

Video-based motion-resilient reconstruction of three-dimensional position for functional near-infrared spectroscopy and electroencephalography head mounted probes

Sagi Jaffe-Dax,^{a,*†} Amit H. Bermano,^{b,c,†} Yotam Erel,^c and
Lauren L. Emberson^a

^aPrinceton University, Psychology Department, Princeton, New Jersey, United States

^bPrinceton University, Computer Science Department, Princeton, New Jersey, United States

^cTel-Aviv University, School of Computer Science, Tel Aviv, Israel

Abstract

Significance: We propose a video-based, motion-resilient, and fast method for estimating the position of optodes on the scalp.

Aim: Measuring the exact placement of probes (e.g., electrodes and optodes) on a participant's head is a notoriously difficult step in acquiring neuroimaging data from methods that rely on scalp recordings (e.g., electroencephalography and functional near-infrared spectroscopy) and is particularly difficult for any clinical or developmental population. Existing methods of head measurements require the participant to remain still for a lengthy period of time, are laborious, and require extensive training. Therefore, a fast and motion-resilient method is required for estimating the scalp location of probes.

Approach: We propose an innovative video-based method for estimating the probes' positions relative to the participant's head, which is fast, motion-resilient, and automatic. Our method builds on capitalizing the advantages and understanding the limitations of cutting-edge computer vision and machine learning tools. We validate our method on 10 adult subjects and provide proof of feasibility with infant subjects.

Results: We show that our method is both reliable and valid compared to existing state-of-the-art methods by estimating probe positions in a single measurement and by tracking their translation and consistency across sessions. Finally, we show that our automatic method is able to estimate the position of probes on an infant head without lengthy offline procedures, a task that has been considered challenging until now.

Conclusions: Our proposed method allows, for the first time, the use of automated spatial co-registration methods on developmental and clinical populations, where lengthy, motion-sensitive measurement methods routinely fail.

© The Authors. Published by SPIE under a Creative Commons Attribution 4.0 Unported License. Distribution or reproduction of this work in whole or in part requires full attribution of the original publication, including its DOI. [DOI: [10.1117/1.NPh.7.3.035001](https://doi.org/10.1117/1.NPh.7.3.035001)]

Keywords: functional near-infrared spectroscopy; infant neuroimaging; spatial co-registration; photogrammetry; convolutional neural network.

Paper 20009R received Feb. 6, 2020; accepted for publication Jul. 6, 2020; published online Jul. 20, 2020.

1 Introduction

Functional near-infrared spectroscopy (fNIRS) and electroencephalography (EEG) require knowledge of the positioning of probes on subjects' scalps. Knowing the position of probes is an essential prerequisite for key analytic steps: aggregating results from a group of subjects,

*Address all correspondence to Sagi Jaffe-Dax, E-mail: jaffedax@princeton.edu

†These authors contributed equally to this work.

comparing between groups of participants (particularly from different developmental populations where head size systematically varies with group), and estimating the source of the recorded signal (in fNIRS^{1,2}). Desired positions of probes are determined *a priori* in accordance with either the standard positioning system (e.g., 10 to 20 in EEG³) or a predefined array that target specific cortical regions (typically in fNIRS⁴). However, it is not possible to place the probes exactly on the desired scalp locations for every experimental session. This is due to various reasons, including variability in head size and shape between participants, different practices between researchers, and even head growth between sessions in longitudinal studies with infants.⁵ Therefore, there would be substantial improvement in analytic methods if there was a reliable, quick, and robust method for calibrating the three-dimensional (3-D) positions of the probes relative to the subject's head.

Existing methods for 3-D estimation of probes' location are not suitable for early developmental or clinical populations and reduce the portability of these methods [one of their major benefits over magnetic resonance imaging (MRI)]. Using the most popular approach, the 3-D digitizer (e.g., Polhemus' Fastrak or Patriot, Colchester, VT) to measure even a small number of points on the participant's scalp typically requires a few minutes while a participant remains completely motionless.⁶ Even with highly skilled developmental cognitive neuroscientists, early developmental populations (and many clinical populations) cannot meet this requirement, thus there are no published studies using a 3-D digitizer for probe location estimation in early developmental populations (for older ages and limitations see Refs. 7 and 8). Furthermore, measurements with a 3-D digitizer are often confounded by interference from metal objects in the participant's surroundings. Thus, even in populations that can comply with these methods, it is difficult to move this method between experimental contexts. This is problematic as one of the major benefits of scalp-based neuroimaging is its portability (e.g., recording in rural Africa⁹). The 3-D digitizer itself can cause interference with sensitive medical devices, such as cochlear implants, which make it unsuitable for this key clinical population.

Alternative co-registration methods have been used in developmental populations with some success but each has major limitations. Recently proposed photogrammetry methods also rely on tightly restricted movements of the subject and allow only a specific movement trajectory of the camera around the subject.^{10,11} These methods require a lengthy recording of the participant (typically longer than two minutes, e.g., Delscan's Artec or GeoScan¹²) or impose motionless requirements in an extremely specific environment (e.g., Philips' GPS). Faster, more developmentally friendly, coregistration methods are based on manual photogrammetry. These manual methods require expertise in image annotation,^{2,13} are extremely time-intensive, and are prone to inter-researcher variability and biases. In these manual photogrammetry methods, several photos of the infant wearing the fNIRS cap are taken from a few different angles. The photos are then manually scaled according to a known size marker on the cap. The researcher then measures a few of the physical distances in the photos [e.g., nasion (Nz) toinion, right ear to left ear, and head circumference] and estimates the distances of the fNIRS channels from predefined fiducials.²

Semiautomatic coregistration methods require extensive physical measurements (i.e., using a measuring tape¹⁴) or MRI scans of each subject wearing the fNIRS cap in the scanner.¹⁵⁻¹⁷ These methods are also not suitable for developmental studies for the same reasons specified above, as they involve motion restrictions for a few minutes and unpleasant environments. Other semimanual photogrammetry methods use commercial 3-D scanners, but still require manual annotation of probe locations on the output images.¹⁸ Despite their disadvantages, these photogrammetric methods are currently being used for developmental studies,^{19,20} since they offer the only reliable coregistration method in these populations. Our proposed video-based method takes the basic principles of these manual and semimanual methods and extends it to an automatic video-based one.

