity. Same as panel A except binning is performed with respect to bias-corrected EPI intensity. *C*, Timecourses derived by TDM. The Early and Late timecourses derived by TDM are normalized to peak at 1. For comparison, we also show the first PC of the timecourses, also normalized to peak at 1. *D*, Summary of results across datasets. Timecourse metrics obtained for individual subjects (Datasets D1–D12) and corresponding group averages are shown.

3.2. TDM provides a method for visualizing timecourse variation

As we have seen, there is systematic variation in the hemodynamic timecourses observed in our data. Given that capturing timecourse variation lies at the core of the TDM method, we are now ready to examine results of applying TDM to our data.

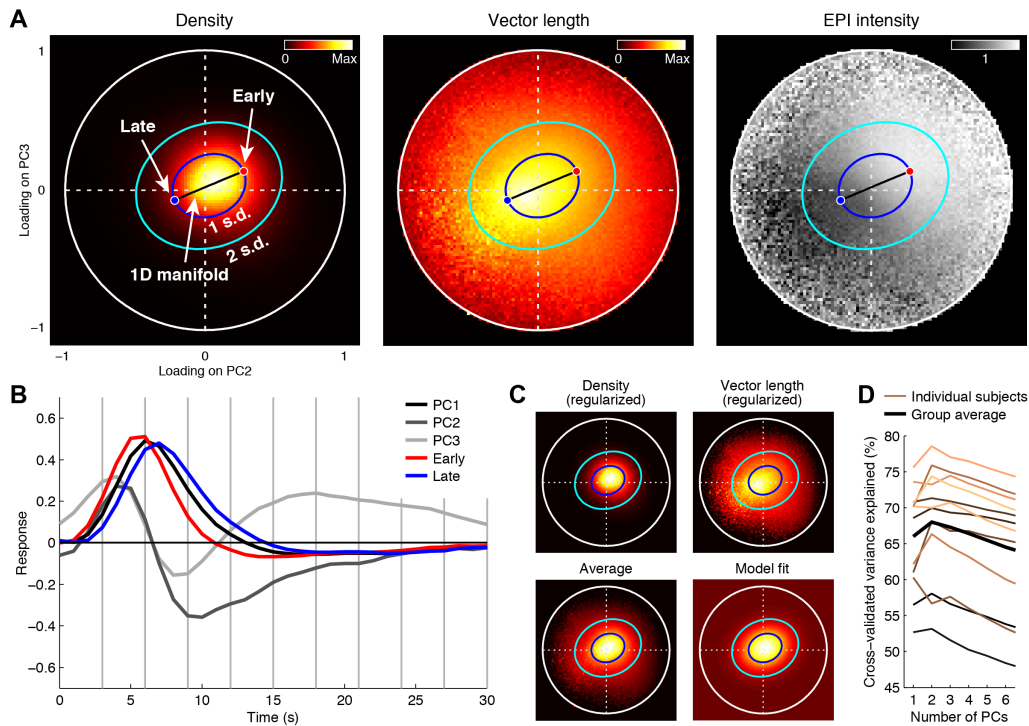


1
2
3
4 **Figure 3. Schematic of the TDM method.** Time-series data are fit with a finite impulse response
5 model to estimate response timecourses. PCA is performed on the timecourses to reduce their
6 dimensionality to 3. Using orthographic projection in the direction of the first PC, a 2D histogram
7 image is calculated (*density*). The same projection and binning scheme is used to calculate an
8 image representing timecourse amplitudes (*vector length*). The two images are combined and fit
9 with a 2D Gaussian in order to determine an early timecourse and a late timecourse that together
10 summarize the principal axis of variation. Finally, the time-series data are re-fit with a model
11 incorporating the two timecourses.

12 The first step in TDM is to visualize distributions of response timecourses in a low-dimensional space
13 (**Figure 3**). Specifically, we use principal components analysis (PCA) to determine the three orthogonal
14 timecourses that account for the most variance in the given timecourses. We then use these three
15 timecourses as axes of a 3D space (PC1, PC2, PC3) in which each of the original timecourses
16 corresponds to one point, or vector, in this space. To visualize the results, we map the timecourse vectors
17 to the unit sphere (by normalizing them to unit length) and use an orthographic projection to visualize the
18 density of the timecourses. Such a visualization reveals typically occurring timecourse shapes,
19 independent of timecourse amplitude (see 'Density' image in **Figure 3**). We separately visualize the
20 amplitude of the timecourses by computing the length of the original timecourse vectors and repeating the
21 orthographic visualization (see 'Vector length' image in **Figure 3**).

22
23 Applying these visualization procedures to a representative dataset (similar results are obtained in every
24 dataset; see **Supplementary Figures 2–3**), we find that timecourse shapes typically reside near the pole
25 of the unit circle where PC1 is maximal, with some variability around this pole (**Figure 4A, left**). We find
26 that timecourse amplitudes are large in a similar portion of the space except for a small extension towards
27 the lower left (**Figure 4A, middle**). A separate plot shows the actual timecourses associated with the
28 three PCs that define the axes of the space (**Figure 4B**). This plot shows that timecourse shapes
29 generally resemble a canonical hemodynamic response timecourse (**Figure 4B, black line**), with a major
30 axis of variation corresponding to different loadings on a timecourse that shifts the peak of the canonical
31 timecourse either earlier or later in time (**Figure 4B, dark gray line**). Note that the general shapes of the
32 PC timecourses are similar to those found in other applications of PCA to fMRI timecourses (d'Avossa et
33 al., 2003; Woolrich et al., 2004).

1



2
3

4 **Figure 4. TDM captures timecourse variation along a 1D manifold.** Panels A–C show results
5 for Dataset D9 (results for all datasets are shown in **Supplementary Figures 2 and 3**). A,
6 Quantities of interest. TDM calculates a density image indicating frequently occurring timecourse
7 shapes (left), a vector-length image indicating timecourse amplitudes (middle), and an EPI-intensity
8 image indicating bias-corrected EPI intensities (right). The black line indicates the identified 1D
9 manifold that connects the Early and Late timecourses. (Ellipses, dots, and lines are identical in all
10 plots.) B, Timecourses. All timecourses are unit-length vectors. C, Gaussian fitting procedure.
11 Density and vector-length images are DC-subtracted, scaled, and truncated, producing regularized
12 images (upper left, upper right). These images are then averaged (lower left) and fit with an
13 oriented 2D Gaussian (lower right). D, Dimensionality of response timecourses. To confirm that
14 three dimensions are sufficient to capture relevant variation in response timecourses, we perform
15 PCA on split-halves of each dataset and assess how well a limited number of PC timecourses from
16 one half reconstruct timecourses measured in the other half (i.e. cross-validation). Results for
17 individual subjects (Datasets D1–D12) and the group average are shown.

18

19 We emphasize the usefulness of the visualizations performed in TDM: the visualizations summarize large
20 amounts of data (response timecourses) in a manner that highlights robust signals and clearly delineate
21 effects of interest in the data. We suggest that the visualizations may be generally useful for investigating
22 timecourse variations across voxels, brain areas, and/or subjects. Although the use of three dimensions
23 in the visualizations is appealing because it is practical to create pictorial representations of a small
24 number of dimensions, visualizing timecourses in only three dimensions might provide an incomplete
25 characterization of the full diversity of timecourses. To investigate this issue, we performed a cross-
26 validation analysis to determine the number of timecourse dimensions necessary to capture signals of
27 interest in the response timecourses (**Figure 4D**). We find that in every dataset, cross-validation
28 performance is maximized using three or fewer PCs. Thus, it appears that using three PC dimensions is
29 sufficient and that additional dimensions would likely be dominated by measurement noise.

30

31 3.3. TDM identifies an axis of timecourse variation

32

1 The second step in TDM is to identify an axis that captures the major variation in the observed response
2 timecourses. As seen in the earlier visualizations (**Figure 4A**), response timecourses empirically occupy a
3 small portion of the 3D space (i.e., timecourse density and vector length are confined to a small section of
4 the unit sphere). Furthermore, the timecourses can be approximated in the 3D space by a simple line
5 segment (arc) that describes variation with respect to timecourse delay, i.e. early vs. late. The
6 interpretation that we adopt here is that (i) the endpoints of the line segment correspond to latent
7 hemodynamic timecourses associated with the microvasculature and the macrovasculature and (ii) any
8 single observed timecourse is simply a mixture of these two latent timecourses plus measurement noise
9 (which causes deviation away from the line). There are a variety of reasons why different mixtures might
10 be observed across voxels in a given dataset. One simple reason is heterogeneity in the spatial structure
11 of the vasculature: some voxels might be in close proximity to a large vein whereas other voxels may
12 largely comprise capillaries and small venules. (Note that one might have expected to find discrete
13 clusters or modes in the timecourse distribution, but this does not appear to be the case in our
14 measurements. It is possible that clustering might be observed at higher spatial resolutions.)
15

16 To calculate the axis of variation, TDM combines density and vector length (**Figure 3, lower left**), fits a
17 2D Gaussian to the result (**Figure 3, bottom**), and extracts points positioned at plus and minus one
18 standard deviation along the major axis of the fitted Gaussian (**Figure 3, bottom, red and blue points**).
19 Examining results obtained on the representative dataset, we see that the procedure works well: the
20 combined image resembles both density and vector length (**Figure 4C, bottom left**), the fitted Gaussian
21 approximates the data (**Figure 4C, bottom right**), and the extracted points reside in sensible locations
22 (**Figure 4A, first two plots**). We reconstruct timecourses corresponding to the two extracted points
23 (**Figure 4B, red and blue lines**), and then label the timecourses 'Early' and 'Late' based on the time-to-
24 peak of the reconstructed timecourses. Note that different mixtures of the Early and Late timecourses
25 trace out an arc (line) on the unit sphere (**Figure 4A, black line**) and result in a continuum of timecourse
26 shapes (**Figure 3, bottom**). This arc is the 1D manifold that describes timecourse variation in the dataset.
27

