

TOPICAL REVIEW • OPEN ACCESS

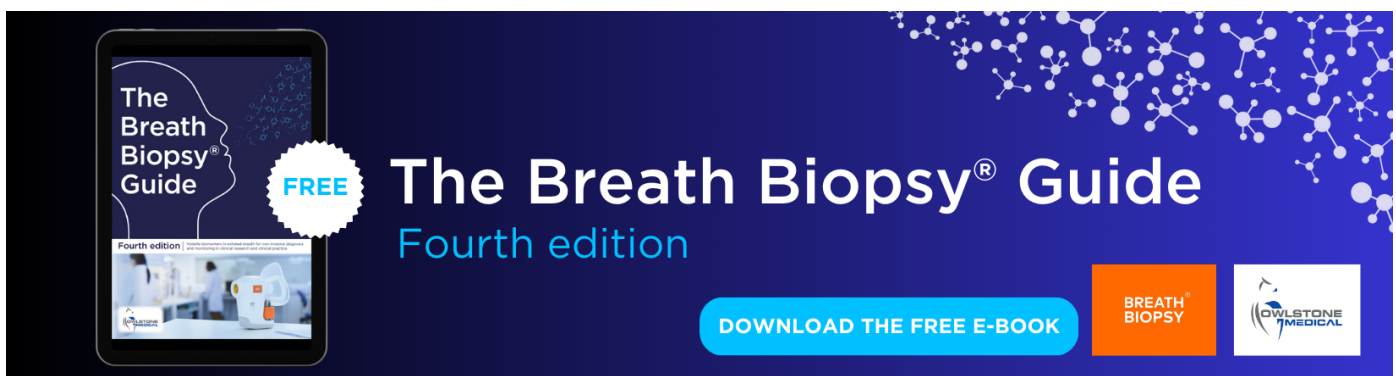
A review of rapid serial visual presentation-based brain–computer interfaces

To cite this article: Stephanie Lees *et al* 2018 *J. Neural Eng.* **15** 021001

View the [article online](#) for updates and enhancements.

You may also like

- [Multi-objective optimization approach for channel selection and cross-subject generalization in RSVP-based BCIs](#)
Meng Xu, Yuanfang Chen, Dan Wang et al.
- [Multi-source domain adaptation based tempo-spatial convolution network for cross-subject EEG classification in RSVP task](#)
Xuepu Wang, Bowen Li, Yanfei Lin et al.
- [A deep learning method for single-trial EEG classification in RSVP task based on spatiotemporal features of ERPs](#)
Boyuan Zang, Yanfei Lin, Zhiwen Liu et al.



The Breath Biopsy® Guide
Fourth edition

DOWNLOAD THE FREE E-BOOK

BREATH BIOPSY

OWLSTONE MEDICAL

Topical Review

A review of rapid serial visual presentation-based brain–computer interfaces

Stephanie Lees¹ , Natalie Dayan¹, Hubert Cecotti², Paul McCullagh¹, Liam Maguire¹, Fabien Lotte³ and Damien Coyle¹

¹ Faculty of Computing and Engineering, Ulster University, Belfast, United Kingdom

² Department of Computer Science, College of Science and Mathematics, California State University, Fresno, 2576 E. San Ramon MS ST 109 Fresno, CA 93740–8039, United States of America

³ Inria Bordeaux Sud-Ouest/LaBRI/CNRS/Université de Bordeaux/IPB, 200 Avenue de la Vieille Tour, 33405 Talence, France

E-mail: dh.coyle@ulster.ac.uk

Received 5 September 2017

Accepted for publication 3 November 2017

Published 24 January 2018



Abstract

Rapid serial visual presentation (RSVP) combined with the detection of event-related brain responses facilitates the selection of relevant information contained in a stream of images presented rapidly to a human. Event related potentials (ERPs) measured non-invasively with electroencephalography (EEG) can be associated with infrequent targets amongst a stream of images. Human–machine symbiosis may be augmented by enabling human interaction with a computer, without overt movement, and/or enable optimization of image/information sorting processes involving humans. Features of the human visual system impact on the success of the RSVP paradigm, but pre-attentive processing supports the identification of target information post presentation of the information by assessing the co-occurrence or time-locked EEG potentials. This paper presents a comprehensive review and evaluation of the limited, but significant, literature on research in RSVP-based brain–computer interfaces (BCIs). Applications that use RSVP-based BCIs are categorized based on display mode and protocol design, whilst a range of factors influencing ERP evocation and detection are analyzed. Guidelines for using the RSVP-based BCI paradigms are recommended, with a view to further standardizing methods and enhancing the inter-reliability of experimental design to support future research and the use of RSVP-based BCIs in practice.

Keywords: rapid serial visual presentation, brain–computer interface, event related potentials, electroencephalography, visual evoked potentials

(Some figures may appear in colour only in the online journal)

1. Introduction

Rapid serial visual presentation (RSVP) is the process of sequentially displaying images at the same spatial location at high presentation rates with multiple images per second,

e.g. with a stimulus onset asynchrony no greater than 500 ms but often lower than 100 ms, i.e. >10 stimuli presented per second. Brain–computer interfaces (BCIs) are communication and control systems that enable a user to execute a task via the electrical activity of the user's brain alone (Vidal 1973). RSVP-based BCIs are a specific type of BCI that are used to detect target stimuli, e.g. letters or images, presented sequentially in a stream, by detecting brain responses to such targets. RSVP-based BCIs are considered as a viable approach



Original content from this work may be used under the terms of the [Creative Commons Attribution 3.0 licence](https://creativecommons.org/licenses/by/3.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

to enhance human–machine symbiosis and offers potential for human enhancement.

To date, the literature on RSVP-BCIs has not been comprehensively evaluated, therefore it is timely to review the literature and provide guidelines for others considering research in this area. In this review we: (1) identify and contextualize key parameters of different RSVP-BCI applications to aid research development; (2) document the growth of RSVP-based BCI research; (3) provide an overview of key current advancements and challenges; (4) provide design recommendations for researchers interested in further developing the RSVP-BCI paradigm.

This review is organized as follows: section 2 presents background information on the fundamental operating protocol of RSVP-BCIs. Section 3 details results of a bibliometric analysis of the key terms ‘rapid serial visual presentation’, ‘RSVP’, ‘electroencephalography’, ‘EEG’, ‘brain–computer interface’, ‘BCI’, ‘event-related potentials’, ‘ERP’ and ‘oddball’ found within authoritative bibliographic resources. Section 4 provides an overview of performance measures. Section 5 outlines existing RSVP-based BCI applications, presenting inter-application study comparisons, and undertakes an analysis of the design parameters with inter-application study comparisons. Section 6 provides a summary, discussion of findings and ongoing challenges.

2. Background

RSVP-based BCIs have been used to detect and recognize objects, scenes, people, pieces of relevant information and events in static images and videos. Many applications would benefit from an optimization of this paradigm, for instance counter intelligence, policing and health care, where large numbers of images/information are reviewed by professionals on a daily basis. Computers are unable to analyze and understand imagery as successfully as humans and manual analysis tools are slow (Gerson *et al* 2005, Mathan *et al* 2008). In studies carried out by Sajda *et al* (2010), Poolman *et al* (2008) and Bigdely-Shamlo *et al* (2008), a trend of using RSVP-based BCIs for identifying targets within different image types has emerged. Research studies show the ability to use RSVP-based BCIs to drive a variety of visual search tasks including, in some circumstances, skills learned for visual recognition. Although the combination of RSVP and BCI has proven successful on several image sets, other research has attempted to establish whether or not greater efficiencies can be reached through the combination of RSVP-based BCIs and behavioral responses (Huang *et al* 2007).