Another existing method involves positioning probes on the scalp relative to skull landmarks (e.g., EEG positions^{21,22}). These EEG positions can be then linked to the underlying cortical areas using existing atlases.²³ This method has the benefit of not requiring offline annotation as well as having the ability to locate central probes (i.e., the ones localized to EEG locations) across infants and populations. However, there are major limitations of this method. First, this method still requires minutes-long measurements with the participant, rendering it challenging for infant studies. Second, given the fixed sensor distances of fNIRS (compared to EEG),

even if one or a handful of probes are positioned carefully using skull landmarks, other probes are likely to vary across populations based on head size and shape. In other words, if a central probe position is positioned at Cz, the adjacent probes have a fixed configuration that typically is not changed based on head size or shape. Thus, this method (in combination with anatomical atlases) only allows alignment for a subset of probes.

In this paper, we present a video-based method that is both easy to implement by novice experimenters and is robust to participant's head movements. Our method requires only ~ 20 s of video using widely available photographic equipment, recorded around the participant's head, with the probes already mounted on the scalp. During acquisition, the participant can move his or her head freely without jeopardizing the accuracy of the measurement. We report both the validity and the reliability of our video-based method compared to the traditional 3-D digitizer on a group of adult participants and show it to be reliable as the 3-D digitizer golden standard. Moreover, we also demonstrate the feasibility of this approach with early developmental populations—to which the 3-D digitizer cannot be applied.

2 Methods

Our automatic video-based method is depicted in Fig. 1 and includes a preparatory step and four subject-specific steps. In a preprocessing step, we capture the cap in perfect conditions—well aligned on a plastic head. We then manually mark the positions of all fiducials and probes on the reconstructed 3-D model, to create our reference [model cap, Fig. 1(a), see below fiducials' definition]. The model cap is used to learn the spatial relationship between fiducial points and probe positions. This allows us to calculate the latter through interpolation of the former. This preregistration step is required for each new configuration of probe positions (i.e., for new regions of interest, caps, or populations). Then for each subject, we do the following: first, the participant's head, with mounted cap and probes, is captured through a short video [Fig. 1(c)], using an off-the-shelf mid- to high-end camera [e.g., GoPro Hero6; GoPro, San Mateo, CA, see Fig. 1(b)]. Then using a convolutional neural network (CNN) [Fig. 1(d)], we crop away the background of each frame in the video, leaving only the cap and a part of the subject's head [Fig. 1(e)]. This facilitates handling head motion, since it focuses the next steps solely on the cap, completely eliminating distractions. Next, we reconstruct a 3-D model of the head using computer vision techniques for object 3-D reconstruction [structure from motion (SfM);²⁴ Fig. 1(f)]. Finally, we extract the coordinates of specific fiducials—an intuitive identification task once we already have a reconstructed model, in order to find all probe positions in MNI coordinates [Montreal Neurological Institute; Fig. 1(g)].

2.1 Cap Preparation

The basic requirement of the method is good tracking of the cap-mounted head throughout the SfM process, which is facilitated by feature-rich captured objects. To enrich the cap-mounted head with features, we replace the solid-color connector sheets of the fNIRS cap with multiple two-color patterned plastic sheets [red and blue in the example in Fig. 1(a)]. We use Perlin noise²⁵ to generate distinct patterns. A script for generating this specific pattern is in the github.com/sagijaffedax/VideoRecon, but any distinct nontrivial color pattern will suffice, e.g., drawings or cartoons. In addition, we place six solid-color stickers [green in Fig. 1(a)] at specific positions on the cap: around the vertex, around theinion, between the vertex and theinion, and on the left, right, and front edges of the cap. These positions roughly correspond to Cz, Iz, Pz, T7, T8, and Fpz of the standard 10–20 system, respectively.

2.2 Participant Preparation

We fit the cap on the subject and place three additional solid-color stickers [green in Fig. 1(c)] at specific anatomical positions on the participant fiducials: Nz, left preauricular, and right preauricular. In order to ensure other objects in the scene do not confuse our cropping method, we clear the surrounding of the participants from objects with similar color to the two-color plastic sheets and fiducial markers (e.g., no toys with similar colors, if the parent's shirt had similar colors,