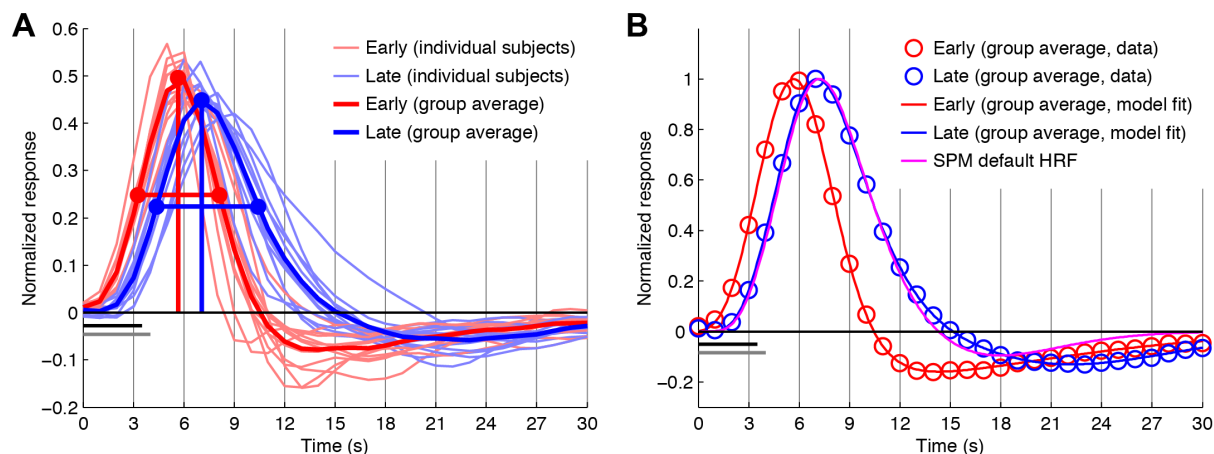
28 We interpret the Early and Late timecourses as reflecting the microvasculature (e.g. capillaries and
29 venules) and macrovasculature (e.g. veins), respectively. Does this interpretation have face validity? We
30 offer several lines of reasoning that suggest validity. First, we find that the Late timecourse is consistently
31 associated with large vector length (see **Figure 4A, middle plot** and **Supplementary Figure 2**),
32 indicating that response timecourses resembling the Late timecourse tend to have large BOLD
33 amplitudes. This makes sense given that veins exhibit large percent BOLD signal changes (Lee et al.,
34 1995; Menon et al., 1993). Second, if we use the same visualization methods (orthographic projection of
35 the unit sphere) to examine the relationship between timecourse shape and bias-corrected EPI intensity,
36 we find that the Late timecourse is consistently positioned in a zone of low EPI intensity (see **Figure 4A,**
37 **right plot** and **Supplementary Figure 2**). Since the TDM procedure does not make use of EPI
38 intensities, this is an empirical finding that provides further evidence of validity, as it is known that veins
39 cause static susceptibility effects in EPI images (Kay et al., 2019; Menon et al., 1993; Ogawa et al.,
40 1990). Third, the idea that veins exhibit delayed BOLD responses is consistent with several previous
41 experimental studies (de Zwart et al., 2005; Kim and Ress, 2017; Lee et al., 1995; Siero et al., 2011).
42 Finally, a biophysical model of vascular dynamics has been proposed, and this model provides a potential
43 explanation for why temporal delays occur in veins (see Fig. 6 in Havlicek and Uludağ, 2020).
44

45 As an additional sanity check, let us compare the Early and Late timecourses identified by TDM against
46 the timecourse inspections presented earlier. We see that the Early and Late timecourses (**Figure 2C**)
47 resemble the ones found through binning by cortical depth (**Figure 2A**) and binning by EPI intensity
48 (**Figure 2B**), but are somewhat more extreme in nature. For example, time-to-peak is earlier for the Early
49 timecourse and later for the Late timecourse compared to the binning-based timecourses. These effects
50 can be seen more clearly in the quantitative summary (**Figure 2D, right column**). Thus, TDM appears to
51 be extracting sensible timecourses (see also **Supplementary Figure 1**). Finally, independent

1 components analysis (ICA) is a commonly used statistical method for deriving latent structure in a set of
2 data, and is commonly applied to fMRI data. We show that applying an ICA-based procedure yields
3 results similar to those obtained from TDM in some datasets, but divergent and unsatisfactory results in
4 other datasets (**Supplementary Figure 2**). We provide additional thoughts regarding ICA in the
5 Discussion.

6 7 3.4. Summary of TDM timecourses across subjects

8
9 For a comprehensive assessment of the TDM results, we plot the TDM-derived Early and Late
10 timecourses obtained for each dataset (**Figure 5A, thin lines**). We see that on the whole, the
11 timecourses are stereotyped and largely consistent across datasets. Keep in mind that TDM is a
12 completely data-driven technique: it makes no assumptions about timecourse shapes that are extracted
13 from the data, and assumes only that timecourses can be characterized in a low-dimensional subspace
14 and can be summarized by a one-dimensional manifold. Thus, the fact that the derived timecourses
15 resemble typical hemodynamic response functions is an empirical finding, and indicates that the TDM
16 technique successfully derives reasonable timecourses when applied to individual datasets.
17



18
19
20 **Figure 5. Early and late timecourses found by TDM.** A, Summary of results. We plot unit-length-
21 normalized Early and Late timecourses found in Datasets D1–D12 (thin lines) and their average
22 (thick lines). Dots mark timecourse peaks and timecourse rise and fall times. Horizontal black and
23 gray bars below the x-axis indicate the stimulus duration in the eccentricity and category
24 experiments (3.5 s and 4 s, respectively). B, Group-average results and parametric model fit.
25 Group-average Early and Late timecourses are normalized to peak at 1 (red and blue circles), and
26 are fit with a double-gamma function as implemented in SPM's *spm_hrf.m* (red and blue lines). The
27 fitting was achieved by convolving a double-gamma function with a 4-s square wave and
28 minimizing squared error with respect to the data. Estimated double-gamma parameters for the
29 group-average Early and Late timecourses were [7.21 17.6 0.5 4.34 1.82 -3.09 50] and [5.76 21.6
30 1.11 1.72 3.34 0.193 50], respectively. As a comparison, we show the timecourse obtained using
31 the default SPM parameters [6 16 1 1 6 0 32] (magenta line).
32

33 While there are qualitative similarities of the Early and Late timecourses across datasets, there is also
34 substantial quantitative variability, consistent with the well-established observation that BOLD
35 timecourses exhibit variability across subjects (Handwerker et al., 2012). It is important to note that the
36 observed variability in timecourses is not due to measurement noise: split-half analyses show that the
37 TDM-derived timecourses are highly stable across splits of each dataset (**Supplementary Figure 1, right
38 column**). These results suggest that in order to achieve the most accurate characterization of responses,
39 one should tailor timecourses to what is empirically observed in individual datasets.

1
2 To summarize overall timecourse results, we compute the group average for the Early and Late
3 timecourses (**Figure 5A, thick lines**). We see that these timecourses are similar in overall shape and
4 differ primarily in their delay and width. We furthermore fit each group-average timecourse using a
5 double-gamma function, and find that smooth parametric functions characterize the empirical results quite
6 well (**Figure 5B**). Finally, as a point of comparison, we plot the predicted timecourse for a 4-s event using
7 the default double-gamma hemodynamic response function implemented in SPM. Interestingly, this
8 timecourse (**Figure 5B, magenta line**) coincides extremely well with the group-average Late timecourse
9 (**Figure 5B, blue line**). This makes sense, given that the default timecourse parameters in SPM were
10 derived from fMRI measurements conducted at low (2T) magnetic field strength (Friston et al., 1998a)
11 where the BOLD response is dominated by contributions from large vessels (Haacke et al., 1994; Uludağ
12 et al., 2009).

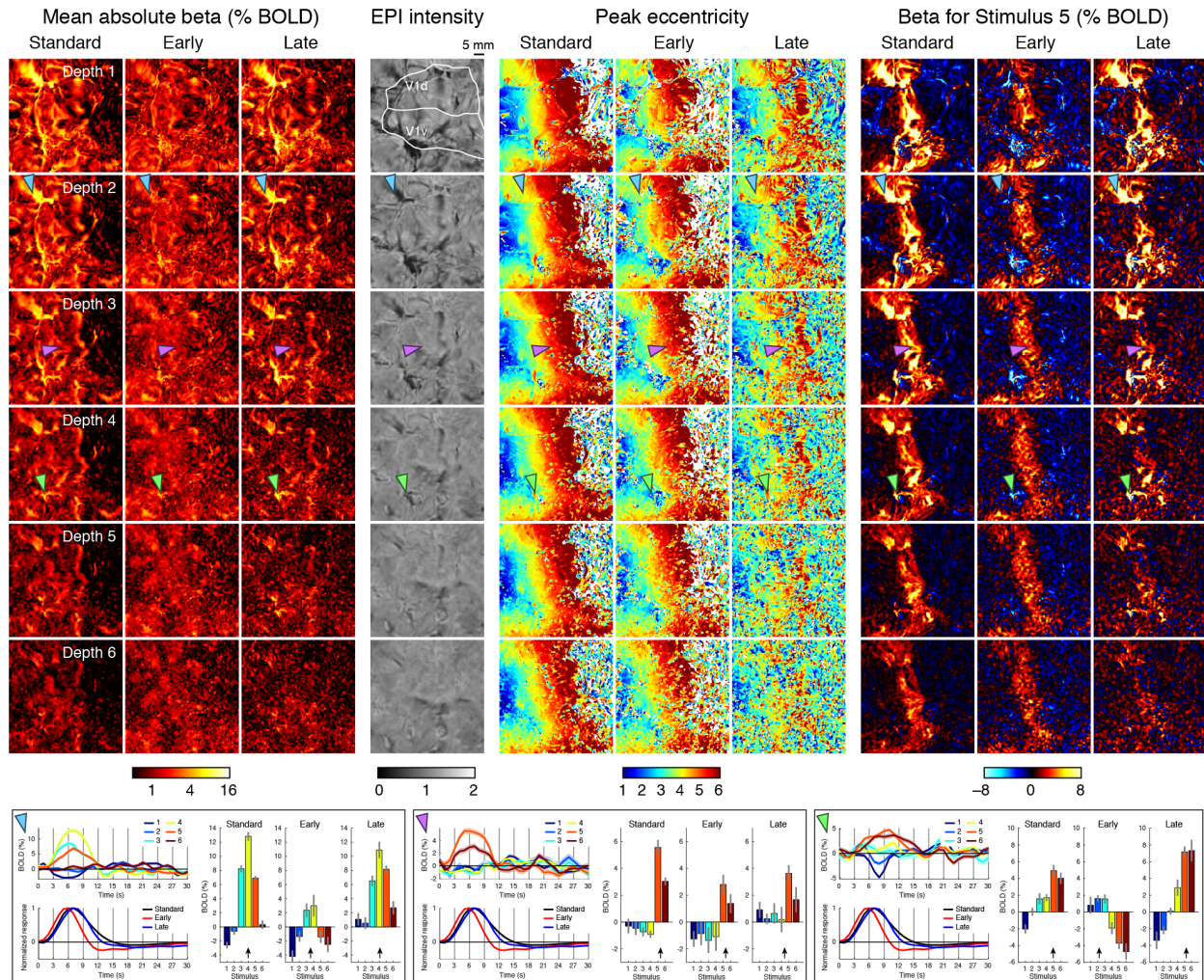
13 14 3.5. Decomposition using TDM timecourses removes artifacts from cortical maps

15
16 To summarize thus far, we have used TDM to derive Early and Late timecourses in each dataset, and we
17 interpret these as reflecting responses from the microvasculature and macrovasculature, respectively.
18 The final step of TDM involves re-analyzing the time-series data using a GLM that incorporates the Early
19 and Late timecourses time-locked to the onset of each experimental condition. Fitting this GLM produces,
20 for each vertex (or voxel) and condition, an estimate of the BOLD response amplitude from the
21 microvasculature and an estimate of the BOLD response amplitude from the macrovasculature. These
22 response amplitudes, or betas, can then be used in subsequent analyses according to the goals of the
23 researcher; for example, one can choose to ignore the beta associated with the Late timecourse and
24 focus on the beta associated with the Early timecourse. Note that the Early and Late timecourses are not
25 necessarily orthogonal and are in fact often quite correlated (in **Figure 4A**, notice the close proximity of
26 the Early and Late timecourses on the unit sphere). Thus, in the GLM model, the two timecourses are, in
27 a sense, competing to account for the response timecourses observed in the data.