2.1. Event related potentials and their use in RSVP-based BCIs

Event-related potentials (ERPs) are electroencephalography (EEG) signal amplitude variations in the EEG associated with the onset of a stimulus (usually auditory or visual) presented to a person. ERPs are typically smaller in amplitude ($<10 \mu\text{V}$) in comparison to the ongoing EEG activity ($\sim 50\text{--}100 \mu\text{V}$)

they are embedded within (Huang *et al* 2008, Acqualagna and Blankertz 2011). As ERPs are locked in phase and time to specific events, they can be measured by averaging epochs over repeated trials (Huang *et al* 2011, Cecotti *et al* 2012, 2014). Shared EEG signal features are accentuated and noise attenuated (Luck 2005, Cohen 2014). The outcome is represented by a temporal waveform with a sequence of positive and negative voltage deflections labeled as ERP components. ERPs are representative of summated cortical neural processing and behavioral counterparts, such as attentional orientation (Wolpaw and Wolpaw 2012, Cohen 2014).

The stream of images presented within an RSVP paradigm comprise frequent non-target images and infrequent target images; different ERP components are associated with target and non-target stimuli (Bigdely-Shamlo *et al* 2008, Cohen 2014, Sajda *et al* 2014). BCI signal processing algorithms are used to recognize spatio-temporal electrophysiological responses and link them to target image identification, ideally on a single trial basis (Manor *et al* 2016).

The most commonly exploited ERP in RSVP-based BCI applications is the P300. The P300 appears at approximately 250–750 ms post target stimulus (Polich and Donchin 1988, Leutgeb *et al* 2009, Ming *et al* 2010, Zhang *et al* 2012). As specified by Polich and Donchin (1988) during the P300 experiment (commonly referred to as the ‘oddball’ paradigm), participants must classify a series of stimuli that fall into one of two classes: targets and non-targets. Targets appear more infrequently than non-targets (typically $\sim 5\text{--}10\%$ of total stimuli in the RSVP paradigm) and should be recognizably different. It is known that P300 responses can be suppressed in an RSVP task if the time between two targets is $<0.5\text{ s}$, which is known as attentional blink (Raymond *et al* 1992, Kranczioch *et al* 2003). The amplitude and the latency of the P300 are influenced by the target discriminability and the target-to-target interval in the sequence. The latency of the P300 is affected by stimulus complexity (McCarthy and Donchin 1981, Luck *et al* 2000). The P300 amplitude can vary as a result of multiple factors (Johnson 1986), such as:

- subjective probability—the expectedness of an event;
- stimulus meaning—comprised of task complexity, stimulus complexity and stimulus value;
- information transmission—the amount of stimulus information a participant registers in relation to the information contained within a stimulus.

2.2. RSVP-based BCI amongst the BCI classes

BCIs can be of three different types: active, reactive or passive (Zander *et al* 2010). An active BCI is purposefully controlled by the user through intentional modulation of neural activity, often independent of external events. Contrastingly, reactive BCIs generate outputs from neural activity evoked in response to external events, enabling indirect control by the user. Passive BCI makes use of implicit information and generate outputs from neural activity without purposeful control by the user. Active/reactive BCIs are commonly aimed at users with restricted movement abilities who intentionally try

to control brain activity, whereas implicit or passive BCIs are more commonly targeted towards applications that are also of interest to able-bodied users (Zander and Kothe 2011, Sasane and Schwabe 2012).

2.3. RSVP-based BCI presentation modes

RSVP-based BCIs have two presentation modes: static mode in which images appear and disappear without moving; and moving mode where targets within short moving clips have to be identified (Sajda *et al* 2010, Cecotti *et al* 2012, Weiden *et al* 2012). Both presentation modes can be used with or without a button press. With a button press, users indicate manually, by pressing a button, when they observe a target stimulus. A button press is used to establish baseline performance, reaction time and/or to enhance performance (discussed further in section 5.1).

2.3.1. Static. In ‘static mode’, images displayed have identical entry and exit points—the images are transiently presented on screen (typically for 100–500 ms) and then disappear. One benefit of static mode is that images occupy the majority of the display and, therefore, identification of targets is likely even if they are only presented briefly. There are a number of different possible instructions a participant may be given.

- Prior to presentation, a target image may be shown to participants and participants are asked to identify this image in a sequence of proceeding images. Target recognition success rates can be achieved with presentation rates as high as 10 s^{-1} (Cecotti *et al* 2012).
- Participants may be asked to identify a *type of target* e.g. an animal within a collection of images. In this mode, the rate of presentation should be slowed down (4 s^{-1}) (Wang *et al* 2009).
- Immediately after image sequence presentation, the participant may be shown an image and asked: ‘did this image appear in the sequence you have just seen?’ (Potter *et al* 2002).

2.3.2. Moving. There has been relatively little research regarding neural signatures of a target and/or anomalies in real world or simulated videos. In ‘moving mode’, short video clips are shown to participants, and within one video clip participants may be asked to identify one or more targets. It is important that these targets are temporally ‘spread out’ to avoid P300 suppression. There are different possible instructions a participant may be given:

- Prior to presentation, participants may be given a description of a target, i.e. asked to identify, say a ‘person’ or ‘vehicle’ in a moving scene (Weiden *et al* 2012).
- Participants can be asked to identify a target event; in this case, the target is identified across space and time. The participant is required to integrate features from both motion and form to decide whether a behavior constitutes a target, for example, Rosenthal *et al* (2014) defined the target as a person leaving a suspicious package in a train station.

2.4. Cognitive blindness

When designing an RSVP-based BCI, three different types of cognitive blindness should be considered namely, the attentional blink, change blindness and saccadic blindness. Generally, RSVP is a paradigm used to study the *attentional blink*, which is a phenomena that occurs when a participant’s attention is grabbed by an initial target image and a further target image may not be detectable for up to 500 ms after the first (Raymond *et al* 1992). Depending upon the duration of stimuli presentation the ratio of target images/total images will change (e.g. if images are being presented at a duration of 100 ms then there must be a minimum of five images between targets 1 and 2. In a sequence of 100 images there can be a maximum of 20 target images. Whereas if images are presented at 200 ms this limits the maximum number of targets to 10/100 images in total).

Change blindness occurs when a participant is viewing two images that vary in a non-trivial fashion, and has to identify the image differences. Change blindness can occur when confronted by images, motion pictures, and real world interactions. Humans have the capacity to get the gist of a scene quickly but are unable to identify particular within-scene features (Simons and Levin 1997, Oliva 2005). For example, when two images are presented for 100 ms each and participants are required to identify a non-trivial variation as the images are interchangeably presented, participants can take between 10–20 s to identify the variation. This latency period in identifying non-trivial variations in imagery can be augmented through use of distractors or motion pictures (Rensink 2000). In the context of designing an RSVP paradigm change blindness is of interest, as it will take longer for a user to identify a target within an image if it does not pop out from the rest of the image. Distractors within the image or cluttered images, will increase the time it takes a user to recognize a target, reducing the performance of the RSVP paradigm.

Saccadic blindness is a form of change blindness described by Chahine and Krekelberg (2009) where ‘*humans move their eyes about three times each second. Those rapid eye movements called saccades help to increase our perceptual resolution by placing different parts of the world on the high-resolution fovea. As these eye movements are performed, the image is swept across the retina, yet we perceive a stable world with no apparent blurring or motion*’. Saccadic blindness thus refers to the loss of image when a person saccades between two locations. Evidence shows that saccadic blindness can occur 50 ms before saccades and up to 50 ms after saccades (Diamond *et al* 2000). Thus, it is important that stimuli have a duration greater than 50 ms to bypass saccadic blindness, unless participants are instructed to attend a focus point and the task is gaze independent and thus does not demand saccades (such as during the canonical RSVP paradigm (section 5.4)).

Having considered some of the factors influencing RSVP-based BCI designs, the remainder of the paper focuses on a bibliometric study of the RSVP literature highlighting the key methodological parameters and study trends. Studies are

compared and contrasted on an intra- and inter-application basis. Later sections focus on study design parameters and provide contextualized recommendations for researchers in the field.