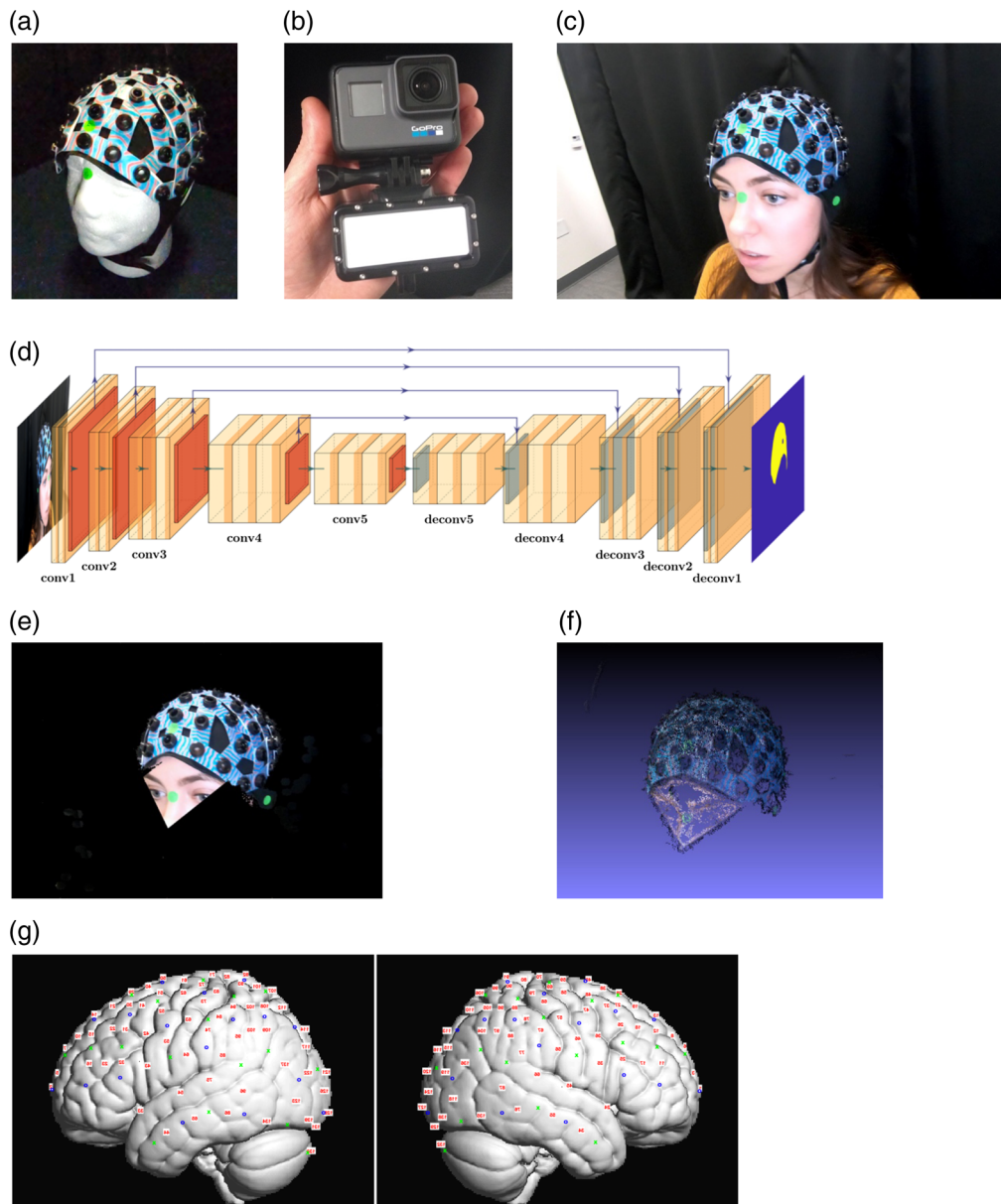


Fig. 1 Method steps overview. (a) The model cap is captured in perfect conditions against a solid background with no movements. (b) The short video is captured using an off-the-shelf camera with diffused flashlight. (c) fNIRS cap is assembled with two-color pattern sheets and solid-color stickers placed in the fiducial positions to be estimated. (d) The architecture of the CNN used to create masks for the images. (e) The frames are cropped using the mask from the neural net, a facial landmarks detector, and a sticker classifier. (f) All cropped frames are combined to a 3-D model through the popular SfM computer vision technique. (g) Fiducials are extracted from the model and channel positions are projected onto MNI space using SPM-fNIRS. Sources in blue circles; detectors in green crosses; and channels in red numbers.

we cover it with a dark barber cape, etc.). When the environment cannot be cleared of these objects (e.g., infant clothing or pacifier), we allow the experimenter to choose the color of the solid-color stickers to be distinct from the subject's clothing, and manually mark the chosen color after the fact (see next steps).

With this visual information about the spatial positions of the cap relative to the participant's head, our problem boils down to finding the spatial relation between the points on the head to those on the cap. Finding this relation is sufficient for exact positioning of the cap on the participant's head.

2.3 Video Capture

As discussed below, the core of our method is based on finding the 3-D positions of the sticker-marked fiducials using a computer vision algorithm. Since the 3-D relationships between the fiducials are required, it is imperative all of them are clearly visible in the captured video, and that the computer vision system's tracking is not lost during video acquisition. Along with patient comfort, these are the driving concerns for the process of capturing the video. Hence, we suggest the following procedure to help ensure successful tracking of each fiducial. We record the cap placement on the subject's head using a GoPro Hero 6 Black. We use the slow-motion (240 fps) feature to reduce motion blurring. Of course, any camera with a similar slow-motion option will suffice for the requirements of our method (e.g., iPhone 8 or higher). Furthermore, we attach a diffused flashlight to the camera to compensate for uneven lighting conditions across all head sides [see Fig. 1(b)]. Finally, we start and end the video recording at approximately the same position relative to the subject's head. Note that depending on environmental conditions not all of, or even none of, the steps are necessary for success. These steps, however, have proven to be reliable and simple to implement. An example for such a video recording is available in Ref. 26.

2.4 Cropping the Video Frames

Passing irrelevant information about the background of the given object-to-be-reconstructed to the SfM method often yields bad reconstructions, partial reconstructions, or falsely merged objects. These undesired reconstruction results are especially common in the presence of motion, i.e., when the object is moving during the video acquisition. Head movement and subsequent bad reconstruction are the major challenges that existing automated photogrammetry coregistration methods face with early developmental populations.

In order to avoid passing irrelevant information about the background to the SfM stage and to prevent the issues that arise from participants' movements, we crop everything but the head and cap from each video frame. First, we maintain color consistency as much as possible, through white-balance correction of the video frames using OpenCV²⁷ package for MATLAB[®].²⁸ Then the user manually marks, from a single frame, the color of the stickers (green in Fig. 1) that were placed on the fiducial points. Manually selecting the color of the solid-color stickers allows flexibility regarding the choice of the used stickers—one could determine a fixed color for the fiducial stickers, which would render the last aforementioned step unnecessary. However, the sticker color should be distinctly different from any other color in the scene (see participant preparation stage).

After these steps, we use the following automatic processes to produce a cropping mask for each frame.

1. *Identify the pixels corresponding to the fiducial stickers.* Each cluster of at least five connected pixels with a hue within $\pm 0.15\%$ of the sticker hue [measured in hue-saturation-value (HSV) space] is classified as a sticker. A circle is defined around each cluster with an additional margin of 20% of the radius.
2. *Identify the cap in each video frame.* We designed and trained CapNet (available in the github.com/sagijaffedax/VideoRecon), a CNN based on the segmentation neural network,²⁹ pretrained on the CamVid³⁰ dataset, and fine-tuned on our dataset. We found that manually annotating 100 images with pixel-wise classes of “cap” and “background” creates a sufficient training set. The convolutional layers used in the architecture composing the encoder part are [(width \times height \times depth) \times stack]:

$$(360 \times 480 \times 64) \times 2, \quad (180 \times 240, 128) \times 2, \quad (90 \times 120 \times 256), \\ (45 \times 60 \times 512) \times 3, \quad (22 \times 30, 512) \times 3.$$

The encoder is followed by a symmetrical stack of deconvolutional layers composing the decoder [Fig. 1(d)]. Batch-normalization layers are placed after each convolutional layer, followed by a rectified linear unit (ReLU) activation layer. Max-pooling layers are used between each stack of convolutional layers, and each one of them additionally has a skip

connection to its max-unpooling counterpart. Images are resized to fit the input layer using bicubic interpolation. Network training was done using the stochastic gradient descent with momentum optimizer, with a momentum value of 0.9, for 100 epochs. Network training stopped after reaching the very satisfactory global accuracy of 98% and weighted average (across classes) intersection over union of 90%.

3. *Identify participants' face.* We found that keeping a part of the subject's face in the image, in addition to the cap, improves the 3-D reconstruction of the Nz region. The participant's face is identified using an off-the-shelf neural-network-based facial detector, found in OpenCV.³¹ We define a polygon connecting the nose, external eye-edges, and forehead, and add it to the cap in the region of the frame that is kept for further analysis. The application of such an example polygon can be seen in Fig. 1(e).