28
29 To assess the quality of the Early and Late betas, we generate cortical surface visualizations and
30 compare these against visualizations of betas obtained using a standard GLM that incorporates a single
31 canonical hemodynamic response function time-locked to each condition. We focus specifically on
32 visualizations for the main datasets D1–D5 which involved presentation of stimuli that vary in eccentricity.
33 This is because studies of the visual system provide well-established 'ground truth' expectations for
34 neural activity patterns elicited by stimuli varying in eccentricity: in brief, neurons in early visual cortex
35 respond selectively to stimuli at specific eccentricities, and the preferred eccentricity varies smoothly from
36 foveal to peripheral eccentricities along the posterior-to-anterior direction (Wandell and Winawer, 2011).

37
38 Inspecting results for a representative dataset (**Figure 6**; other datasets shown in **Figure 7**), we first
39 consider the overall magnitudes of the betas obtained with the different GLMs. The Standard betas
40 appear to be roughly the sum of the Early and Late betas. We find that the Early and Late betas are both
41 relatively large in magnitude (**Figure 6, first three columns**), indicating that both Early and Late
42 timecourses contribute substantially to measured BOLD responses. However, the Early and Late betas
43 show strikingly different patterns. The Early betas are relatively homogeneous across the cortical surface,
44 relatively flat across cortical depth, and are moderate in size at around 1–4% signal change. In contrast,
45 the Late betas are quite heterogeneous across the cortical surface (sparsely distributed), heavily biased
46 towards outer cortical depths, and are sometimes quite large in size, reaching 10% or more signal
47 change. The observation of additional activations arising in the Late betas is consistent with the fact that
48 the BOLD point-spread function appears larger when sampling late in the BOLD response (Shmuel et al.,
49 2007). Furthermore, comparing the spatial pattern of the Late betas against the spatial pattern of bias-
50 corrected EPI intensities (**Figure 6, fourth column**), we see a general correspondence between the
51 sparsely distributed locations where very large BOLD responses are observed and regions with dark EPI

1 intensity. This co-localization of Late betas and dark EPI intensity is consistent with the earlier
 2 observation that timecourse shapes resembling the Late timecourse originate from vertices with dark EPI
 3 intensity (see **Figure 4A, right column**). These several observations (sparse distribution, outer depth
 4 bias, large BOLD signal change, dark EPI intensity) strongly suggest that the Late betas reflect responses
 5 from the macrovasculature.
 6



7
 8
 9 **Figure 6. TDM decomposes brain activity patterns into early and late components.** Here we
 10 show detailed results for Dataset D1; summary results for all datasets are provided in **Figure 7**.
 11 Rows correspond to different cortical depths for a small patch covering primary visual cortex
 12 (flattened surface, left hemisphere). Three versions of results are shown: the first (Standard)
 13 reflects betas from a GLM incorporating a single canonical HRF, while the next two (Early, Late)
 14 reflect betas from a GLM incorporating the timecourses found by TDM. Four types of images are
 15 shown from left to right: (1) absolute value of betas averaged across conditions, (2) bias-corrected
 16 EPI intensities, (3) peak eccentricity quantified as the center-of-mass of the six betas observed at
 17 each vertex, and (4) single-condition activity patterns for the fifth stimulus. Blue, purple, and green
 18 arrows mark vertices illustrated in greater detail in the inset boxes. Each inset box shows FIR
 19 timecourses (upper left) with ribbon width indicating standard error across condition-splits;
 20 canonical and TDM-derived timecourses (lower left); and the three versions of the betas (right)
 21 with error bars indicating standard error across condition-splits and black arrows indicating peak
 22 eccentricity. Notice that at the level of individual vertices, TDM decomposes BOLD responses into
 23 two sets of betas (Early, Late) that exhibit different stimulus selectivity.

1
2 We next consider maps of peak eccentricity tuning, calculated as the center-of-mass of the betas
3 corresponding to different stimulus eccentricities (**Figure 6, fifth through seventh columns**). All three
4 versions of the betas (Standard, Early, Late) exhibit the expected smooth large-scale progression from
5 foveal (blue) to peripheral (red) eccentricities as one moves posterior (left) to anterior (right) in early visual
6 cortex. However, the quality or robustness of the eccentricity map is highest for the Standard betas,
7 moderately high for the Early betas, and relatively low for the Late betas. Moreover, for the Late betas,
8 there is a substantial decrease in quality moving from outer to inner cortical depths; this is consistent with
9 the sharp fall-off in the magnitude of betas moving from outer to inner depths, as seen previously (see
10 **Figure 6, third column**). We also observe that although large-scale eccentricity patterns are similar
11 across the three versions of the betas, the maps show divergence at a finer scale. In particular, there are
12 artifacts in eccentricity tuning that are present in the Standard and Late betas but absent in the Early
13 betas. We illustrate a few example vertices in further detail (**Figure 6, arrows and inset boxes**); in these
14 vertices, the tuning derived from the Early betas is better matched to expectations (based on inspection of
15 neighboring sections of cortex and the large-scale eccentricity map). This is an important outcome and
16 suggests that the Early betas isolate a component of the data that is more closely matched to responses
17 from the microvasculature and avoids unwanted responses from the macrovasculature.

18
19 Finally, for further clarification of these results, we examine activity patterns elicited by a single
20 experimental condition (**Figure 6, eighth through tenth columns**). Based on known tuning properties of
21 early visual cortex (Wandell and Winawer, 2015, 2011), we expect a relatively compact 'stripe' of positive
22 activity extending along the superior-inferior direction in cortex. All versions of the activity pattern
23 (Standard, Early, Late) indeed show evidence of a stripe. However, there are differences, recapitulating
24 effects observed earlier. Specifically, the Early beta is relatively homogeneous within the extent of the
25 stripe, relatively flat across cortical depth, and moderate in size within the stripe at around 4% signal
26 change. In contrast, the Late beta is heterogeneous across the cortical surface, heavily biased to outer
27 cortical depths, and often very large in size, exceeding 8% signal change. These results demonstrate that
28 TDM is capable of decomposing single-condition activity patterns into Early and Late components, with
29 the former component more closely matching ground-truth expectations. Note that the ground-truth
30 expectation is not necessarily a perfectly flat depth profile. This is because vascular density varies, to
31 some degree, across cortical depth (Schmid et al., 2019). Moreover, our 0.8-mm measurements sampled
32 within gray matter are certainly susceptible to partial volume effects from cerebrospinal fluid and white
33 matter.

34

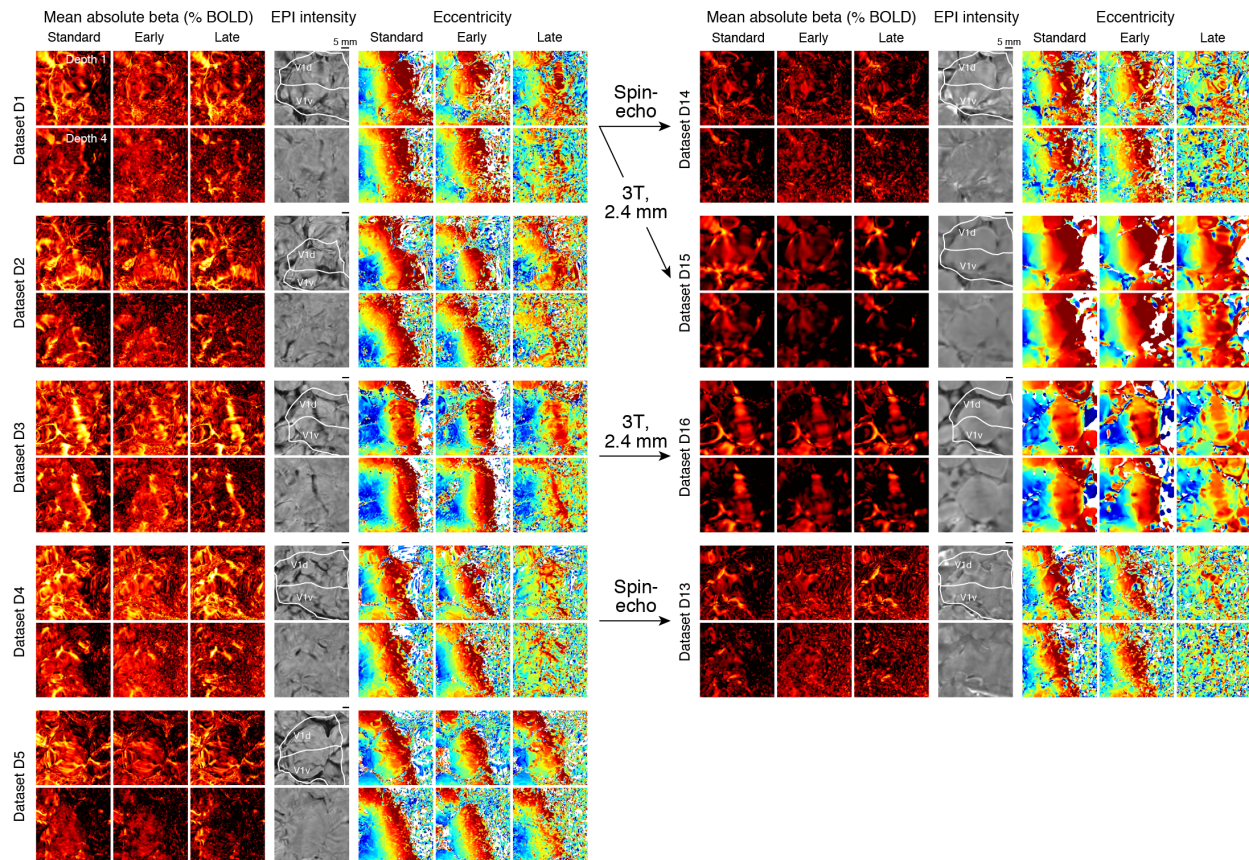


Figure 7. Decomposition of brain activity patterns across datasets and acquisition protocols. Same format as Figure 6, except only two cortical depths (Depths 1 and 4) are displayed. On the left are results obtained using high-resolution (0.8-mm) 7T gradient-echo (Datasets D1–D5). On the right are results obtained using high-resolution (1.05-mm) 7T spin-echo (Datasets D13–D14) and low-resolution (2.4-mm) 3T gradient-echo (Dataset D15–D16). These alternative acquisition protocols were conducted in the same subjects as the high-resolution gradient-echo protocol (correspondence indicated by arrows).