3. Bibliometric study of the RSVP related literature

A bibliometric review of the RSVP-based BCIs was conducted. The inclusion criteria for this review were studies that focused on EEG data being recorded while users were performing visual search tasks using an RSVP paradigm. The studies involved various stimulus types presented using the RSVP paradigm where participants had to identify target stimuli. All reported studies were not simply theoretical and had at least one participant. One or more of the keywords BCI, RSVP, EEG or ERP appeared in the title, abstract or keyword list. Only papers published in English were included. The literature was searched, evaluated and categorized up until August 2017. The databases searched were Web of Science, IEEE, Scopus, Google Scholar, and PubMed. The search terms used were: ‘rapid serial visual presentation’, ‘RSVP’, ‘electroencephalography’, ‘EEG’, ‘brain–computer interface’, ‘BCI’, ‘event-related potentials’, ‘ERP and ‘oddball’.

Papers were excluded for the following reasons: 1. the research protocol had insufficient detail; 2. key aspects needed to draw conclusive results were missing; 3. the spectrum of BCI research reported was too wide (i.e. review papers not specific to RSVP), 4. a ‘possible’ research application was described but the study was not actually carried out; 5. the study was a repeated study by original authors with only minor changes. Due to the immaturity of RSVP-based BCI as research topic, conference papers were not excluded. Inclusion of conference papers was considered important in order to provide a comprehensive overview of the state-of-the-art and trends in the field. Fifty-four papers passed initial abstract/title screening; these were then refined to the 45 most relevant papers through analysis of the entire paper contents. The date of the included publications ranged from 2003–2017.

The relevant RSVP-based BCI papers are presented in table 1 when a button press was required, and table 2 when no button presses were conducted. RSVP-based BCIs were evaluated in terms of the interface design. Tables 1 and 2 show that there is considerable variation across the different studies in terms of the RSVP-BCI acquisition paradigm, including the total number of stimuli employed, percentage of target stimuli, size of on-screen stimuli, visual angle, stimulus presentation duration, and the number of study participants. Performance was measured using a number of metrics: the area under the receiver operating characteristic (ROC) curve (Fawcett 2006), classification accuracy (%) and information transfer rate. ROC curves are used when applications have an unbalanced class distribution, which is typically the case with RSVP-BCI, where the number of target stimulus is much smaller than that of non-target stimuli. Many studies report different experimental parameters and some aspects of the studies have not been comprehensively reported. From tables 1 and 2, it can be seen that the majority of applications using a button press as a baseline may be classified as surveillance applications while

applications that do not use a button press are more varied. This may be because often surveillance applications have an industry focus, and quantified improvement relative to manual labeling alone is crucial for acceptance. In the majority of the applications where a button press was used, participants undertake trials with and without a button press and the difference in latency of response between the two is calculated to compare neural and behavioral response times. The results of the bibliometric analysis are further discussed in sections 4–6, following the analysis of key papers identified in the following section.

4. Validating inter-study comparison through performance measures

When comparing RSVP studies it is important to acknowledge that researchers use different measures of performance. Before going into depth about signal processing techniques (section 5.7) it is important to discuss, firstly, the variations in approaches used to measure performance. To encourage valid inter-study comparison within and across RSVP application types, it is crucial to emphasize that we are, on the whole, reporting classification accuracy when it is calculated in terms of the number of correctly classified trials. Classification accuracy can be swayed by the imbalanced target and non-target classes, with targets being infrequently presented, e.g. with a 10% target prevalence; if all trials are classed as non-targets, correct classification rate would be 90%. Hence, ROC values are also reported in this review where relevant information was provided in the publications reviewed.

In the literature, there are many variations on how performance is estimated and reported. The studies cited in the current section provide examples of performance measure variations from the literature. The intention of Files and Marathe (2016) with respect to the reference list provided. Please check.] was to develop a regression-based method to predict hit rates and error rates whilst correcting for expected mistakes. There is a need for such methods, due to uncertainty and difficulty in correctly identifying target stimuli. The regression method developed by Files and Marathe (2016), had relatively high hit rates, which spanned 78.4%–90.5% across all participants. Contrastingly, as a measure of accuracy, Sajda *et al* (2010) used hit rates expressed as a fraction of total targets detected per minute. Sajda *et al* (2010) discussed an additional experiment that employed ROC values as an outcome measure. In Alpert *et al* (2014), where the RSVP application was categorization based, accuracy was defined as the number of trials in which the classifier provided the correct response divided by the total number of available trials, with regards to target/non-target classification. Yazdani *et al* (2010) were concerned with surveillance applications of RSVP-based BCI and used the F-measure to evaluate the accuracy of the binary classifier in use. Precision (fraction of occurrences flagged that are of relevant) and recall (fraction of relevant occurrences flagged) were reported, as the F-measure considers both these values.

Different variations in ROC value calculations were also discovered across the studies evaluated. Variability in the

distribution of accuracy outcome measures is also founded upon whether the dataset is non-parametric, e.g. median AUC is reported as opposed to the mean AUC (Matran-Fernandez and Poli 2014). As a measure of accuracy, Rosenthal *et al* (2014) conducted a bootstrap analysis: to show the sampled distribution of AUC values for HDCA classifiers were 1000 times over, labels were randomized, classifiers were trained and AUC values calculated through a ‘leaving-one-out cross-validation’ technique. Cecotti *et al* (2012) presented a comparison of three class classifiers in a ‘one versus all’ strategy. The focus of Cecotti *et al* (2012) was to compare the AUC to the volume under the ROC hyper-surface and the authors found an AUC of 0.878, which is suggestive of the possibility for discrimination between greater than two types of ERPs using single-trial detection. Huang *et al* (2006) reported the AUC for session one of two experiments during button press trials. This paper demonstrates that the three classifiers approach produces a similar performance with AUC of >0.8 across the board (Huang *et al* 2006). Moreover, accuracy reportedly increases through collating evidence from two BCI users, and reportedly yielded a 7.7% increase in AUC compared to a single BCI user (Matran-Fernandez and Poli 2014) using collaborative BCIs. This process was repeated 20 times to achieve an average accuracy measurement that would not be relatable to other studies included in the bibliometric analysis that involved average performance over single trial test. Cecotti *et al* (2011) carried out a study where they compared varying target stimuli probability. Target probability has a significant effect on both behavioral performance and target detection. The best mean AUC is achieved with target probability of 0.10 AUC = 0.82. The best target stimuli probability for optimal detection performance were $5\% = 78.7\%$.

This above review exemplifies how performance measures are used. The variability of accuracy analytics limits the extent to which inter-study comparability is feasible, nonetheless a high proportion of studies use AUC values and percentage accuracy as outcome measures, therefore these measures provide the basis for comparisons in section 5. In the RSVP-based BCI application sections that follow, we provide additional information about the values reported in tables 1 and 2, the intention being to validate why these performance metrics were selected when a number of different results are reported by the specified study, and to highlight inter-study idiosyncrasies that may need to be considered whilst comparing findings. In the next section, the different design parameters for the studies identified in tables 1 and 2 are reviewed and a number of recommendations are suggested for the parameters that should be considered for RSVP-based BCI applications.

5. Design parameters

RSVP-based BCI applications to date can be grouped into surveillance, data categorization, RSVP speller, face recognition and medical image analysis applications. Often EEG-based RSVP-BCI system studies are multifactorial by design and report numerous results in the form of different outcome

measures. In the RSVP-based BCI application section that follows, we provide examples of the different application types and examples of their design parameters.

When designing an RSVP paradigm, there are eight criteria that we recommend be taken into consideration.