The rest of the frame is blackened [Fig. 1(e)], and frames in which more than 98% of the frame are blackened are rejected from further processing. These frames are either blurred frames due to fast motion of the camera or frames in which the participant's head is out of frame.

2.5 3-D Surface Reconstruction

In order to find the 3-D coordinates of the fiducial points, a 3-D model is automatically reconstructed. The cropped frames are passed to an SfM³² software (visual SfM³³). This widely used algorithm searches for shared features between frames (or correspondences³⁴) and estimates the camera position for each frame accordingly. Note that understanding this mechanism had enabled our key head-motion robustness contribution—when searching for correspondences between frames, a static background along with a moving head confuses the camera-position-estimation process. Hence, by removing the background in the cropping step, a moving head can be simply interpreted as additional camera movement, rendering accurate and successful reconstructions of our object of interest. We instruct the software to only search for correspondences that are less than two seconds apart. Searching for temporally distant correspondences is very computationally expensive and might yield faulty matches, since the video does not show the same parts for long periods. In addition, since the videos are taken such that they start and end the same way, matches are also searched for between the first and last frames of the video. This closed-loop helps in reducing accumulated drift and in cases where matches the tracking is lost midway through the video. A 3-D model can then be built [Fig. 1(f)] using the correspondences, camera positions, and simple geometric triangulations.

2.6 Fiducial Extraction and Probe Positions Estimation

Our basic assumption in this part is that the cap undergoes limited deformation when worn by the participant (e.g., the cap does not stretch according to differences in head size as an EEG cap would). Instead, differences in head size, shape, cap placement, etc. are captured by the positions of the markers on the cap relative to the head (cap points and fiducial points, respectively, see below). After these points are extracted, we assume the cap can only rotate globally with respect to the head, and scale along the anterior-to-posterior and left-to-right axes. In this way, head size is implicitly considered in our reconstruction of probe positions. Therefore, capturing five points on the cap (along with the ones on the head) are sufficient for complete cap-head registration. The position of the probes can be then interpolated similarly to common photogrammetry-based methods.² In practice, we used six points for redundancy, with two in the difficult-to-capture back region (Pz and Iz), because infants are sitting on their parent's lap.

From the reconstructed 3-D model, we identify the position of the fiducials by clustering nearby reconstructed vertices that hold the solid sticker color. Using RANSAC,³⁵ we determine the fiducial name associated with each aforementioned cluster (by looking at the expected relative position of all fiducials). In other words, we select three of the identified points and assign a labeling to them that is taken from the initial cap model. We then deform the head so that the three points match their counterparts and measure the distance of the rest of the markers from the ones they end up closest to. After choosing the labeling that minimizes the latter distances, we use the head fiducials (Nz, AR, AL, and Cz) to project the cap points (Cz, Pz, Iz, front, left,

right of the cap; we are assuming Cz is relatively in the correct position) into MNI space by rotation and scaling.³⁶ The channel positions are then interpolated in MNI space based on the preregistered model cap relative to the cap points using “SPM for fNIRS” [Fig. 1(g)].³⁷ The interpolation is performed using barycentric coordinates: we represent each manually marked probe p_j as a weighted sum of all fiducials $f_i, i = 1 \dots 9$. In other words, we find all $w_{i,j}$'s such that $p_j = \sum_{i=1}^9 w_{i,j} \cdot f_i$ in the preregistration step and simply use the same formula with the newly found fiducials to extract probe positions.

2.7 3-D Digitizer Measurement

For comparison to our automated video-based method of 3-D reconstruction, we collected the 3-D coordinates of the head points and cap points using a 3-D digitizer (Fastrak, Polhemus, Colchester, VT) twice for each adult subject in each session. We then projected the channel positions into MNI space using the same predefined model cap and SPM for fNIRS.

2.8 Estimation of Validity and Reliability of Channel Position Estimation

We compared the channel positions (139 channels; LABNIRS, Shimadzu inc., Kyoto, Japan) determined based on our new video-based method to the positions found using the dominant method in adult participants, the 3-D digitizer (Fastrak, Polhemus, Colchester, VT). We did the comparison twice for each subject ($N = 10$) for the same session and compared the positions of the same channels both between the two methods (intermethod validity) and within each one (intramethod reliability). The intramethod comparison was used to estimate the error of each method independently (test–retest reliability).

In addition, we measured the channel position for a separate group of adult subjects ($N = 10$) in two separate sessions (separate days, multiple cap placements) to estimate whether the video-based method captured the shifts in cap positioning on the same subject as compared to the 3-D digitizer.

2.9 Feasibility for Early Developmental Populations

Finally, we show how this method can be used to estimate probe positions with early developmental populations. Specifically, we present results from a 6-month-old infant. The example video is available here.³⁸ Although it would be ideal to compare our video-based reconstruction methods as conducted in infants to other methods, we cannot conduct the same comparison as adults because it is not practically possible to collect the digitized positions of the fiducials and probes using the 3-D digitizer. After months of effort on this front and communication with other developmental cognitive neuroscience labs, we reached the conclusion that this method is too sensitive to motion to be reliable, and thus, is unusable for this population.

3 Results

3.1 Video-Based Method, Validated through the Standard 3-D Digitizer

We found that the validity of our new video-based method was highly comparable to the field's standard method—the 3-D digitizer [Figs. 2(a) and 2(b) right]. Namely, the distances between the positions of the same channel as estimated by the two methods were similar, albeit slightly larger, to the distances between the positions of the same channel as measured by the 3-D digitizer twice (intermethod validity of 3.4 ± 0.9 mm compared to intramethod reliability of 2.6 ± 0.6 mm; mean \pm STD; $t_9 = 3.1, P < 0.05$). Thus, there is a high correspondence between the channel positions estimated by our video-based methods and the 3-D digitizer.