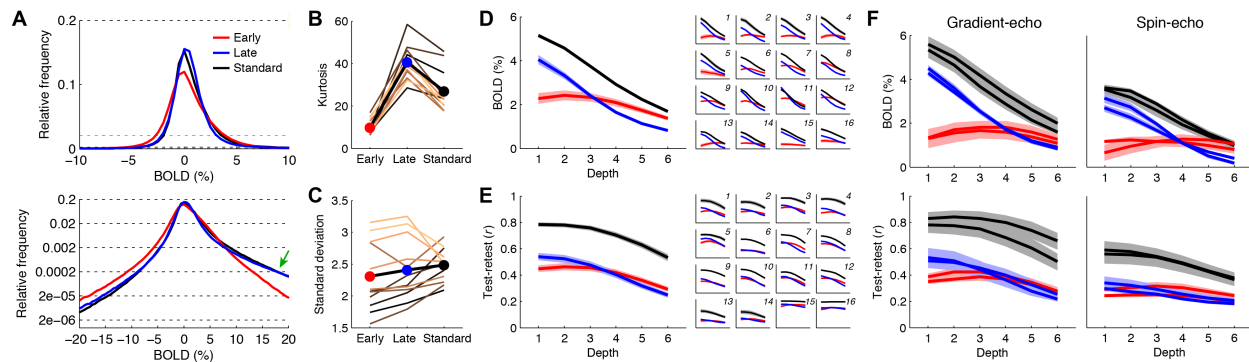
3.6. Quantification of TDM improvements

Although surface visualizations convey a large amount of information, they are qualitative. We therefore calculated quantitative metrics to more rigorously assess the results. First, we constructed histograms of the distributions of betas obtained under the different GLMs (**Figure 8A**). These distributions have long tails (**Figure 8A, bottom**) that appear to correspond to very large BOLD responses from the macrovasculature (see **Figure 6**). We quantified the magnitude of these tails using kurtosis, a metric that is large for heavy-tailed distributions. We find that Late betas have very high kurtosis, whereas Early betas have relatively low kurtosis (**Figure 8B**). Standard betas have an intermediate level of kurtosis, consistent with the interpretation that Standard betas behave essentially like an average or mixture of the Early and Late betas. In addition to quantifying the tails of the distributions, we quantified the overall magnitude of the betas by calculating the standard deviation of each distribution. We find that unlike kurtosis, the standard deviations of the beta distributions are similar across the three versions of the betas (**Figure 8C**). These results confirm that Early and Late betas are generally large in magnitude (indicating that BOLD responses are observed for both Early and Late timecourses), but Late betas have the unique feature of sometimes reaching extremely high values. This is consistent with the interpretation of Late betas as reflecting macrovascular responses.

1
 2 Next, we constructed depth profiles in order to understand how the magnitude of BOLD responses
 3 change with cortical depth. For these analyses, we restricted our quantification of BOLD responses to
 4 portions of cortex that have at least some substantive BOLD response for each given experimental
 5 condition. Specifically, for each condition, we selected vertices whose Standard beta exceeds 3% signal
 6 change at any cortical depth. We then averaged the BOLD response across these vertices at each depth
 7 (thus, a common set of vertices is used across depths). Consistent with earlier inspections of surface
 8 maps, we find that the magnitudes of both Standard betas and Late betas exhibit strong depth
 9 dependence (**Figure 8D**). For example, the BOLD signal change for Late betas is about 4 times greater
 10 at Depth 1 (outer) than at Depth 6 (inner). In contrast, the BOLD signal change for Early betas is fairly
 11 constant across depth. This effect is observed both in the group average (**Figure 8D, main plot**) as well
 12 as in each individual dataset (**Figure 8D, inset plots**). Thus, TDM successfully removes late
 13 components, which likely reflect unwanted macrovascular effects that bias BOLD responses towards
 14 outer cortical depths.

15
 16 Finally, an important aspect of betas is their reliability, i.e. robustness or consistency across repeated
 17 measurements. We observed earlier that the eccentricity maps vary substantially in quality (see **Figure 6**,
 18 **fifth through seventh columns**), suggesting that betas can suffer from low reliability. To quantify these
 19 observations, we calculated the correlation of betas across split-halves of each dataset (test-retest
 20 reliability). We find that the Standard betas are the most reliable, and there is a drop in reliability for the
 21 Early and Late betas (**Figure 8E**). The decrease in reliability for the Early and Late betas is unfortunate
 22 but not surprising: Early and Late timecourses are typically highly correlated, and the accuracy of
 23 regression estimates in the case of correlated regressors is expected to be somewhat degraded.

24



25

26

27 **Figure 8. Quantitative assessment of BOLD amplitude estimates provided by TDM.** A,
 28 Histogram of BOLD amplitudes. The top plot shows distributions of BOLD amplitudes aggregated
 29 across Datasets D1–D12; the bottom plot shows the same results on a log scale and with a wider
 30 x-axis range. B, Kurtosis of BOLD amplitudes. Results are shown for individual datasets (thin lines)
 31 and the group average (thick black line). C, Standard deviation of BOLD amplitudes. Same format
 32 as panel B. D, Cortical depth profiles of BOLD amplitudes. The main plot shows the average depth
 33 profile observed in Datasets D1–D12, with ribbons indicating standard error across datasets; the
 34 inset plots show results for individual datasets (D1–D16), with ribbons indicating standard error
 35 across conditions (same axis limits as the main plot). E, Test-retest reliability of BOLD amplitudes.
 36 Same format as panel D. F, Gradient-echo versus spin-echo. We re-plot results from panels D and
 37 E, directly comparing the two datasets acquired using gradient-echo against the two datasets
 38 acquired using spin-echo (conducted in the same subjects).

39

40 3.7. TDM is compatible with other acquisition methods

41

1 To gain further insight into the nature of TDM and its generalizability, we repeated the eccentricity
2 experiment using low-resolution fMRI (3T, 2.4-mm; Datasets D15–D16) as well as spin-echo fMRI (7T,
3 1.05-mm; Datasets D13–D14). The acquisition was performed in subjects who also participated in the
4 high-resolution gradient-echo acquisition, enabling direct comparison of results.

5

6 We find that the TDM method successfully applies to both acquisition styles. In the low-resolution data,
7 we see that the distribution of timecourse shapes tightens (**Supplementary Figure 3, fifth and seventh**
8 **rows**), which likely reflects a combination of averaging diverse timecourses within individual voxels and
9 reduction of thermal noise. Moreover, we find that TDM derives Early and Late timecourses that closely
10 resemble those found in the high-resolution data (**Supplementary Figure 3, fourth column**). This
11 implies that even at a resolution of 2.4 mm, there is sufficient diversity of timecourses to support data-
12 driven discovery of latent timecourses. Examining the surface visualizations (**Figure 7**), we observe
13 results for the low-resolution measurements that are consistent with the high-resolution measurements.
14 However, the differences between the spatial patterns of the Early and Late betas are reduced, indicating
15 that the differences between microvascular and macrovascular effects tend to wash out at low spatial
16 resolutions.

17

18 In the spin-echo data, we find patterns of results that look remarkably similar to the gradient-echo data.
19 Early and Late timecourses are identified (**Supplementary Figure 3, first and third rows**), and large
20 betas are found for Late timecourses (**Figure 7**). Standard GLM analysis of spin-echo data yields betas
21 that exhibit depth-dependent bias in BOLD signal change (**Figure 8F, upper right, black lines**), and this
22 bias is largely eliminated after applying TDM (**Figure 8F, upper right, red lines**). Importantly, the spin-
23 echo measurements suffer from a decrease in sensitivity compared to the gradient-echo measurements
24 (**Figure 8F, bottom plots, black lines**), even though larger voxels were used for the spin-echo
25 measurements. Overall, these results indicate that spin-echo measurements still contain substantial
26 contributions from the macrovasculature and that TDM is able to identify and remove the
27 macrovasculature-related effects.

28

4. Discussion

4.1. Nature of the TDM method

The core idea underlying TDM is that early and late timecourses can be derived from task-based fMRI data and reflect microvasculature-related and macrovascular-related responses, respectively. This idea is conceptually simple and draws its roots from critical observations made in prior studies (de Zwart et al., 2005; Lee et al., 1995; Siero et al., 2011). The value of the work presented here lies not so much in the discovery of a new phenomenon, but rather in the design and validation of new analysis methodology. Specifically, we have designed analysis procedures that are simple and robust and produce readily interpretable visualizations and results. Furthermore, we provide substantial empirical validation that the technique, when applied to fMRI measurements, produces sensible results: late components co-vary with dark EPI intensity (**Figures 6–7, Supplementary Figure 2**), depth-dependent response profiles are substantially flattened after removal of late components (**Figure 6, Figure 8D**), and artifacts in cortical maps are eliminated (**Figures 6–7**).

The TDM method has a few prerequisites:

- *Data acquisition.* The TDM method is likely compatible with a broad range of acquisition styles. As we have shown, TDM can be applied to data from standard spatial resolutions (2–3 mm; see **Figure 7, Supplementary Figure 3**) or data from high spatial resolutions (<1 mm; see **Figure 6**). Presumably, the spatial resolution needs only to be high enough to allow diverse sampling of vasculature in the brain. With respect to temporal resolution, we have shown that TDM can be successfully applied to data acquired at even fairly slow rates, such as the 2.2-s sampling rate used in Datasets D1–D14; this was likely aided by the fact that we jittered the acquired time points with respect to the experimental conditions (see **Figure 1B**).
- *Experimental design.* The TDM method requires a task-based experiment in which neural events occur at prescribed times, and is therefore inapplicable to resting-state paradigms. We suspect that TDM will be most effective for event-related paradigms where experimental events are somewhat short (e.g. 4 s or less). Block designs involving prolonged events (e.g. 16–30 s) or designs involving continuously changing experimental parameters (e.g. sinusoidal variation of a stimulus property) are likely to generate microvasculature-related and macrovascular-related timecourses that are more similar and therefore harder to disambiguate.
- *Amount of data.* Because TDM is a data-driven technique, it is necessary to acquire sufficient data to support the method. For example, there must be sufficient data to estimate response timecourses from the voxels in a dataset. If a given dataset is overly noisy or if not enough data are collected, timecourse estimates may be noisy and nearly isotropic in their distributions in the 3-dimensional PCA space (e.g. see Dataset D6 in **Supplementary Figure 2**), making it difficult to extract the underlying manifold structure of the data.

Given that these prerequisites are fairly minimal (i.e., event-related designs with reasonable fMRI acquisition parameters and reasonable amounts of data), we suspect that the TDM method may be widely applicable to different kinds of fMRI data.

Algorithmically, TDM uses a manifold-fitting method to characterize latent structure in timecourse variations. There are other methods that can characterize latent structure; two widely used methods are principal components analysis (PCA) and independent components analysis (ICA). Could these methods have been used instead? Due to the orthogonality constraint in PCA, it is necessarily the case that the PC timecourses returned by PCA are orthogonal (i.e. the dot product between each pair of timecourses equals zero). Although TDM does make use of PCA to determine the 3-dimensional space within which to perform further analyses, the PCA timecourses themselves do not constitute good candidates for latent timecourses. This is simply because there is no reason to expect hemodynamic timecourses in the brain

1 to be orthogonal. Indeed, the bulk of empirically measured timecourses tend to reside in a small portion of
2 the 3-dimensional space, and the early and late timecourses returned by TDM are nearby in this space
3 and highly correlated (see **Figure 4** and **Supplementary Figure 2**).