- (1) The type of target images and how rapidly these can be detected, e.g. picture, number of words.
- (2) The differences between target and non-target images and how these influence the discrimination in the RSVP paradigm.
- (3) The display mode—static or moving stimuli and the background the images are presented on, e.g. single color white, mixed, textured.
- (4) The response mode—consideration should be given as to whether a button press is used or not to confirm if a person has identified a target.
- (5) The number of stimuli/the percentage of target stimuli—how many are presented throughout the duration of a session and the effect this could have on the ERP.
- (6) The rate at which stimuli are presented on screen throughout the duration of a session and the effect this has on the ERP.
- (7) The area (height \times width), visual angle and the overt or covert attention requirement of the stimuli.
- (8) The signal processing pipeline—determine the features, channels, filters, and classifiers to use.

5.1. Display and response modes

A button press may be used in conjunction with either of the aforementioned presentation modes (section 2.2), and entails users having to click a button when they see a target. This mode is used as a baseline to estimate the behavioral performance and the difficulty of the task. In most research studies, participants undergo an experimental trial without a button press and a follow-on trial with a button press.

A button press can be used in RSVP-based BCI research in combination with the participant’s EEG responses in order to monitor attention (Marathe *et al* 2014). The combination of EEG and button press can lead to increased performance in RSVP-based BCIs. Tasks that require sustained attention can cause participants to suffer from lapses in vigilance due to fatigue, workload or visual distractors (Boksem *et al* 2005). The button press can be used to determine if there is a tipping point during the presentations when participants are unable to consciously detect target stimuli, while still identifying targets via EEG recordings (Potter *et al* 2014). However, the core advantage of the RSVP-based BCIs is the enhanced speed of using a neural signature instead of a behavioral response to determine if a user has detected an intended image of interest.

Forty of the studies reported use static mode as a method of presentation; six of these papers used moving mode in conjunction with static mode while one study exclusively used moving mode. Moving mode is more complex than static mode as participants have to take in an entire scene rather than specific images. Moving mode uses motion onset in conjunction with the P300 for scenes in which the targets are moving,

Table 1. Design parameters reviewed, mode: button press = yes. Table acronyms: SVM (support vector machine), SFFS (sequential forward feature selection), N/A (not available), BLDA (Bayesian linear discriminant analysis), CSP (common spatial pattern), BCSP (bilinear CSP), CCSP (composite CSP), LDA (linear discriminant analysis), C (EEG channel), FDA (fisher discriminant analysis), FDM (finite difference model), LLC (linear logistic classifier), RBF SVM (radial basis function SVM), PCA (principle component analysis), LP (Laplacian classifier), LN (linear logistic regression), SP (spectral maximum mutual information projection), FDA (fisher discriminant analysis), ACSP (analytic CSP), HT (human target), NHT (non-human target), ST (single trial), DT (dual trial), BDA (bilinear discriminant analysis), ABDA (analytic BDA), DCA (directed components analysis), HDCA (hierarchical discriminant component analysis), TO (only background distractors), TN (non-target distractor stimuli and background and target stimuli), TvB (target versus background distractor), Tv[B + NT] (target versus both background distractor and non-target).

	Reference	Mode	Application	Stimuli	Targets (%)	Duration (ms)	Size (px)	Visual angle (°)	Participants	Data analysis	ROC performance	Accuracy (%)
1	Healy and Smeaton (2011)	Static	Categorization	4800	1.25	100	N/A	N/A	8	SVM linear kernel, SFFS	N/A	N/A
2	Cecotti <i>et al</i> (2011)	Static	Categorization/ face recognition	12000 trials	5 10 25 50	500	N/A	N/A	8	XDAWN + BLDA	0.768 ± 0.074 0.821 ± 0.063 0.815 ± 0.068 0.789 ± 0.070	78.7 76.4 77.0 71.5
3	Yu <i>et al</i> (2011)	Static	Categorization	>4000	~1.5	150	N/A	N/A	20	BCSP, SVM CCSP, SVM CSP, SVM	N/A	83.0 ± 8.0 75.4 ± 8.3 71.8 ± 9.9
4	Ušćumlić <i>et al</i> (2013)	Static	Categorization	1382	10	250	N/A		15	Gaussian Ensemble LDA (8C) Ensemble LDA (41C) Gaussian Ensemble LDA (8C) Ensemble LDA (41C) Gaussian Ensemble LDA (8C) Ensemble LDA (41C)	0.66 0.78 0.80 0.75 0.80 0.91 0.65 0.68 0.73	90.0 94.8 90.1
5	Mohedano <i>et al</i> (2015)	Static	Categorization	3000	5	100 or 200	N/A	N/A	8	SVM	$0.564-0.863$	N/A

6	Acqualagnav <i>et al</i> (2010)	Static	RSVP speller	30	User dependent	83 or 133	N/A	1	9	LDA	N/A	~70 ~85–90
7	Touryan <i>et al</i> (2011)	Static	Face recognition	470–480	N/A	500	256 × 320	7 horizontally 9 vertically	22	PCA	0.868–0.991	60.4–92.0
8	Sajda <i>et al</i> (2003)	Static	Surveillance	330	50	200 100 50	768 × 512	12.4 by 15.3	2	Spatial linear discriminator	0.79–0.96 0.74–0.80 0.84–0.79	N/A
9	Gerson <i>et al</i> (2006)	Static	Surveillance	284	2	100	640 × 426	33 ± 3 × 25 ± 35	35	Spatial linear discriminator	N/A	74–96
10	Erdogmus <i>et al</i> (2006)	Static	Surveillance	N/A	50	100 50	N/A	N/A	1	LP LN LP LN SP	0.90/0.95 (100/50 ms) 0.37/0.66 (100/50 ms) 0.87–0.83 0.87–0.82 0.89–0.86	N/A
11	Bigdely-Shamlo <i>et al</i> (2008)	Static	Surveillance	24394	40–60	~83	N/A	1.6 by 1.6	7	Bayes fusion of FDA	0.78–0.95	N/A
12	Poolman <i>et al</i> (2008)	Static	Surveillance	8300	4 or 1	100	500 × 500	2	3	DCA FDM	0.70–0.82	72–84
13	Huang <i>et al</i> (2011)	Static	Surveillance	N/A	~1	60–150	500 × 500	22 × 22	33	RBF SVM Linear SVM LLC	0.848–0.941 0.846–0.927 0.753–0.834	N/A
									4	RBF SVM Linear SVM LLC	0.909–0.961 0.887–0.944 0.625–0.866	N/A
14	Weiden <i>et al</i> (2012)	Static/ moving Moving	Surveillance	2500	2	234	512 × 512	N/A	8	SVM	0.50–0.78 (static) 0.89–1.00 (video) 0.72–0.94 (video) 0.58–0.94 (video) 0.55–0.91 (video)	42 (static) 97 (video) N/A
15	Cecotti <i>et al</i> (2012a)	Static/ moving	Surveillance	30000	10	100	N/A	N/A	15	XDAWN, BLDA	~0.874–0.931 (Static HT) ~0.675–0.937 (Video NHT) ~0.875–0.926 (Video HT)	N/A
16	Cecotti <i>et al</i> (2012b)	Static	Surveillance	300	10	200	683 × 384	~13	10	XDAWN, BLDA	0.837 (ST) 0.838 (DT)	N/A N/A

(Continued)

Table 1. (Continued)

Reference	Mode	Application	Stimuli	Targets (%)	Duration (ms)	Size (px)	Visual angle (°)	Participants	Data analysis	ROC performance	Accuracy (%)
17 Yu <i>et al</i> (2014)	Static	Surveillance	>4472	~1.61.6	150	400 × 400	N/A	22	CSP ACSP BDA ABDA	N/A	83.8 ± 6 85.8 ± 5 87.2 ± 4 89.7 ± 5
18 Marathe <i>et al</i> (2014)	Moving	Surveillance	N/A	10	200	N/A	N/A	15	HDCA Sliding HDCA	0.8691 ± 0.0359 0.9494 ± 0.9610	N/A
19 Marathe <i>et al</i> (2015a)	Static	Surveillance	N/A	5 6	500	960 × 600	36.3 × 22.5	17	XDAWN, BLDA	~0.984 (TvB, TO) ~0.971 (TvB, TN) ~0.959 (Tv[B + NT], TN)	N/A
20 Files and Marathe (2016)	Static/ moving	Surveillance	N/A	10	100	N/A	N/A	15	Linear classifiers	N/A	78.4–90.5
21 Bamgrover <i>et al</i> (2016)	Static	Surveillance	4384	4	200	100 × 50	N/A	19	SVM with Haar-like feature classifier	N/A	>70
22 Marathe <i>et al</i> (2015b)	Static/ moving	Intelligence	N/A	10	100 or 500	N/A	N/A	15	HDCA CSP XDAWN, BLDA	>0.9	>70