Moreover, the reliability of our video-based method was better than the reliability of the 3-D digitizer. The distances between the estimated positions of the same channel as measured twice by our video-based method [2.0 ± 0.5 mm; Fig. 2(c) right] was smaller than that of the 3-D digitizer ($t_9 = 2.6, P < 0.05$; paired t -test). Overall, the magnitude of the errors that we found

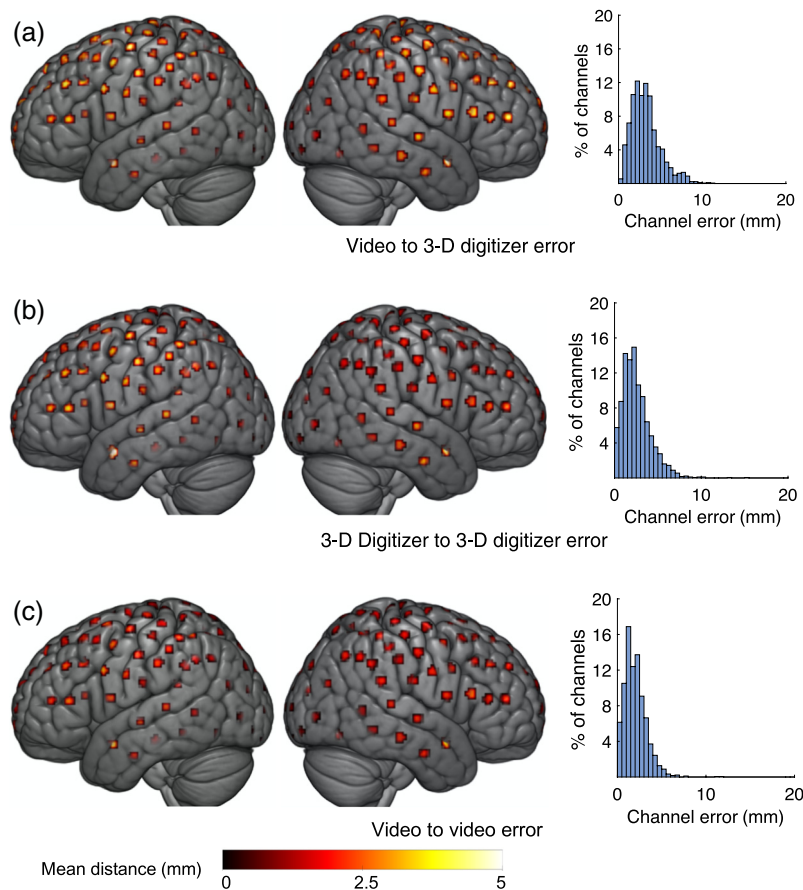


Fig. 2 Spatial distribution of intermethod validity and intramethod reliability, comparing our method and the digitizer. (a) Distance from our (video) method to the 3-D digitizer one (intermethod validation; 3.4 ± 0.9 mm). (b) Distance between two measurements of the 3-D digitizer (intramethod reliability; 2.6 ± 0.6 mm). (c) Distance between two measurements of video-based method (intramethod reliability; 2.0 ± 0.5 mm). Distances are represented by the color and the diameter of the patches for each channel position.

is much smaller than the radius of the channels, whose positions were estimated (typically on the order of centimeters).

Although the discrepancies between the two position estimation methods were spread evenly through the scalp, we did observe a tendency for larger distances in the anterior temporal channels [Fig. 2(a) left and middle]. This increased error at this position stems from smaller reliability of both methods in these areas [Figs. 2(b) and 2(c) left and middle]. The lower reliability of both methods in the anterior temporal lobe is probably a result of the low angle of camera or digitizer that is required to acquire the position of the fiducials in this region.

3.2 Video-Based Method Captures Cap Positioning Variability

We now considered how well these two localization methods identify any change in position of the cap (i.e., shift or displacement of each channel) for the same participant across sessions. We found a significant regression coefficient for the size of the estimated shift in channel position (in mm) between the two methods [red line in Fig. 3(a); $F_{1,1388} = 5.86$ $p < 0.05$; linear mixed effect model]. Namely, the shift size as estimated by the video-base method predicts the shift size as estimated by the 3-D digitizer method. As Fig. 3(a) clearly shows, there were many cases where the two estimation methods did not agree on the shift size between the two sessions [e.g., a high shift size was estimated by the 3-D digitizer and a small shift size was estimated by the video-based method for the points on the right-hand side; channels below the diagonal—black dashed line in Fig. 3(a)]. We manually examined a few cases of large discrepancies

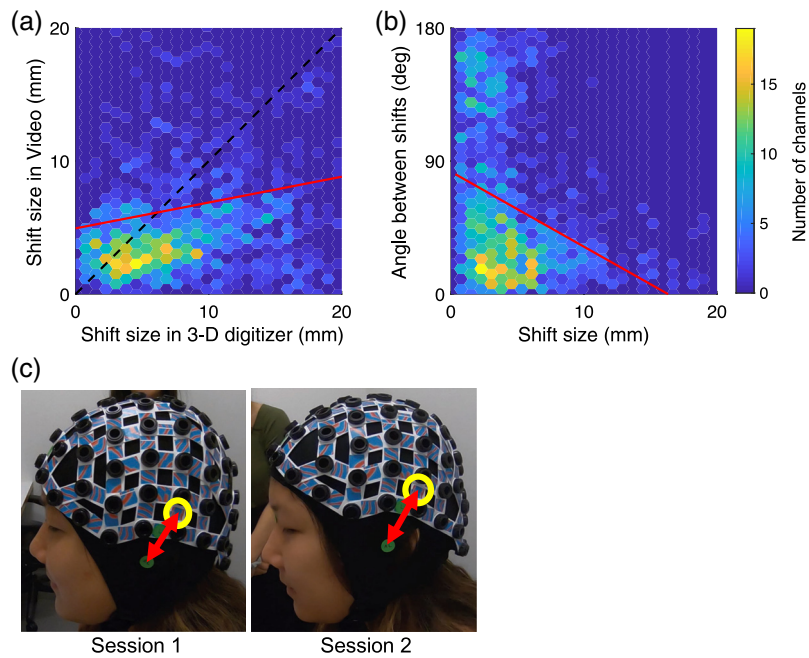


Fig. 3 Comparison of estimation of cap shifting between session through the video-based and the 3-D digitizer methods. (a) Estimation of shift size between the sessions of the video-based method versus the 3-D digitizer one. The color of each hexagon represents the number of channels in its range. We found a high correlation between the shifts as estimated by the two methods ($F_{1,1368} = 21.8$ $p < 10^{-5}$; red line). (b) Angles between the shift as estimated by the video-based method and the shift as estimated by the 3-D digitizer method versus the size of the shift between sessions. The distribution of angles is skewed toward zero ($K = 0.3$, $p < 10^{-25}$; distribution along the ordinate). We found a negative correlation between the size of the angle between estimated shifts and the size of these shifts. When the estimated shifts are larger, the direction of the shifts as estimated by the two methods were similar (smaller angle). (c) Manual estimation of a case of large discrepancy between the two estimation methods. Left: photo of participant 1 in session 1. The distance of channel 66 (circled in yellow) from the left tragus (lower green sticker) was 49.3 mm (the length of the red arrow). Right: photo of participant 1 in session 2. The distance between channel 66 and the left tragus was 53.3 mm. The shift of the channel position between the two sessions was ~ 4 mm and was closer to the video-based estimation (3.6 mm) than to the digitizer-based estimation (24.5 mm).

between the shift estimation of the two methods. In most cases, the video-based method was actually closer to estimating the shift magnitude correctly. For example, in Fig. 3(c), the shift between the sessions of channel 66 (circled in yellow) for participant 1 was estimated as 3.6 mm by the video-based method and 24.5 mm by the 3-D digitizer. Our manual estimation, based on image-space measurements scaled by known distances between channels, suggests a shift of 7 mm—a number much closer to the video-based estimation than to the 3-D digitizer-based one [Fig. 3(c)].