4
5 In contrast to PCA, ICA does not impose the constraint of orthogonality. Instead, ICA optimizes
6 timecourses with respect to statistical independence, often through some measure of non-Gaussianity
7 (e.g. kurtosis). We have demonstrated that it is possible to construct an ICA-based procedure that can
8 potentially derive early and late timecourses. Though the derived timecourses from the ICA-based
9 procedure are sometimes similar to those produced by TDM (see x's in **Supplementary Figure 2**), there
10 are clear advantages of the TDM method. First, the data visualization and explicit modeling performed by
11 TDM allow the user to evaluate and confirm the data features that give rise to the derived timecourses.
12 ICA, without further analysis, remains a 'black box' and it is difficult to understand the specific features of
13 the data that give rise to its results. Second, there is no *a priori* reason to think that loadings on early and
14 late hemodynamic timecourses must necessarily conform to statistical independence. Thus, relying on a
15 procedure that is predicated on independence seems risky. Third, ICA alone does not identify the early
16 and late timecourses; rather, we found it necessary to couple the results of ICA with several post-hoc
17 procedures that are somewhat heuristic in nature and thus unsatisfying (see Section 2.8.4). On the whole,
18 we suggest that TDM is more explicit, more direct, and more interpretable than ICA. Indeed, historically,
19 the first method that we developed was the ICA-based procedure, and the shortcomings described above
20 are what prompted us to develop the TDM method.

21
22 GLM-based analyses of fMRI time-series data sometimes allow flexible modeling of timecourse shape
23 through the inclusion of a canonical hemodynamic response timecourse and its temporal derivative
24 (Friston et al., 1998a) or some other basis function decomposition such as PCA (Woolrich et al., 2004).
25 While TDM shares the common feature of providing a means to capture timecourse variation, the key
26 difference with respect to these alternative approaches lies in the specific timecourses that are chosen by
27 TDM. The Early and Late timecourses found by TDM are often quite correlated (unlike a timecourse and
28 its derivative or those returned by PCA). Moreover, the Early and Late timecourses have specific
29 biological meanings, and so the beta loadings found for these timecourses have specific value. It is
30 possible that alternative timecourse models can yield fits to a set of data that are as good as the fit
31 achieved by TDM, but the beta loadings associated with these models cannot be interpreted in terms of
32 microvasculature- and macrovasculature-related components.

33 34 4.2. Validation of the TDM method

35
36 In this study, we demonstrated that TDM delivers robust and meaningful results in each of the 16 fMRI
37 datasets collected (11 unique subjects). The main lines of validation include sensible timecourses (the
38 shapes of the Early and Late timecourses are plausible and consistent with simple inspections of
39 response timecourses; see **Figure 2**, **Supplementary Figure 1**), co-variation of loadings on the Late
40 timecourse with dark EPI intensities (see **Figure 6**) and with kurtotic BOLD amplitudes (see **Figure 8B**),
41 flattening of depth-dependent response profiles (see **Figure 8D**), and elimination of artifacts in cortical
42 maps for which we have ground-truth expectations (see **Figures 6–7**). Furthermore, we make available
43 data and analysis code to ensure that the TDM method is transparent and reproducible (Poline and
44 Poldrack, 2012).

45
46 While we believe a solid case has been made for TDM, additional validation would nonetheless be useful.
47 Further work could be directed at assessing and optimizing the technique with respect to experimental
48 design characteristics such as the duration of experimental conditions, the spatial and temporal resolution
49 of the acquisition, and the amount of data acquired. In addition, it would be worthwhile to test the
50 technique on other types of experiments (other sensory, cognitive, and/or motor experiments) and in
51 other brain areas. In particular, it would be interesting to assess how well TDM can resolve fine-scale

1 variation in neural representations, such as ocular dominance columns (Cheng et al., 2001). Since TDM
2 makes no restrictions on the spatial loadings of the Early and Late timecourses, the technique should in
3 principle be applicable not only to large-scale neural representations like eccentricity but also fine-scale
4 representations like ocular dominance. Finally, it would be quite valuable to validate results obtained
5 using TDM against direct neural activity measurements, such as electrophysiological recordings with
6 laminar resolution (Maier et al., 2010; Self et al., 2019).

7 8 4.3. Strengths and limitations of the TDM method

9
10 All methods have strengths and limitations. To make an informed decision of whether to pursue a given
11 method, it is important to understand its strengths and limitations.

12
13 Strengths of the TDM method include the following:

- 14 • *Flexibility*. As discussed in Section 4.1, TDM has fairly minimal prerequisites. It can be applied to
15 different types of acquisitions: gradient-echo, spin-echo, high-resolution, low-resolution, etc.
16 Moreover, as a data-driven method, TDM makes no assumptions about shapes of hemodynamic
17 timecourses that might be found in a dataset and naturally adapts to different brain areas,
18 subjects, and/or datasets. Furthermore, since the response decomposition is performed
19 independently for each voxel, the technique makes no assumptions regarding the spatial
20 distribution of BOLD responses in a given experiment. In other words, the loadings on the Early
21 and Late timecourses can vary from voxel to voxel, and the technique can, in theory, capture
22 these variations. Finally, TDM applies to single-condition activity maps, and therefore avoids the
23 assumptions required in differential paradigms where unwanted non-specific effects are assumed
24 to be removed through subtraction.
- 25 • *Transparency*. An integral aspect of TDM is direct visualization of distributions of timecourses.
26 Thus, it is easy for the user to understand the nature of the data and whether the derived
27 timecourses are meaningful. This stands in contrast to ‘black box’ methods that might produce
28 unusual results without explanation.
- 29 • *Simplicity and robustness*. TDM is simple in its design and does not, as far as we have seen,
30 require fine-tuning of parameters to be effective. Our results establish more than just a proof of
31 principle: we show that the method without any modification performs robustly across different
32 datasets, subjects, and experiments.
- 33 • *Analysis not acquisition*. Since TDM is an analysis method, it can be retrospectively applied to
34 datasets that have already been acquired. Moreover, since TDM does not place major constraints
35 on acquisition, the user is not burdened with making difficult decisions regarding optimal
36 acquisition parameters (e.g., choosing between a standard acquisition scheme that is guaranteed
37 to produce reasonably strong signals versus a specialized acquisition scheme that might suffer
38 from low sensitivity).

39
40 Limitations of the TDM method include the following:

- 41 • *Sensitivity loss*. TDM involves decomposing fMRI responses using two timecourses that are often
42 highly correlated. Thus, from a regression perspective, one expects to incur a penalty in terms of
43 high variance in beta estimates. One can view the use of TDM as approximately doubling the
44 number of experimental conditions while keeping the number of data points constant. As such,
45 the ensuing model will be more difficult to reliably estimate compared to a basic GLM model with
46 a single set of experimental conditions. Thus, if sensitivity (i.e. reliability of beta estimates) is the
47 sole priority and specificity (i.e. accurate estimates of local neural activity) is not critical, the TDM
48 method is not recommended. If, on the other hand, specificity is of utmost importance, TDM is
49 likely to be a valuable method. In short, TDM does not deliver more robust fMRI maps (e.g. large

1 blobs of statistically significant activations), but aims to deliver more spatially accurate and
2 neurally meaningful maps.

- 3 • *Intrinsic physiological limitations.* The TDM method attempts to disambiguate BOLD contributions
4 from the microvasculature and macrovasculature based on their respective associated
5 timecourses. The more similar these timecourses, the more difficult it will be to estimate the
6 distinct contributions of the timecourses. Thus, the intrinsic physiology of the subject places limits
7 on the overall effectiveness of the TDM method. For example, in our data (see **Supplementary**
8 **Figure 2**), we find that Subject S5 exhibits Early and Late timecourses that are widely separated
9 in the 3-dimensional PCA space and are highly distinct, whereas Subject S4 exhibits Early and
10 Late timecourses that are close together in the 3-dimensional PCA space and are fairly similar.
11 This may be the reason why some datasets (such as Subject S4's Dataset D4) experience a
12 larger reduction in reliability when using TDM compared to other datasets (see **Figure 8E**). One
13 possible approach to achieve optimal results is to screen subjects according to the temporal
14 separability of their microvasculature- and macrovasculature-related timecourses.
- 15 • *Complexity of the vasculature.* TDM proposes a simple two-component model to decompose
16 BOLD timecourses. The vasculature is certainly more complex than this simple characterization
17 (Uludağ and Blinder, 2018). Thus, it may be fruitful to develop a more nuanced characterization of
18 vasculature types and their dynamics.

20 4.4. Comparison to other methods

21
22 How does TDM compare to other approaches for removing or avoiding venous effects in fMRI? Some
23 researchers have proposed simple heuristic selection methods. For example, sampling BOLD responses
24 at only deep cortical depths (Polimeni et al., 2010) can help avoid the influence of large draining veins
25 near the pial surface. However, this comes at the cost of not being able to infer response properties in
26 superficial cortical depths; moreover, it is still possible for veins to penetrate deep into cortex (see **Figure**
27 **6**). Another example is masking out voxels with very high percent BOLD signal change (e.g. Shmuel et
28 al., 2007), low EPI intensity (e.g. Olman et al., 2007), and/or low temporal signal-to-noise ratio (e.g.
29 Fracasso et al., 2018). While these are reasonable heuristics for removing voxels that are most
30 egregiously affected by large veins, it is not clear what principle can be used to set the threshold to be
31 used. Moreover, similar to the approach of sampling only deep cortical depths, this approach fails to
32 recover usable signals from the removed voxels. Finally, one suggestion found in older work (Goodyear
33 and Menon, 2001; Shmuel et al., 2007) and more recent work (Blazejewska, Nasr, Polimeni, ISMRM
34 2018 abstract) is to sample early time points in the BOLD response. This is certainly consistent with the
35 spirit of TDM and may produce a response snapshot that is more weighted towards the microvasculature.
36 However, our results indicate that Early and Late timecourses are highly overlapping (see **Figure 5**). If
37 the chosen time point is not sufficiently early, this incurs the risk of the late component “bleeding” into the
38 analysis results. Moreover, choosing only one or a few time points does not make efficient use of all of
39 the available data. Compared to these various heuristic selection methods, we believe that TDM has
40 substantial appeal: *TDM can recover signals at all depths and even in voxels that have substantial*
41 *venous influence; it makes efficient use of all of the fMRI data collected; and, as a data-driven technique,*
42 *it naturally identifies appropriate timecourse parameters for each dataset.*

43
44 Recently, a method has been proposed that first constructs a forward model characterizing the mixing of
45 hemodynamic signals from different cortical depths due to blood drainage towards the pial surface and
46 then uses this model to invert observed BOLD depth profiles (Heinzle et al., 2016; Markuerkiaga et al.,
47 2016; Marquardt et al., 2018). This method bears a parallel to TDM in the sense that it is a spatial
48 deconvolution approach whereas TDM is a temporal decomposition approach. However, the accuracy of
49 the spatial deconvolution approach may be dependent on the correctness of the model parameters, which
50 might vary across brain regions and/or subjects. In addition, the approach deals only with vascular effects

1 that vary across depths. In contrast, TDM is a data-driven technique that adapts to each given dataset
2 and compensates for vascular effects present at every voxel.