Table 2. Design parameters reviewed, mode: button press = no. Table acronyms: FDA (fisher discriminant analysis), N/A (not available), SWFP (spatially weighted fisher linear discriminant—principal component analysis), CNN (convolutional neural network), HDPCA (hierarchical discriminant principal component analysis algorithm), HDCA (hierarchical discriminant component analysis), SVM (support vector machine), RBF (radial basis function) kernel, RDA (regularized discriminant analysis), HMM (hidden Markov model), PCA (principal component analysis), BDCA (bilinear discriminant component analysis), BFBP (bilinear feature-based discriminants), BLDA (Bayesian linear discriminant analysis), SWLDA (step-wise linear discriminant analysis), MLP (multilayer perceptron), LIS (locked in syndrome), CV (computer vision), STIG (spectral transfer with information geometry), MSS (max subject-specific classifier), L1 (ℓ_1 -regularized cross-validation), MV (majority vote), PMDRM (pooled Riemannian mean classification algorithm), AWE (accuracy weighted ensemble), MT (multi-task learning), CALIB (within-subject calibration), RF (random forest), BHCN (Bayesian human vision-computer vision retrieval).

Reference	Mode	Application	Stimuli	Targets (%)	Duration (ms)	Size (px)	Visual angle (°)	Participants	Data analysis	ROC performance	Accuracy (%)
1 Hope <i>et al.</i> (2013)	Static	Medical	166	~1.1	100	189 × 189	N/A	2	FDA	0.75–0.78	N/A
2 Alpert <i>et al.</i> (2014)	Static	Categorization	725	20	90–110	360 × 360	6.5 × 6.5	12	SWFP HDPCA HDCA	0.64–0.85 N/A N/A	66–82 66–81 57–70
			290					4	SWFP HDPCA	0.58–0.99 /0.99 ± 0.55 0.99 ± 0.67	91 N/A N/A
3 Mohamedano <i>et al.</i> (2014)	Static	Categorization	4224	15	200	N/A	N/A	5	SVM RBF	0.63–0.78	N/A
4 Manor and Geva (2015)	Static	Categorization	N/A	20	90–110	360 × 360	6.5 × 6.5	15	SWFP Deep CNN	0.652–0.850 0.692–0.858	70.0–83.1 66.2–82.5
5 Huang <i>et al.</i> (2017)	Static	Categorization	N/A	12.5	200	N/A	N/A	7	LDA + RF BHCN	0.873 0.987	N/A
6 Orhan <i>et al.</i> (2011)	Static	RSVP Speller	26	~3.8	150	N/A	N/A	2	RDA	0.948–0.973	N/A
7 Hild <i>et al.</i> (2011)	Static	RSVP Speller	26	~3.6 (User dependent)	400	N/A	N/A	2 (1 LIS)	RDA	N/A	N/A
8 Orhan <i>et al.</i> (2012a)	Static	RSVP Speller	26	~3.8	150	N/A	N/A	2	RDA HMM	N/A	N/A
9 Orhan <i>et al.</i> (2012b)	Static	RSVP Speller	28	~3.6 (User dependent)	400 or 150	N/A	N/A	3 (1 LIS)	RDA PCA	N/A	Healthy controls = 95 LIS = 85
10 Chennu <i>et al.</i> (2013)	Static	RSVP Speller Matrix P300 Speller	25 25	4 4	133	N/A	N/A	11	SWLDA	0.82 0.84	86.02 88.58
11 Orhan <i>et al.</i> (2013)	Static	RSVP Speller	28	~3.8	150	N/A	N/A	2	PCA RDA	0.812–0.998	N/A
12 Oken <i>et al.</i> (2014)	Static	RSVP Speller	28	~3.6 (semi-user dependent)	400	N/A	3.8	15 (6 LIS)	PCA RDA	Healthy controls = 0.81–0.86 LIS = 0.73–0.92	N/A

(Continued)

19	Manor <i>et al</i> (2016)	Static	Surveillance	N/A	~10	100 or 200	400 × 400	N/A	2	Supervised multimodal network Semi-supervised multimodal network	N/A	88.1–93.9 81.4–90.3
20	Waytowich <i>et al</i> (2016)	Moving/ static	Surveillance	N/A	~11	100	N/A		32	Offline	STIG MSS L1 MV PMDRM STIG AWE MT CALIB	N/A
									17	Real-time feedback	STIG MSS L1 MV PMDR M STIG AWE MT CALIB	
21	Yazdani <i>et al</i> (2010)	Static	Other	52	~2	500	N/A	N/A	5	SVM with radial basis function kernel	N/A	35 ± 10.4– 71.1 ± 9.0 (F-measure range)
22	Huang <i>et al</i> (2017)	Static	Categorization	96	12.5	200	N/A	N/A	7	Adaboost Bagging ANN RF SVM LR	0.873 0.987	0.887
23	Lin <i>et al</i> (2017)		Categorization	2000	10	250	N/A	N/A	7 8	SWLDA HDCA	0.7837–0.9148 0.9082–0.9522	N/A

yielding a more realistic setting to validate RSVP-based BCIs (Weiden *et al* 2012). All papers employing moving mode were found within the surveillance application category; this is unsurprising as the moving mode offers the opportunity to detect targets in realistic surveillance situations where movements of people or vehicles are of interest. For the other application areas, i.e. medical, categorization, etc the static mode is likely to be the most appropriate.

Won *et al* (2017) compared motion RSVP to standard RSVP, with the motion-type RSVP being the rapid presentation of letters of the alphabet, numbers 1–9 and a hyphen ‘-’ used to separate words, in six different color groups in one of six directions in line with the hands of a clock, i.e. 2, 4, 6, 8, 10, 12, whilst participants focused on a central point. An increase in performance accuracy with motion-type RSVP versus static-type was demonstrated, which could be accounted for by the shorter latency and greater amplitudes of ERP components in the motion-type variation (Won *et al* 2017).

Out of the studies found, 22 used a button press while 23 did not. 70% of surveillance applications used a button press. In categorization studies and face recognition studies the majority of applications used a button press. 89% of RSVP-speller applications did not use a button press. Typically, the BCI studies that involve spellers focus on movement-free communication and high information transfer rates. Having a button press for confirmation of targets is not standard practice in such applications (Orhan *et al* 2012, Oken *et al* 2014). In many of the studies that did not utilize a button press, researchers were focused on different aspects of the RSVP paradigm other than reaction time. For example, researchers focused on the comparison of two classification methods, image duration, etc (Sajda *et al* 2010, Cecotti *et al* 2014). Combining EEG responses with button press can improve accuracy although more signal processing is required in order to remove noise that occurs as a result of participant movement (Healy and Smeaton 2011). Button press confirmation is unnecessary unless an assessment of physical reaction time is an important aspect of the study.

Maguire and Howe (2016) instructed participants to use a button press following image blocks to indicate if a target was consciously perceived as present or absent. Such an approach is useful when studying RSVP-based parameters and the limits of perception. However, button press responses might be less useful than EEG responses during RSVP for data labeling or image sorting, where the focus is to label individual images within the burst. Nonetheless, Bigdely-Shamlo *et al* (2008) applied an image burst approach where a button press at the end of the image burst was used to determine if the participant saw a target image or not. The authors showed that airplanes could be detected in aerial shots with image bursts lasting 4100ms and images presented at 12 Hz. The button press served well in determining correct and incorrect responses. In practice, however, a button press may be superfluous or infeasible.