We also found that the shifts in channel positions were estimated in similar directions using the two methods. We measured the spatial angle between the direction of the shifts as estimated by the video method and the direction of the shift as estimated by the 3-D digitizer. The estimated shift of each channel between sessions can be represented as a line between the two 3-D estimated positions. Each method thus estimates such a line and we can measure the level of agreement between the methods by measuring the angle between these lines. The angle between the direction of the shifts, as estimated by the two methods, was small [closer to zero than a uniform distribution; distribution along the ordinate in Fig. 3(b); $K = 0.3$, $p < 10^{-25}$; Kolmogorov–Smirnov test; mean angle \pm STD: 60.6 deg \pm 50.1 deg]. Crucially, the similarity in shift direction was found to be greater in the larger shifts [red line in Fig. 3(a); $F_{1,1368} = 21.8$ $p < 10^{-5}$; linear mixed effect model], where this similarity in direction matters most (i.e., when the cap shifts substantially and may potentially result in a change in the cortical region that the channel is assigned to).

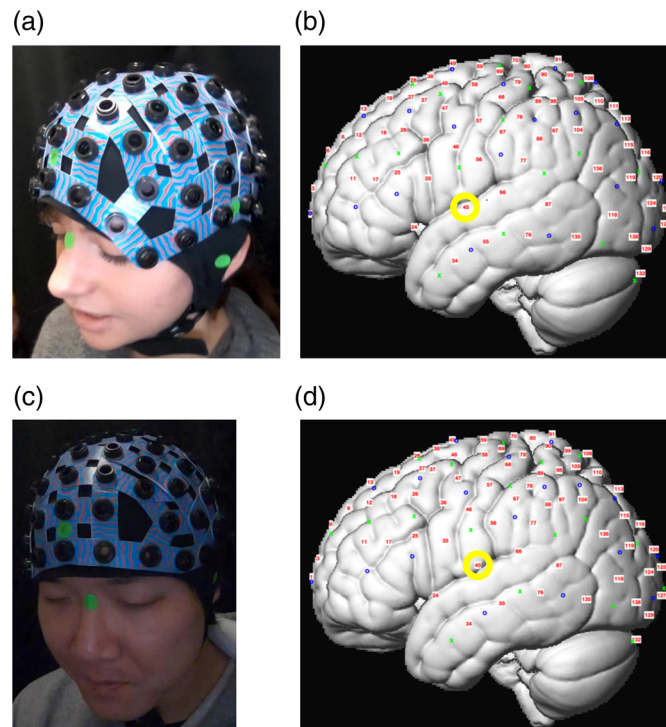


Fig. 4 Comparison of two subjects with relatively distinct head sizes. The subject with the smaller head (a) had a mean error of 2.7 mm. The subject with the larger head (b) had a mean error of 2.3 mm. Due to the rigidity of the fNIRS cap, the probe positions on the smaller head (c) were shifted slightly downward relative to the larger head (d). For example, channel 45 (circled in yellow) was estimated in STG on the smaller head size, but in postcentral gyrus on the larger head size.

Investigating these shifts in cap placement across sessions reveals how crucial it is for any coregistration to be accurate. Even if recordings of the same subject take place with the same cap mere days apart, these shifts in cap placement are large enough to change the cortical region under a channel. In our sample, more than 8% of the channels were shifted to a different lobe between sessions.

Additionally, our automatic method does not rely on the specific head size or model it specifically. Instead, it automatically adjusts for the head size by measuring the distance between the fiducial points and the cap directly, taking the head size implicitly into account (see Sec. 2.6). Thus, we expect no correlation between the error magnitude and head sizes. Unfortunately, we did not take independent measurements of our subjects' head sizes (e.g., circumference, inion to Nz distance, etc.). However, we qualitatively compare the errors that were obtained for a distinctly larger head with the errors that were obtained for a smaller one. We did not find any systematic difference (Fig. 4). We also find in the visual comparison between these two head sizes that the probes are shifted down on the smaller head and up on the larger head reflecting the differences in the relative size of the cap to the participant's head size. This logical transformation reflects the fact that head size is accommodated by our method through differences in the fiducial points and cap points and does not need to be included explicitly.

3.3 Automatic Estimation of Channel Position in Infants

We successfully reconstructed a 3-D model from video recordings of infant participants and estimated the channel positions on their scalp through our automated video-based method (using methods described above, Fig. 5 for one sample infant). Thus, here we present a proof of concept that the same method can be used for early developmental populations. In practice, we have employed this method for many dozens of infants at different ages. Since there can be no comparison with a 3-D digitizer, we focus this paper of the validity of the overall method in adult participants in comparison to established methods.

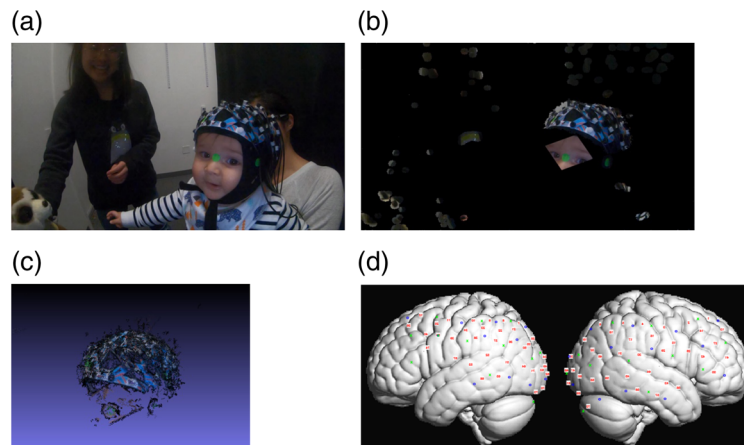


Fig. 5 Infant channel position estimation using the video-based method. (a) The captured video. (b) The images were cropped to eliminate confounds of movements. (c) A 3-D model was reconstructed using visual SfM. Coordinates of fiducial points were extracted from the model. (d) Channel positions were interpolated from fiducial coordinates onto MNI space.