3

4 Besides analysis methods, one can consider using acquisition methods to avoid venous effects.
5 Switching from conventional gradient-echo pulse sequences to spin-echo pulse sequences (or related
6 techniques such as GRASE (De Martino et al., 2013; Moerel et al., 2018; Olman et al., 2012)) has
7 traditionally been considered the standard approach for mitigating venous effects in fMRI (Yacoub et al.,
8 2008). While the refocusing of T2* effects by the 180° RF pulse in spin-echo eliminates sensitivity to
9 extravascular effects around large veins, it is important to note this holds only for a specific point in time
10 (typically the center of the readout window). The remainder of the image acquisition incurs T2* effects
11 (Goense and Logothetis, 2006). Furthermore, spin-echo does not eliminate intravascular effects in large
12 vessels (Budde et al., 2014; Duong et al., 2003). Thus, spin-echo does not provide full elimination of
13 venous effects. In addition, spin-echo acquisitions suffer from increased energy deposition, limits on
14 spatial coverage, lower temporal resolution, and loss of signal-to-noise ratio.

15

16 In this study, we have provided a direct comparison of TDM and spin-echo. In order to maintain
17 sensitivity, we acquired spin-echo data at a lower resolution (1.05-mm vs. 0.8-mm) but maintained the
18 same TR and the same overall experiment duration as the gradient-echo data. Our results show that the
19 spin-echo data analyzed using a standard GLM (single canonical hemodynamic timecourse) is more
20 robust than gradient-echo data decomposed using TDM (see **Figure 8F**). One reason for the increased
21 robustness of the spin-echo data is its lower spatial resolution, providing the data with some advantage
22 over the gradient-echo data. However, even if the two types of data were matched in resolution, spin-
23 echo should not be viewed as a complete solution since it does not fully suppress venous effects. Indeed,
24 we demonstrate that TDM can be applied to the spin-echo data in order to remove venous influences
25 present in those data (see **Figure 7**). When comparing gradient-echo data and spin-echo data that have
26 both been decomposed using TDM, gradient-echo has greater robustness (see **Figure 8F**). Thus, we
27 suggest that if removal of venous effects is top priority and one plans to use TDM, there is little benefit to
28 spin-echo acquisition over conventional gradient-echo acquisition.

29

30 A promising alternative to spin-echo is vascular space occupancy (VASO), a pulse sequence that is not
31 sensitive to the BOLD effect but rather to changes in cerebral blood volume (Lu et al., 2003). This
32 approach avoids the impurities that linger in spin-echo acquisitions, and has been shown to generate
33 highly specific measures of task-driven hemodynamic responses (Huber et al., 2017). A promising
34 direction for future work is to perform direct comparison of an gradient-echo acquisition optimized for use
35 with TDM against an optimized VASO acquisition. Through rigorous evaluations, we hope that the field
36 will eventually reach consensus with respect to the most effective approach for achieving accurate fine-
37 scale measurements of brain activity.

38

1 **5. Author Contributions**

2

3 K.K. designed the experiment. R.Z. collected the data. K.K. and K.J. developed techniques and analyzed
4 the data. K.K. wrote the paper. All authors discussed and edited the manuscript.

5

6 **6. Acknowledgements**

7

8 We thank E. Margalit and N. Petridou for helpful discussions. This work was supported by NIH Grants
9 P41 EB015894, P41 EB027061, P30 NS076408, S10 RR026783, S10 OD017974-01, U01 EB025144,
10 and the W. M. Keck Foundation.

11

12 **7. Competing Interests**

13

14 The authors confirm that there are no competing interests.

15

References

- 1
2
3 Attwell, D., Iadecola, C., 2002. The neural basis of functional brain imaging signals. *Trends Neurosci.* 25,
4 621–625. [https://doi.org/10.1016/s0166-2236\(02\)02264-6](https://doi.org/10.1016/s0166-2236(02)02264-6)
- 5 Bianciardi, M., Fukunaga, M., van Gelderen, P., de Zwart, J.A., Duyn, J.H., 2011. Negative BOLD-fMRI
6 signals in large cerebral veins. *J. Cereb. Blood Flow Metab.* 31, 401–412.
7 <https://doi.org/10.1038/jcbfm.2010.164>
- 8 Brainard, D.H., 1997. The Psychophysics Toolbox. *Spat Vis* 10, 433–436.
- 9 Budde, J., Shajan, G., Zaitsev, M., Scheffler, K., Pohmann, R., 2014. Functional MRI in human subjects
10 with gradient-echo and spin-echo EPI at 9.4 T. *Magn Reson Med* 71, 209–218.
11 <https://doi.org/10.1002/mrm.24656>
- 12 Charest, I., Kriegeskorte, N., Kay, K.N., 2018. GLMdenoise improves multivariate pattern analysis of fMRI
13 data. *NeuroImage* 183, 606–616. <https://doi.org/10.1016/j.neuroimage.2018.08.064>
- 14 Cheng, K., 2018. Exploration of human visual cortex using high spatial resolution functional magnetic
15 resonance imaging. *NeuroImage* 164, 4–9. <https://doi.org/10.1016/j.neuroimage.2016.11.018>
- 16 Cheng, K., Waggoner, R.A., Tanaka, K., 2001. Human ocular dominance columns as revealed by high-
17 field functional magnetic resonance imaging. *Neuron* 32, 359–374.
- 18 d’Avossa, G., Shulman, G.L., Corbetta, M., 2003. Identification of cerebral networks by classification of
19 the shape of BOLD responses. *J. Neurophysiol.* 90, 360–371.
20 <https://doi.org/10.1152/jn.01040.2002>
- 21 Dale, A.M., 1999. Optimal experimental design for event-related fMRI. *Hum Brain Mapp* 8, 109–114.
- 22 De Martino, F., Yacoub, E., Kemper, V., Moerel, M., Uludağ, K., De Weerd, P., Ugurbil, K., Goebel, R.,
23 Formisano, E., 2018. The impact of ultra-high field MRI on cognitive and computational
24 neuroimaging. *NeuroImage* 168, 366–382. <https://doi.org/10.1016/j.neuroimage.2017.03.060>
- 25 De Martino, F., Zimmermann, J., Muckli, L., Ugurbil, K., Yacoub, E., Goebel, R., 2013. Cortical depth
26 dependent functional responses in humans at 7T: improved specificity with 3D GRASE. *PLoS*
27 *ONE* 8, e60514. <https://doi.org/10.1371/journal.pone.0060514>
- 28 de Zwart, J.A., Silva, A.C., van Gelderen, P., Kellman, P., Fukunaga, M., Chu, R., Koretsky, A.P., Frank,
29 J.A., Duyn, J.H., 2005. Temporal dynamics of the BOLD fMRI impulse response. *NeuroImage* 24,
30 667–677.
- 31 Dumoulin, S.O., Fracasso, A., van der Zwaag, W., Siero, J.C.W., Petridou, N., 2018. Ultra-high field MRI:
32 Advancing systems neuroscience towards mesoscopic human brain function. *NeuroImage* 168,
33 345–357. <https://doi.org/10.1016/j.neuroimage.2017.01.028>
- 34 Duong, T.Q., Yacoub, E., Adriansy, G., Hu, X., Ugurbil, K., Kim, S.-G., 2003. Microvascular BOLD
35 contribution at 4 and 7 T in the human brain: gradient-echo and spin-echo fMRI with suppression
36 of blood effects. *Magn Reson Med* 49, 1019–1027. <https://doi.org/10.1002/mrm.10472>
- 37 Fischl, B., 2012. FreeSurfer. *NeuroImage* 62, 774–781. <https://doi.org/10.1016/j.neuroimage.2012.01.021>
- 38 Fracasso, A., Luijten, P.R., Dumoulin, S.O., Petridou, N., 2018. Laminar imaging of positive and negative
39 BOLD in human visual cortex at 7T. *NeuroImage* 164, 100–111.
40 <https://doi.org/10.1016/j.neuroimage.2017.02.038>
- 41 Friston, K.J., Fletcher, P., Josephs, O., Holmes, A., Rugg, M.D., Turner, R., 1998a. Event-related fMRI:
42 characterizing differential responses. *Neuroimage* 7, 30–40.
43 <https://doi.org/10.1006/nimg.1997.0306>
- 44 Friston, K.J., Josephs, O., Rees, G., Turner, R., 1998b. Nonlinear event-related responses in fMRI. *Magn*
45 *Reson Med* 39, 41–52. <https://doi.org/10.1002/mrm.1910390109>
- 46 Goense, J.B.M., Logothetis, N.K., 2006. Laminar specificity in monkey V1 using high-resolution SE-fMRI.
47 *Magnetic resonance imaging* 24, 381–392. <https://doi.org/10.1016/j.mri.2005.12.032>
- 48 Goodyear, B.G., Menon, R.S., 2001. Brief visual stimulation allows mapping of ocular dominance in visual
49 cortex using fMRI. *Hum Brain Mapp* 14, 210–217. <https://doi.org/10.1002/hbm.1053>