A body of researchers is of the opinion that RSVP-related EEG accuracy must surpass button press accuracy in order to be useful. However, this need not be the case as Gerson *et al*

(2006) report no significant differences in triage performance based on EEG recordings or button presses. Nevertheless button-based triage performance is superior for participants that correctly respond to a high percentage of target images. Conversely, EEG-based triage alone is shown to be ideal for the subset of participants who responded correctly to fewer images Gerson *et al* (2006). Hence, the most reliable strategy for image triaging in an RSVP-based paradigm may be through reacting to the target image by real-time button presses in conjunction with an EEG-based detection method. Target identification reflected in EEG responses can be confirmed by a button press, and through signal processing techniques both reported and missed targets can be identified.

Studies such as Marathe *et al* (2014) proposed methods for integrating button press information with EEG-based RSVP classifiers to improve overall target detection performance. However, challenges arise when overlaying ERP and behavioral responses, such as issues concerning stimulation presentation speed and behavioral latency (Files and Marathe 2016). Crucially Files and Marathe (2016) demonstrated that techniques for measuring real-time button press accuracy start to fail at higher presentation rates. Given evidence of human capacity for semantic processing during 20 Hz image streams (approximately 50 ms per image) and response times (RTs) often being an order of magnitude greater than EEG responses, button presses may be unsuitable for faster RSVP-based image triaging.

Pending further studies investigating the reliability of fast detection of neural correlates, EEG-based responses have the potential to exceed button press. However, it is not necessary for EEG-based RSVP paradigms to surpass button press performance and evidence suggests that a complement of both modalities at comfortable lower presentation rates may indeed be the best approach. Nevertheless, ideally studies would contain an EEG-only block and EEG plus button press block, where the button press follows the target and not the image burst. This would facilitate more accurate evaluation of differences and correlations between behavioral and neural response times. Interesting, Bohannon *et al* (2017), presented a heterogeneous multi-agent system comprising computer vision, human and BCI agents, and showed that heterogeneous multi-agent image systems may achieve human level accuracies in significantly less time than a single human agent by balancing the trade-off between time-cost and accuracy. In such cases a human–computer interaction may occur in the form of button press if the confidence in the response of other, more rapid agents such as RSVP-BCI agents or computer vision algorithm is low for a particular sequence of stimuli.

5.2. Type of stimuli

Surveillance is the largest RSVP BCI system application reported in this review, reflected as such by the discussion length of this subsection (Sajda *et al* 2003, Erdogmus *et al* 2006, Gerson *et al* 2006, Poolman *et al* 2008, Bigdely-Shamlo *et al* 2008, Sajda *et al* 2010, Huang *et al* 2011, Cecotti *et al* 2012, Weiden *et al* 2012, Matran-Fernandez and Poli 2014, Marathe *et al* 2014, Rosenthal *et al* 2014, Yu *et al* 2014,

Marathe *et al* 2015, Barngrover *et al* 2016, Cecotti 2016, Files and Marathe 2016).

In a surveillance application study carried out by Huang *et al* (2011) the targets were surface-to-air missile sites. Target and non-target images shared low-level features such as local textures, which enhanced complexity. Nonetheless target images were set apart due to large-scale features such as unambiguous road layouts. Another example of surveillance targets denoted by Bigdely-Shamlo *et al* (2008) is where overlapping clips of London satellite images were superimposed with small target airplane images, which could vary in location and angle within an elliptical focal area. Correspondingly, in Barngrover *et al* (2016), the prime goal was to correctly identify sonar images of mine-like objects on the seabed. Accordingly, a three-stage BCI system was developed whereby the initial stages entailed computer vision procedures, e.g. Haar-like feature classification whereby pixel intensities of adjacent regions are summed and then the difference between regions is computed, in order to segregate images into image chips. These image chips were then fed into an RSVP type paradigm exposed to human judgment, followed by a final classification using a support vector machine (SVM).

In the categorization, application type images were sorted into different groups (Cecotti *et al* 2011). Alpert *et al* (2014) conducted a study whereby five image categories were presented: cars, painted eggs, faces, planes, and clock faces (Sajda *et al* 2014). A second study by Alpert *et al* (2014), containing target (cars) and non-target image (scrambled images of the same car) categories, was conducted. In both RSVP experiments, the proposed spatially weighted Fisher linear discriminant–principal component analysis (SWFP) classifier correctly classified a significantly higher number of images than the hierarchical discriminant component analysis (HDCA) algorithm. In terms of categorization, empirical grounds were provided for potential intuitive claims, stating that target categorization is more efficient when there is only one target image type, or distractors are scrambled variations of the target image as opposed to different images all together (Sajda *et al* 2014).

Face recognition applications have been used to seek out whether a recognition response can be delineated from an uninterrupted stream of faces, whereby each face cannot be independently recognized (Touryan *et al* 2011). Two of the three studies evaluated utilized face recognition RSVP paradigm spin offs with celebrity/familiar faces as targets and novel, or other familiar or celebrity faces as distractors (Touryan *et al* 2011, Cai *et al* 2013). Cecotti *et al* (2011) utilized novel faces as targets amongst cars with both stimuli types presented with and without noise. Utilizing the RSVP paradigm for face recognition applications is an unconventional approach; nonetheless the ERP itself has been used exhaustively to study neural correlates of recognition and declarative memory (Yovel and Paller 2004, Guo *et al* 2005, MacKenzie and Donaldson 2007, Dias and Parra *et al* 2011). Specifically, early and later components of the ERP have been associated with the psychological constructs of familiarity and recollection, respectively (Smith 1993, Rugg *et al* 1998). There is thus substantial potential for the utility of the RSVP-based BCI paradigm for applications

in facial recognition. In the future, RSVP-based BCI face recognition may be apposite in a real world setting in conjunction with security-based identity applications to recognize people of interest. Furthermore, Touryan *et al* (2011) claimed that, based on the success of their study, RSVP paradigm-based EEG classification methods could potentially be applied to the neural substrates of memory. Indeed, some studies show augmentation in the posterior positivity of ERP components for faces that are later remembered (Paller and Wagner 2002, Yovel and Paller 2004). That is to say, components of ERPs triggered by an initial stimulus may provide an indication of whether memory consolidation of the said stimulus will take place, which provides an interesting avenue for utilizing RSVP-based BCI systems for enhancing human performance. Based on these studies, it is clear that relatively novel face recognition paradigms have achieved success when used in RSVP-based BCIs.

RSVP-based BCIs that assist with finding targets within images to support clinical diagnosis has received attention (Stoica *et al* 2013), for example, in the development of more efficient breast cancer screening methods (Hope *et al* 2013). Hope *et al* (2013) is the only paper evaluated from the field of medical image analysis and hence described in detail. During an initial sub-study participants were shown mammogram images, where target lesions were present or absent. In a subsequent study, target red or green stimuli were displayed among a set of random non-target blobs. These studies facilitated comparison between ‘masses’ and ‘no masses’ in mammograms, and strong color-based images versus random distractors. Images were presented against a grey background in three-second bursts of 30 images (100 ms per image). A difference in the amplitude of the P300 potential was observed across studies, with a larger amplitude difference between target and non-target images in the mammogram study. The researchers attributed this to the semantic association with mammogram images, in contrast to the lack thereof in the colored image-based study.