As is apparent in this example, the cap was slightly misplaced on the infant, tilting to the right [Fig. 5(a)]. Our probe position estimation method successfully estimated this misplacement [Fig. 5(d)].

4 Discussion

The field of developmental cognitive neuroscience requires an automated reliable method to estimate the spatial position of probes on the scalp. This method is especially important when researchers are drawing conclusions about the development of functional anatomy that underlies the recorded signal from EEG or fNIRS. As we show in this work, even minor unintentional shifts in cap positioning can yield dramatic changes between the cortical regions that are being recorded in different sessions. Such uncertainty in position estimation and changes between sessions may lead to increased levels of noise in neuroimaging results and even to false conclusions regarding the underlying cortical activity.

We present an innovative approach for 3-D localization of probes relative to a participant's head, which is highly needed in scalp-based methods of neuroimaging such as EEG and fNIRS and particularly for developmental and clinical populations where current methods are inadequate. The proposed method requires only a short-video recording during the experimental session. Moreover, during this video, participants can move freely and interact with research staff and caregivers. These differences in methodology make this method suitable for developmental and clinical populations.

We designed this new method to substantially improve the ability to conduct probe localization in developmental and clinical populations. Currently available methods have major disadvantages: some are susceptible to participant's movements and take a relatively long amount of time.¹⁴ Thus they are not useable for developmental populations (3-D digitizer) and are onerous for clinical populations.³⁹ Other methods require laborious and bias-prone manual annotation (manual coregistration²). The video-based method we present here is both motion-resilient and automatic. Additionally, this method does not require a constant power supply and is not sensitive to metal in the environment. Therefore, this method is applicable not only to the population it was designed for, but also to out-of-lab settings and may be widely used in clinical settings and global neuroscience projects. The automatic analysis of the videos does not require the extensive manual annotation of images and thus is less prone to bias and much less laborious.

We present evidence that the proposed method is both valid and reliable compared to the standard method currently employed in the field, the 3-D digitizer. Overall, we find strong agreement between these methods of localization. We find a high level of agreement between methods

(intermethod validity) and within the method in a test–retest repeated measurements configuration (intramethod reliability). Importantly, our video-based method exhibits better intramethod reliability compared to the traditional 3-D digitizer. In the region where these two methods had the lowest agreement, the anterior temporal region, we suspect that it stemmed from difficulty in placing the digitizer in the correct angle, perpendicular to the scalp. Throughout the scalp, the error size of our estimation is an order of magnitude smaller than the size of the channels that are to be estimated.

We further tested the localization abilities of these methods by examining changes in localization of channels between sessions with the same subject. We find that these methods similarly estimate the displacement of the cap across sessions. We estimated the session-to-session variability in cap placement using the video-based and the digitizer methods and found good agreement between the two. As can be expected, we found that the agreement between the two estimation methods about the shift direction was higher when the size of the shift was larger. Where there were disagreements between the methods (e.g., large difference between the size of the estimated shift), we have found that the video localization method was usually more accurate in channel position estimation.

Perhaps most importantly, we have demonstrated the successful measurement of infants' head-mounted probe positions: an estimation that is not possible using currently available online methods (e.g., 3-D digitizers). Of course, successfully estimating misalignments in cap placement is of even greater value when administering experiments with developmental population, since it is a more prevalent problem in such cases.

Using a different model cap, as described at the beginning of the method section, researchers can easily adapt the outline proposed here to their own system and to any cap type.

The method of course has some limitations. First, much like other methods, it is acutely dependent on recovering all the sticker-marked fiducials (excluding the one redundant pair). Unlike other methods though, the missing information can be recovered from the abundance of visual information found in the captured video. Second, while quick to capture, the method takes a few minutes to compute the registration results, meaning the technician is not aware of the quality of the registration while running the experiment. Therefore, future implementation of the proposed video-based method will include a user interface, facilitating easier insertions of new cap models, as mentioned above, and other features such as manual corrections to badly captured or detected positions. This approach potentially can profit from automation wherever possible but from human understanding where it is not. Another interesting direction to take would be to develop a single neural-network-based pipeline producing a real-time indication of the calibration results. This could be done perhaps by collecting online feedback from users of the system regarding specific reconstructions.

Finally, we hope our proposed coregistration method will allow the combination of data from different locations and labs and will enable the creation of a unified anatomical framework for fNIRS analysis.

Disclosures

All authors declare that they have no conflicts of interests.

Acknowledgments

This research was supported by grants from the National Institutes of Health (Grant No. R004R00HD076166-02), McDonnell Foundation (Grant No. 220020505), and the Bill and Melinda Gates Foundation—Modeling Neurodevelopment: Physical Growth above the Neck.

References

1. S. L. Ferradal et al., “Atlas-based head modeling and spatial normalization for high-density diffuse optical tomography: in vivo validation against fMRI,” *Neuroimage* **85**, 117–126 (2014).