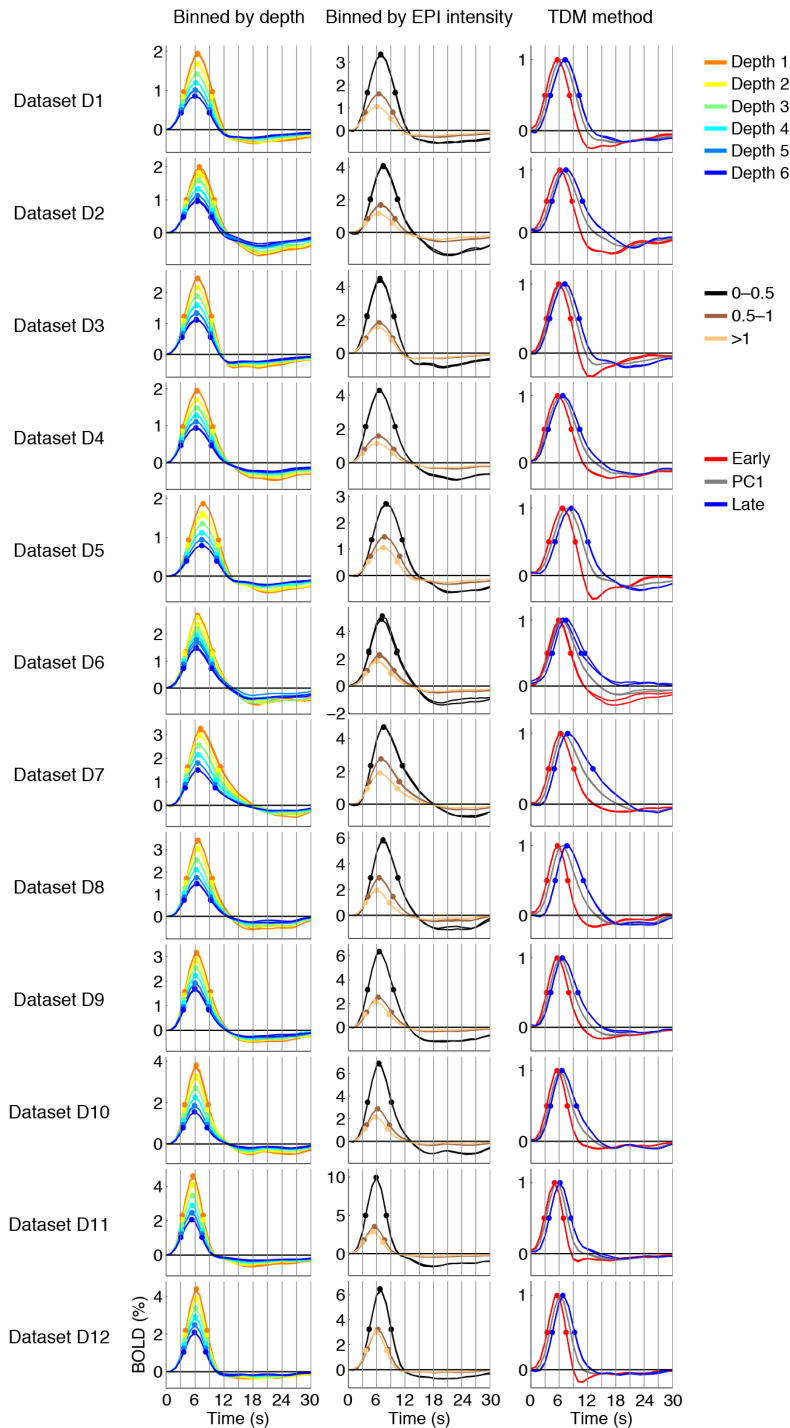
- 1 Haacke, E.M., Hopkins, A., Lai, S., Buckley, P., Friedman, L., Meltzer, H., Hedera, P., Friedland, R.,
2 Klein, S., Thompson, L., 1994. 2D and 3D high resolution gradient echo functional imaging of the
3 brain: venous contributions to signal in motor cortex studies. *NMR in biomedicine* 7, 54–62.
- 4 Handwerker, D.A., Gonzalez-Castillo, J., D'Esposito, M., Bandettini, P.A., 2012. The continuing challenge
5 of understanding and modeling hemodynamic variation in fMRI. *NeuroImage* 62, 1017–1023.
6 <https://doi.org/10.1016/j.neuroimage.2012.02.015>
- 7 Hansen, K.A., Kay, K.N., Gallant, J.L., 2007. Topographic organization in and near human visual area V4.
8 *J. Neurosci.* 27, 11896–11911.
- 9 Havlicek, M., Uludağ, K., 2020. A dynamical model of the laminar BOLD response. *Neuroimage* 204,
10 116209. <https://doi.org/10.1016/j.neuroimage.2019.116209>
- 11 Heinzle, J., Koopmans, P.J., den Ouden, H.E.M., Raman, S., Stephan, K.E., 2016. A hemodynamic
12 model for layered BOLD signals. *NeuroImage* 125, 556–570.
13 <https://doi.org/10.1016/j.neuroimage.2015.10.025>
- 14 Huber, L., Handwerker, D.A., Jangraw, D.C., Chen, G., Hall, A., Stüber, C., Gonzalez-Castillo, J., Ivanov,
15 D., Marrett, S., Guidi, M., Goense, J., Poser, B.A., Bandettini, P.A., 2017. High-Resolution CBV-
16 fMRI Allows Mapping of Laminar Activity and Connectivity of Cortical Input and Output in Human
17 M1. *Neuron* 96, 1253-1263.e7. <https://doi.org/10.1016/j.neuron.2017.11.005>
- 18 Kay, K., Jamison, K.W., Vizioli, L., Zhang, R., Margalit, E., Ugurbil, K., 2019. A critical assessment of data
19 quality and venous effects in sub-millimeter fMRI. *NeuroImage* 189, 847–869.
20 <https://doi.org/10.1016/j.neuroimage.2019.02.006>
- 21 Kay, K.N., David, S.V., Prenger, R.J., Hansen, K.A., Gallant, J.L., 2008a. Modeling low-frequency
22 fluctuation and hemodynamic response timecourse in event-related fMRI. *Hum Brain Mapp* 29,
23 142–156. <https://doi.org/10.1002/hbm.20379>
- 24 Kay, K.N., Naselaris, T., Prenger, R.J., Gallant, J.L., 2008b. Identifying natural images from human brain
25 activity. *Nature* 452, 352–355. <https://doi.org/10.1038/nature06713>
- 26 Kay, K.N., Rokem, A., Winawer, J., Dougherty, R.F., Wandell, B., 2013. GLMdenoise: a fast, automated
27 technique for denoising task-based fMRI data. *Front Neurosci* 7, 247.
28 <https://doi.org/10.3389/fnins.2013.00247>
- 29 Kim, J.H., Ress, D., 2017. Reliability of the depth-dependent high-resolution BOLD hemodynamic
30 response in human visual cortex and vicinity. *Magn Reson Imaging* 39, 53–63.
31 <https://doi.org/10.1016/j.mri.2017.01.019>
- 32 Lawrence, S.J.D., Formisano, E., Muckli, L., de Lange, F.P., 2017. Laminar fMRI: Applications for
33 cognitive neuroscience. *NeuroImage*. <https://doi.org/10.1016/j.neuroimage.2017.07.004>
- 34 Lee, A.T., Glover, G.H., Meyer, C.H., 1995. Discrimination of large venous vessels in time-course spiral
35 blood-oxygen-level-dependent magnetic-resonance functional neuroimaging. *Magn Reson Med*
36 33, 745–754.
- 37 Lu, H., Golay, X., Pekar, J.J., Van Zijl, P.C.M., 2003. Functional magnetic resonance imaging based on
38 changes in vascular space occupancy. *Magn Reson Med* 50, 263–274.
39 <https://doi.org/10.1002/mrm.10519>
- 40 Maier, A., Adams, G.K., Aura, C., Leopold, D.A., 2010. Distinct superficial and deep laminar domains of
41 activity in the visual cortex during rest and stimulation. *Front Syst Neurosci* 4.
42 <https://doi.org/10.3389/fnsys.2010.00031>
- 43 Markuerkiaga, I., Barth, M., Norris, D.G., 2016. A cortical vascular model for examining the specificity of
44 the laminar BOLD signal. *NeuroImage* 132, 491–498.
45 <https://doi.org/10.1016/j.neuroimage.2016.02.073>
- 46 Marquardt, I., Schneider, M., Gulban, O.F., Ivanov, D., Uludağ, K., 2018. Cortical depth profiles of
47 luminance contrast responses in human V1 and V2 using 7 T fMRI. *Hum Brain Mapp* 464, 1155.
48 <https://doi.org/10.1002/hbm.24042>
- 49 Menon, R.S., Ogawa, S., Tank, D.W., Ugurbil, K., 1993. Tesla gradient recalled echo characteristics of
50 photic stimulation-induced signal changes in the human primary visual cortex. *Magn Reson Med*
51 30, 380–386.

- 1 Moerel, M., De Martino, F., Kemper, V.G., Schmitter, S., Vu, A.T., Ugurbil, K., Formisano, E., Yacoub, E.,
2 2018. Sensitivity and specificity considerations for fMRI encoding, decoding, and mapping of
3 auditory cortex at ultra-high field. *NeuroImage* 164, 18–31.
4 <https://doi.org/10.1016/j.neuroimage.2017.03.063>
- 5 Ogawa, S., Lee, T.M., Nayak, A.S., Glynn, P., 1990. Oxygenation-sensitive contrast in magnetic
6 resonance image of rodent brain at high magnetic fields. *Magn Reson Med* 14, 68–78.
- 7 Olman, C.A., Harel, N., Feinberg, D.A., He, S., Zhang, P., Ugurbil, K., Yacoub, E., 2012. Layer-specific
8 fMRI reflects different neuronal computations at different depths in human V1. *PLoS ONE* 7,
9 e32536. <https://doi.org/10.1371/journal.pone.0032536>
- 10 Olman, C.A., Inati, S., Heeger, D.J., 2007. The effect of large veins on spatial localization with GE BOLD
11 at 3 T: Displacement, not blurring. *NeuroImage* 34, 1126–1135.
- 12 Parkes, L.M., Schwarzbach, J.V., Bouts, A.A., Deckers, R.H.R., Pullens, P., Kerskens, C.M., Norris, D.G.,
13 2005. Quantifying the spatial resolution of the gradient echo and spin echo BOLD response at 3
14 Tesla. *Magn Reson Med* 54, 1465–1472. <https://doi.org/10.1002/mrm.20712>
- 15 Pedregosa, F., Eickenberg, M., Ciuciu, P., Thirion, B., Gramfort, A., 2015. Data-driven HRF estimation for
16 encoding and decoding models. *NeuroImage* 104, 209–220.
17 <https://doi.org/10.1016/j.neuroimage.2014.09.060>
- 18 Pelli, D.G., 1997. The VideoToolbox software for visual psychophysics: transforming numbers into
19 movies. *Spat Vis* 10, 437–442.
- 20 Polimeni, J.R., Fischl, B., Greve, D.N., Wald, L.L., 2010. Laminar analysis of 7T BOLD using an imposed
21 spatial activation pattern in human V1. *NeuroImage* 52, 1334–1346.
22 <https://doi.org/10.1016/j.neuroimage.2010.05.005>
- 23 Poline, J.B., Poldrack, R.A., 2012. Frontiers in brain imaging methods grand challenge. *Front Neurosci* 6,
24 96. <https://doi.org/10.3389/fnins.2012.00096>
- 25 Schmid, F., Barrett, M.J.P., Jenny, P., Weber, B., 2019. Vascular density and distribution in neocortex.
26 *NeuroImage* 197, 792–805. <https://doi.org/10.1016/j.neuroimage.2017.06.046>
- 27 Self, M.W., van Kerkoerle, T., Goebel, R., Roelfsema, P.R., 2019. Benchmarking laminar fMRI: Neuronal
28 spiking and synaptic activity during top-down and bottom-up processing in the different layers of
29 cortex. *NeuroImage* 197, 806–817. <https://doi.org/10.1016/j.neuroimage.2017.06.045>
- 30 Shmuel, A., Yacoub, E., Chaimow, D., Logothetis, N.K., Ugurbil, K., 2007. Spatio-temporal point-spread
31 function of fMRI signal in human gray matter at 7 Tesla. *NeuroImage* 35, 539–552.
- 32 Siero, J.C.W., Petridou, N., Hoogduin, H., Luijten, P.R., Ramsey, N.F., 2011. Cortical depth-dependent
33 temporal dynamics of the BOLD response in the human brain. *J. Cereb. Blood Flow Metab.* 31,
34 1999–2008. <https://doi.org/10.1038/jcbfm.2011.57>
- 35 Stigliani, A., Weiner, K.S., Grill-Spector, K., 2015. Temporal Processing Capacity in High-Level Visual
36 Cortex Is Domain Specific. *J. Neurosci.* 35, 12412–12424.
37 <https://doi.org/10.1523/JNEUROSCI.4822-14.2015>
- 38 Thompson, S.K., Engel, S.A., Olman, C.A., 2014. Larger neural responses produce BOLD signals that
39 begin earlier in time. *Front Neurosci* 8, 159. <https://doi.org/10.3389/fnins.2014.00159>
- 40 Turner, R., 2002. How much cortex can a vein drain? Downstream dilution of activation-related cerebral
41 blood oxygenation changes. *NeuroImage* 16, 1062–1067.
- 42 Ugurbil, K., 2016. What is feasible with imaging human brain function and connectivity using functional
43 magnetic resonance imaging. *Philosophical transactions of the Royal Society of London* 371,
44 20150361. <https://doi.org/10.1098/rstb.2015.0361>
- 45 Ugurbil, K., Xu, J., Auerbach, E.J., Moeller, S., Vu, A.T., Duarte-Carvajalino, J.M., Lenglet, C., Wu, X.,
46 Schmitter, S., Van de Moortele, P.F., Strupp, J., Sapiro, G., De Martino, F., Wang, D., Harel, N.,
47 Garwood, M., Chen, L., Feinberg, D.A., Smith, S.M., Miller, K.L., Sotiropoulos, S.N., Jbabdi, S.,
48 Andersson, J.L.R., Behrens, T.E.J., Glasser, M.F., Van Essen, D.C., Yacoub, E., WU-Minn HCP
49 Consortium, 2013. Pushing spatial and temporal resolution for functional and diffusion MRI in the
50 Human Connectome Project. *NeuroImage* 80, 80–104.
51 <https://doi.org/10.1016/j.neuroimage.2013.05.012>