5.3. Total stimuli number and prevalence of target stimuli

The number of stimuli refers to the total number of stimuli, i.e. the same stimulus can be shown several times. An exception to this is RSVP-speller studies where researchers only report on the number of symbols used, i.e. 28 symbols—26 letters of the alphabet, space and backspace (Hild *et al* 2011). In the RSVP-speller studies reviewed, the number of times each symbol was shown was not explicit. RSVP-speller applications are likely to have significantly fewer stimuli than the other aforementioned applications as participants are spelling out a specific word or sentence, which only has a small number of target letters/words. The integration of language models into RSVP-speller applications enables ERP classifiers to utilize the abundance of sequential dependencies embedded in language to minimize the number of trials required to classify letters as targets or non-targets (Orhan *et al* 2011, Kindermans *et al* 2014). Some systems, such as the RSVP keyboard (described in Hild *et al* (2011), Orhan *et al* (2012a), Oken *et al* (2014)) display only a subset of available characters in each

sequence. This sequence length can be automatically defined or be a pre-defined parameter chosen by the researcher. The next letter in a sequence becomes highly predictable in specific contexts, therefore it is not necessary to display every character in the RSVP speller. Studies show that target characters are generally displayed more than once before the character is selected. The length of a sequence and the ratio of target to non-target stimuli can have an effect on the typing rate/performance. In an online study by Acqualagna *et al* (2011), participants were shown 30 symbols that were randomly shuffled ten times before a symbol was selected through classification and presented on screen. Orhan *et al* (2012), carried out an offline study whereby two healthy participants were shown three sequences (consisting of 26 randomly ordered letters of the alphabet). The results of this study showed that the number of correctly identified symbols more than doubled when using three sequences instead of one sequence to identify targets.

Task complexity is enhanced by the multiplicity of target categories. In Poolman, *et al* (2008) there were two blocks of target presentations: a helipad block with a 4% target prevalence; and a surface-to-air missile and anti-aircraft artillery block with a 1% target prevalence. Additionally, in Cecotti *et al* (2012) the targets were 50% vehicles, 50% people, with 50% being stationary and 50% moving. Further to this, (Weiden *et al* 2012) demonstrated that presenting kinetic images during the RSVP paradigm as opposed to stationary images increased the performance of EEG-based detection, and that this is negatively correlated with the cognitive load associated with the presented stimuli. In RSVP-speller applications task complexity varies based on what instructions participants are given, e.g. (1) participants may be asked to 'spell dog'; (2) 'type a word related to weather'; (3) participants can be given a word bank containing 20 words and asked to 'spell a word found within this word bank'. Half of the RSVP-speller-based BCI studies evaluated involved user-defined sequence lengths (instructions 2 and 3) (Acqualagna *et al* 2010, Hild *et al* 2011, Orhan *et al* 2012, Oken *et al* 2014), while the other half involved users being given a target word/sentence to spell (instruction 1). If a participant has to remember the sentence or how to spell a long or unfamiliar word this can increase the complexity of a task (i.e. dog is much easier to spell than idiosyncrasy) (Primativo *et al* 2016). Note however that these different complexities in instructions are only present for evaluation/training tasks with the RSVP-BCI spellers. For their real use, participants choose themselves what they want to spell. The RSVP-based text application allows the number of stimuli before a target stimulus be reduced (i.e. letters such as 'z' that are less commonly used can be shown less frequently).

Excluding RSVP-speller applications, as it is already known that they do not require the same number of stimuli as the other applications, the number of stimuli used typically varied between studies from approximately 800 in the surveillance application study by Sajda *et al* (2010) to 26 100 in a categorization application study by Sajda *et al* (2014). The most common target stimuli percentage range was 1–10% found in 61% of the studies reviewed, followed by 11–20% then >20%. There are a number of studies that focus specifically on the percentage of target stimuli. In a study by Cecotti

Table 3. Variation of image duration in RSVP studies.

Duration (ms)	Number of studies	Accuracy % range
<100	7	66–93
100–199	22	70–92
200–299	11	70–96
300–399	—	—
400–499	2	85–94
500+	8	78.4–90

et al (2011), researchers investigated the influence of target probability when categorizing face and car images. In this study, researchers used spatially filtered EEG signals as the input for a Bayesian classifier. Using eight healthy participants, this method was evaluated using four probabilities of target stimuli conditions, i.e. 0.05, 0.10, 0.25, or 0.50. It was found that the target probability had an effect on the participant's ability to detect targets and on behavioral performance. The best mean AUC (0.82) was achieved using the 0.1 probability condition. The results showed that the percentage of targets shown in an RSVP paradigm has an effect on participants' performance. As number and percentage of target stimuli used can have an effect on the complexity of a task, it is important to keep the percentage of targets to <10% to evoke the P300 and maximize detection rates. This was proposed to be in line with well-established P3 measures, whereby bigger gaps between target trials reduce peak latency and increase amplitude (Gonsalvez and Polich 2002).

5.4. Duration of stimuli presentation

A key factor of the RSVP paradigm is the rate of presentation, as the focus of this paradigm is presenting data at a rapid rate so that large datasets can be analyzed in short periods. The duration for which stimuli were presented varied from 50 to 500 ms (Sajda *et al* 2003, Touryan *et al* 2011, Cai *et al* 2013). The upper limits for the presentation time of stimuli during the RSVP paradigm is ill-defined in the literature; however we found 500 ms per image to be the maximum RSVP duration used across all RSVP studies. The duration of stimuli typically differs between applications. Table 3 shows that the most common duration of stimuli was between 100–199 ms per image. The quickest duration of 50 ms per image was used in a study by Sajda *et al* (2003) where two participants were asked to identify scenes containing people in natural scenes. In each trial, the duration of the stimulus presentation was decreased from 200 to 100 to 50 ms per image. The results of this study showed that both participants had reduced performance for faster stimulus presentations, i.e. 50 ms. This would suggest that the most suitable duration for RSVP-based BCI applications is 100–200 ms, to balance the trade-off between accuracy and speed.

Overall, these limited findings are suggestive of presentation rates of >10 Hz being infeasible for identification of neural correlates that allow successful identification of targets. Despite the low a participant number in Sajda *et al* (2003), validation for this upper cut-off presentation rate may be provided by Raymond *et al* (1992), where the attentional blink was first described. An RSVP paradigm was undertaken

whereby the participant must register a target white letter in a stream of black letters and a second target 'X' amongst this stream. It was found that if the 'X' appeared within ~100–500 ms of the initial target, errors in indicating whether the 'X' was present or not were likely to be made even when the first target was correctly identified (Raymond *et al* 1992). This is not to say that humans cannot correctly process information presented at >10 Hz. Forster (1970), has shown that participants can process words presented in a sentence at 16 Hz (16 words per second). However, the sentence structure may have influenced the correct detection rate, which has an average of four words per second for simple sentence structures and three words for complex sentences. Detection rates improve when presented at a slower pace, e.g. four relevant words per second, with masks (not relevant words) presented between relevant words. Additionally, Fine and Peli (1995) showed that humans can process words at 20 Hz in an RSVP paradigm.

Potter *et al* (2014) assessed the minimum viewing time needed for visual comprehension using RSVP of a series of 6 or 12 pictures presented at between 13 and 80 ms per picture, with no inter-stimulus interval. They found that observers could determine the presence or absence of a specific picture even when the pictures in the sequence were presented for just 13 ms each. The results suggest that humans are capable of detecting meaning in RSVP at 13 ms per picture. However, the finding challenges established feedback theories of visual perception. Specifically, research assert that neural activity needs to propagate from the primary visual cortex to higher cortical areas and back to the primary visual cortex before recognition can occur at the level of detail required for an individual picture to be detected, Maguire and Howe (2016). Maguire and Howe (2016) supported Potter *et al* (2014) in that the duration of this feedback process is likely ≥ 50 ms, and suggest that this is feasible based on work done by Lamme and Roelfsema (2000). Explicitly, Lamme and Roelfsema (2000) estimated that response latencies at any hierarchical level of the visual system are ~10 ms. Therefore, assuming that a minimum of five levels must be traversed as activity propagates from the V1 to higher cortical areas and back again, this feedback process is unlikely to occur in <50 ms. However, Maguire and Howe (2016) suggested a potential confound of Potter *et al* (2014), which was that pictures in the RSVP sequence, on occasion, contained areas with no high-contrast edges and hence may not have adequately masked preceding pictures. Consequently, Maguire and Howe (2016) replicated the study rectifying the edges to ensure high-contrast covering the entire image. They were unable to find any evidence that meaning can be detected in an RSVP stream at 13 ms, or even 27 ms, per image but at 53 and 80 ms this is possible. Upon this basis, the limits of RSVP processing could be reduced to a minimum of ~20 Hz. Nonetheless, further study is needed to investigate the limits of human capability to rapidly distinguish target from non-target information, in comparison to the limit in detecting target related ERPs versus non-target ERPs at 20 Hz presentation rates.