2. S. Lloyd-Fox et al., “Coregistering functional near-infrared spectroscopy with underlying cortical areas in infants,” *Neurophotonics* **1**(2), 025006 (2014).
3. V. Jurcak, D. Tsuzuki, and I. Dan, “10/20, 10/10, and 10/5 systems revisited: their validity as relative head-surface-based positioning systems,” *Neuroimage* **34**(4), 1600–1611 (2007).
4. R. N. Aslin, M. Shukla, and L. L. Emberson, “Hemodynamic correlates of cognition in human infants,” *Annu. Rev. Psychol.* **66**(1), 349–379 (2015).
5. C. Bulgarelli et al., “The developmental trajectory of fronto-temporoparietal connectivity as a proxy of the default mode network: a longitudinal fNIRS investigation,” *Hum. Brain Mapping* **41**(10), 2717–2740 (2020).
6. A. K. Singh et al., “Spatial registration of multichannel multi-subject fNIRS data to MNI space without MRI,” *Neuroimage* **27**(4), 842–851 (2005).
7. G. A. Buzzell et al., “Development of the error-monitoring system from ages 9–35: unique insight provided by MRI-constrained source localization of EEG,” *Neuroimage* **157**, 13–26 (2017).
8. G. Zhou et al., “Development of effective connectivity during own- and other-race face processing: a granger causality analysis,” *Front. Hum. Neurosci.* **10**, 474 (2016).
9. S. Lloyd-Fox et al., “Functional near infrared spectroscopy (fNIRS) to assess cognitive function in infants in rural Africa,” *Sci. Rep.* **4**, 4740 (2015).
10. T. Clausner, S. S. Dalal, and M. Crespo-García, “Photogrammetry-based head digitization for rapid and accurate localization of EEG electrodes and MEG fiducial markers using a single digital SLR camera,” *Front. Neurosci.* **11**, 1–12 (2017).
11. A. Stopczynski et al., “The smartphone brain scanner: a portable real-time neuroimaging system,” *PLoS One* **9**(2), e96652 (2014).
12. P. M. R. Reis and M. Lochmann, “Using a motion capture system for spatial localization of EEG electrodes,” *Front. Neurosci.* **9**, 1–8 (2015).
13. L. L. Emberson, J. E. Richards, and R. N. Aslin, “Top-down modulation in the infant brain: learning-induced expectations rapidly affect the sensory cortex at 6 months,” *Proc. Natl. Acad. Sci. U. S. A.* **112**(31), 9585–9590 (2015).
14. D. Tsuzuki and I. Dan, “Spatial registration for functional near-infrared spectroscopy: from channel position on the scalp to cortical location in individual and group analyses,” *Neuroimage* **85**, 92–103 (2014).
15. X. Wu et al., “Evaluation of rigid registration methods for whole head imaging in diffuse optical tomography,” *Neurophotonics* **2**(3), 035002 (2015).
16. C. Whalen et al., “Validation of a method for coregistering scalp recording locations with 3D structural MR images,” *Hum. Brain Mapping* **29**(11), 1288–1301 (2008).
17. M. Chen et al., “Spatial coregistration of functional near-infrared spectroscopy to brain MRI,” *J. Neuroimaging* **27**(5), 453–460 (2017).
18. S. Homölle and R. Oostenveld, “Using a structured-light 3D scanner to improve EEG source modeling with more accurate electrode positions,” *J. Neurosci. Methods* **326**, 108378 (2019).
19. L. L. Emberson et al., “Deficits in top-down sensory prediction in infants at risk due to premature birth,” *Curr. Biol.* **27**, 431–436 (2017).
20. S. Lloyd-Fox et al., “Habituation and novelty detection fNIRS brain responses in 5- and 8-month-old infants: the Gambia and UK,” *Dev. Sci.* **22**(5), e12817 (2019).
21. T. Wilcox et al., “Hemodynamic response to featural changes in the occipital and inferior temporal cortex in infants: a preliminary methodological exploration: paper,” *Dev. Sci.* **11**(3), 361–370 (2008).
22. A. Blasi et al., “Test-retest reliability of functional near infrared spectroscopy in infants,” *Neurophotonics* **1**(2), 025005 (2014).
23. J. E. Richards et al., “A database of age-appropriate average MRI templates,” *Neuroimage* **124**, 1254–1259 (2016).
24. S. Ullman, “The interpretation of structure from motion,” *Proc. R. Soc. Lond. B* **203**(1153), 405–426 (1979).
25. K. Perlin, “Making noise,” in GDC Talk (1999).
26. S. Jaffe-Dax, “Video based method for fNIRS co-registration,” <https://youtu.be/sD75ctsGID4> (2019).

27. G. Bradski, "OpenCV," <https://opencv.org/> (2000).
28. K. Yamaguchi, "mexopencv," <https://github.com/kyamagu/mexopencv> (2018).
29. V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: a deep convolutional encoder-decoder architecture for image segmentation.," *IEEE Trans. Pattern Anal. Mach. Intell.* **39**(12), 2481–2495 (2017).
30. G. J. Brostow, J. Fauqueur, and R. Cipolla, "Semantic object classes in video: a high-definition ground truth database," *Pattern Recognit. Lett.* **30**(2), 88–97 (2009).
31. W. Liu et al., "SSD: single shot multibox detector," *Lect. Notes Comput. Sci.* **9905**, 21–37 (2016).
32. C. Wu, "Towards linear-time incremental structure from motion," in *Int. Conf. 3D Vision*, pp. 127–134 (2013).
33. C. Wu, "VisualSFM : a visual structure from motion system," 2011, <http://ccwu.me/vsfm/>.
34. D. G. Lowe, "Object recognition from local scale-invariant features," in *Int. Conf. Comput. Vision* (1999).
35. P. J. Huber and E. M. Ronchetti, *Robust Statistics*, 1.1, John Wiley & Sons, New York (1981).
36. A. C. Evans et al., "Anatomical mapping of functional activation in stereotactic coordinate space," *Neuroimage* **1**(1), 43–53 (1992).
37. S. Tak et al., "Sensor space group analysis for fNIRS data," *J. Neurosci. Methods* **264**, 103–112 (2016).
38. S. Jaffe-Dax, "Infant video-based co-registration for fNIRS," <https://youtu.be/ecn-GSGoRh8> (2019).
39. C. A. Anderson et al., "Adaptive benefit of cross-modal plasticity following cochlear implantation in deaf adults," *Proc. Natl. Acad. Sci. U. S. A.* **114**(38), 10256–10261 (2017).

Sagi Jaffe-Dax received his PhD (computational account and neural basis of dyslexia) in the Interdisciplinary Center for Neural Computation at the Hebrew University of Jerusalem. Currently, he is getting postdoc training in the Princeton Baby Lab, working on the mechanisms that support infants' learning during their first year of life. In Summer 2021, he will establish a Developmental Cognitive Neuroscience lab in the School for Psychological Sciences at Tel-Aviv University.

Amit Bermano is a senior lecturer (assistant professor) at the Blavatnik School of Computer Science in Tel-Aviv University since 2018. Previously, he was a postdoctoral researcher at the Princeton GraphicsGroup (2016–2018), and a postdoctoral researcher at Disney Research Zurich (2015). He has conducted his PhD studies at ETH Zurich in collaboration with Disney Research Zurich (2011–2015). His master's and bachelor's degrees were obtained at The Technion–Israel Institute of Technology.

Yotam Erel: Biography is not available.

Lauren L. Emberson is an assistant professor in the Psychology Department, Princeton University. She received her PhD from Cornell University and was a postdoctoral research associate at the University of Rochester. Her work uses fNIRS, as well as behavioral methods, to investigate perceptual development and learning in young infants. Her research consistently pushes both theoretical and methodological or technical boundaries with the ultimate goal of understanding how experience supports development.