- 1 Uludağ, K., Blinder, P., 2018. Linking brain vascular physiology to hemodynamic response in ultra-high
2 field MRI. *NeuroImage* 168, 279–295. <https://doi.org/10.1016/j.neuroimage.2017.02.063>
- 3 Uludağ, K., Müller-Bierl, B., Uğurbil, K., 2009. An integrative model for neuronal activity-induced signal
4 changes for gradient and spin echo functional imaging. *NeuroImage* 48, 150–165.
5 <https://doi.org/10.1016/j.neuroimage.2009.05.051>
- 6 Wandell, B., Winawer, J., 2015. Computational neuroimaging and population receptive fields. *Trends in*
7 *cognitive sciences* 19, 349–357. <https://doi.org/10.1016/j.tics.2015.03.009>
- 8 Wandell, B., Winawer, J., 2011. Imaging retinotopic maps in the human brain. *Vision research* 51, 718–
9 737.
- 10 Wang, L., Mruczek, R.E.B., Arcaro, M.J., Kastner, S., 2015. Probabilistic Maps of Visual Topography in
11 Human Cortex. *Cereb. Cortex* 25, 3911–3931. <https://doi.org/10.1093/cercor/bhu277>
- 12 Woolrich, M.W., Behrens, T.E.J., Smith, S.M., 2004. Constrained linear basis sets for HRF modelling
13 using Variational Bayes. *NeuroImage* 21, 1748–1761.
14 <https://doi.org/10.1016/j.neuroimage.2003.12.024>
- 15 Yacoub, E., Harel, N., Ugurbil, K., 2008. High-field fMRI unveils orientation columns in humans.
16 *Proceedings of the National Academy of Sciences of the United States of America* 105, 10607–
17 10612.
- 18 Yacoub, E., Wald, L.L., 2018. Pushing the spatio-temporal limits of MRI and fMRI. *NeuroImage* 164, 1–3.
19 <https://doi.org/10.1016/j.neuroimage.2017.11.034>
- 20 Zhang, N., Yacoub, E., Zhu, X.-H., Ugurbil, K., Chen, W., 2009. Linearity of blood-oxygenation-level
21 dependent signal at microvasculature. *NeuroImage* 48, 313–318.
22 <https://doi.org/10.1016/j.neuroimage.2009.06.071>
- 23

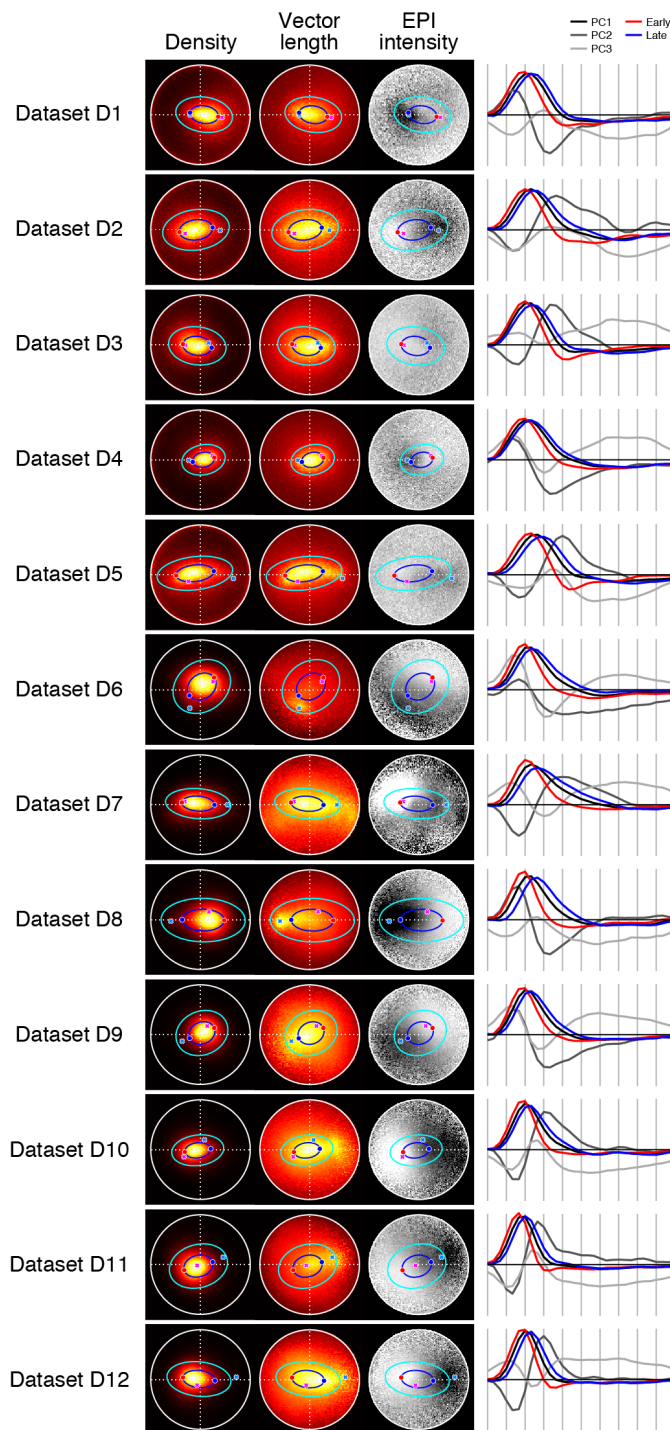
1
2

Supplementary Material



3
4
5
6
7
8
9

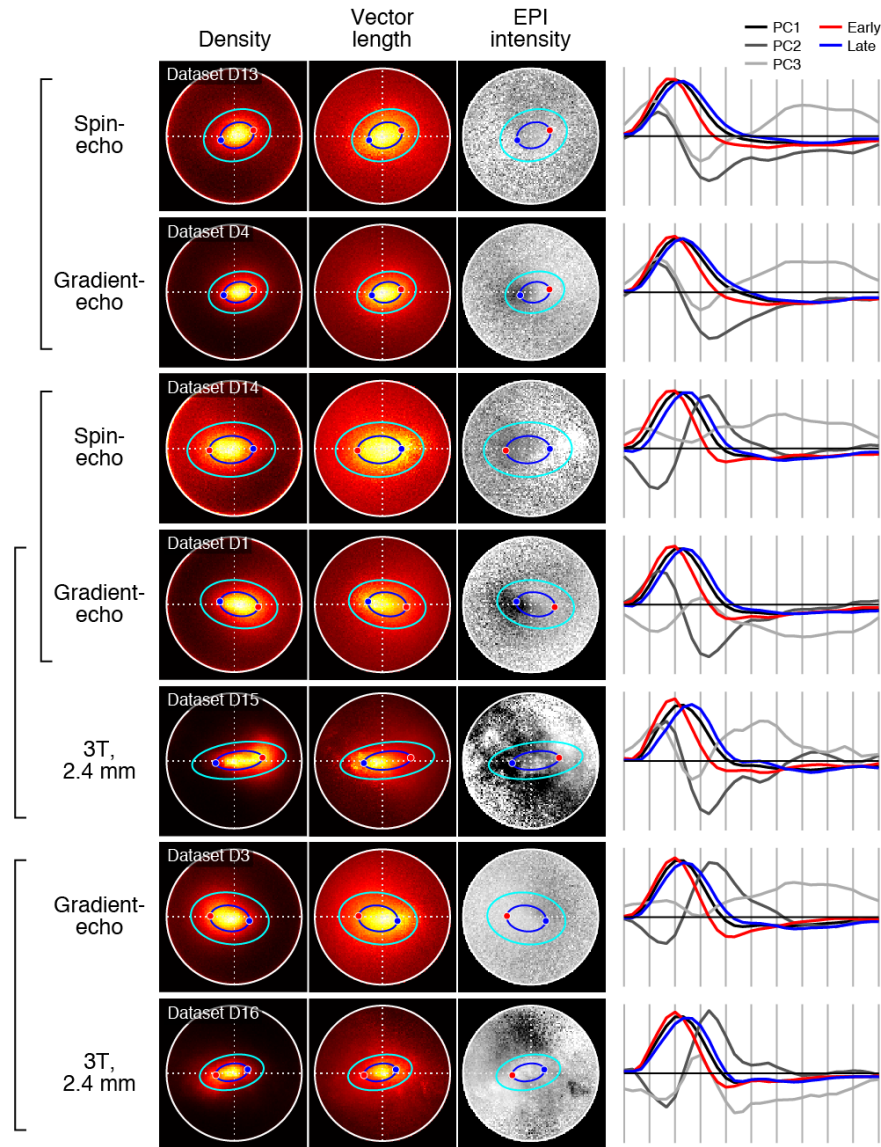
Supplementary Figure 1. Comprehensive summary of BOLD timecourses. Same format as Figure 2. Qualitative patterns of results are highly consistent across datasets (e.g. time-to-peak is delayed at superficial depths). However, there is substantial quantitative variation across datasets (e.g. time-to-peak is short in Dataset D11 but long in Dataset D5). This underscores the importance of tailoring timecourse derivation to individual subjects or scan sessions.



1
2
3
4
5
6
7
8
9

Supplementary Figure 2. TDM results for the high-resolution gradient-echo datasets (D1–D12).

Same format as Figure 4, except for the following addition: magenta and cyan crosses indicate the early and late timecourses derived from the ICA-based procedure (see Section 2.8.4). TDM consistently identifies reasonable early and late timecourses in each dataset. The ICA-based procedure yields similar timecourses in some datasets (e.g. D4), but diverges substantially in others (e.g. D8). It appears that timecourses with very large BOLD responses (see D8) is a major factor that influences the timecourses returned by ICA.



1
2
3
4
5
6

Supplementary Figure 3. TDM results for the alternative acquisition protocols (D13–D16). Same format as Figure 4. To facilitate comparison, we place results obtained using the spin-echo and low-resolution protocols next to results obtained using the high-resolution gradient-echo protocol.