In all three face recognition studies, each face image was displayed for 500 ms (Cecotti *et al* 2011, Touryan *et al* 2011, Cai *et al* 2013). In two of the studies there was no ISI (Cecotti *et al* 2011, Touryan *et al* 2011), and in the other an ISI of

500 ms was given to ensure ample time for image processing (Cai *et al* 2013). The speed at which face images were shown was reduced in comparison to the other RSVP applications. RSVP spellers most commonly use a duration of 400 ms; RSVP-spellers can benefit from slower stimulus duration with the incorporation of a language model to enable the prediction of relevant letters. The estimation of performance can be challenging in the RSVP paradigm when the ISI is small, as assigning a behavioral response (i.e. button press) to the correct image cannot be done with certainty. A solution to this problem is to assign behavioral responses to each image, therefore researchers are able to establish hits or false alarms (Touryan *et al* 2011). When two targets are temporally adjacent with a SOA of 80 ms, participants are able to identify one of the two targets but not both. SOA should be at least 400 ms and target images should not be shown straight after each other (Raymond *et al* 1992). Acqualagnav *et al* (2010), had a four factorial design looking at classification accuracy when the letters presented as no-color or color letters at either 83 or 133 ms with an ISI of 33 ms (Acqualagnav *et al* 2010). The number of sequence stimuli was presented for enhanced accuracy rate in selecting letter of choice. After 10 sequences ~90% mean accuracy was reached in 133 ms color presentation mode (100% for 6/9 participants). After ten sequences in 133 ms no color presentation mode ~80% mean accuracy was reached (100% in 3/9 participants). Whilst at presentation rates of 83 ms mean accuracy rate was ~70% and there was no significant effect of color. This formulation is based on the chance rate of 3.33% (i.e. 1 in 30). This implies that cultured letters enhances performance accuracy but not past a certain speed of stimulus presentation.

There is likely a significant interaction between the difficulty of target identification and presentation rate. For example, the optimal presentation rate for a given stimulus set is highly dependent on the difficulty of identifying targets within that set (Ward *et al* 1997). Image sets with low clutter, high contrast, no occlusion, and large target size are likely amenable to faster presentation rates; while image sets with high clutter, low contrast, high levels of occlusion, with small target sizes will require slower presentation rates (Rousselet *et al* 2004, Serre *et al* 2007, Hart *et al* 2013, Liu and Kwon 2016). A more conclusive analysis of the effect of stimulus presentation duration for each application type could be derived by varying the presentation rate duration between 100, 200, and 500 ms, whilst other parameters remain fixed. With regards to temporal proximity of target images, 500 ms should be taken to be the minimum to maximize performance.

5.5. Image size/visual angle

Another RSVP design aspect to be considered is stimulus size. There is a large variation in image sizes ranging from 256×256 pixels in a categorization application to 960×600 pixels in a surveillance applications. In general, surveillance applications use larger images than the other applications described. The most common image size used is 500×500 pixels. This is only used in static surveillance applications and all surveillance studies using this image size achieved a high

accuracy (>80%). The other applications used smaller image sizes such as 360×360 pixels and achieved high accuracies (i.e. 91% and 89.7%). Therefore, it can be concluded that for surveillance studies, image size should be at least 500×500 pixels, although for all other applications the image size may be smaller. A more complex task is where a target stimulus is presented in the background of a larger image eliciting the N2 ERP. Early components such as the P1 and N2 are sensitive to the spatial location of the stimuli (Saavedra and Bougrain 2012).

One issue with reporting only image size is that it is always relevant to the distance viewed from screen and its location on the screen with respect to the viewer, i.e. the visual angle. The visual angle is the angle an image subtends at the eye, reported in degrees of arc. In a study by Dias and Parra (2011) it was shown that participants performed best (90%) when the target stimulus was centered. Performance consistently decreased to 50% in all participants as target stimulus were placed further away from the center (4° of visual angle), this dropped further when target stimulus was placed at 8° of visual angle. Although performance drops significantly participants are still able to detect target stimulus shown in their peripheral visual field even at such rapid paces. Many papers report that the visual angle of the stimuli can have an effect on performance. As a general principle, targets must appear larger or be more distinct for detection at the outer edge of the visual field. The visual angle can thus be deemed the most important measure as it accounts for distance from screen, image location on screen and image size. Authors are therefore encouraged to report visual angle, as reporting image size alone is not useful without the availability of distance from the screen. For RSVP-speller studies, none of the papers found reported on the size of the image or font, however some reported the visual angle.

5.6. Target versus non-target stimuli

Many different types of target images have been identified within this review. The majority of research focuses on a two-class problem, i.e. detecting target images in sequences of non-target images that are completely different from each other. However, in real-life situations, non-target images are likely to share some of the same characteristics as target images (Marathe *et al* 2015). These presentation sequences appear to be more like moving images than static images. In Marathe *et al* (2015) a more complex surveillance task was carried out where, in the first task, participants were required to detect targets when targets are the only infrequent image whilst, in the second task, targets were presented with non-targets (i.e. the target image could be found in the background of a larger image). Participants were required to ignore everything else in the image, a much more difficult task, and consequently the amplitude of the P300 was reduced. The results of this study found that the introduction of the infrequent non-target stimuli in the scene yielded a substantial slowing of the reaction time. Surveillance applications commonly use stimuli that are more complex where trained participants, such as intelligence analysts, outperform novice participants, as they are able to

give meaning to the stimuli. The RSVP-speller applications present their letters as images one at a time on screen (Hild *et al* 2011). Due to the nature of the RSVP paradigm, it is important that these letters are shown in a random order as participants pre-empting a target can have an effect on ERP responses (Oken *et al* 2014). Data categorization applications had the most variance between the different types of stimuli presented to a participant. However, these stimuli tend to be everyday items that participants can easily recognize.

5.7. Signal processing

All applications have certain requirements in terms of speed and type of images displayed, which, as outlined above, can influence the ERP and therefore also variations in performance as measured by detection accuracy. The signal processing framework plays an important role in being able to cope with variations in ERP and maximizing performance. There is a likely tradeoff between the design parameters used as described above and the level of sophistication built into the signal processing framework, which often varies across studies. Here we review some of the approaches applied.

5.7.1. Pre-processing. To extract the relevant features, data is first pre-processed to improve the signal to noise ratio (SNR). The signal is pre-processed using varying band pass filters, depending on the application, in order to remove high frequency noise or artifacts (such as muscle activity). Generally, lower and upper cut-off frequencies of around 0.1 Hz and 30–40 Hz are used, respectively. The data is then often down-sampled, and, for offline analyses, electrodes with substantial noise are removed through visual inspection of the EEG data or automated approaches based on thresholding or correlating artifacts in EEG channels with simultaneously recorded electrooculography or electromyography. Data is then epoched into segments typically lasting ~600 ms, from 100 ms prior to stimulus onset and the 500 ms post-stimulus onset. The starting point and duration of the epochs selected for further analysis vary from study to study.

5.7.2. Feature extraction. Feature extraction is applied to the data for dimensionality reduction and to extract discriminant and non-redundant features. It can be difficult to carry out feature extraction due to the low SNR in single trial analysis. Conventionally, averaging over multiple repeated trials is often used to overcome this. Many studies employ spatial filtering to extract ERPs from EEG. Some of the spatial filtering methods used include principal component analysis (PCA) (Sajda *et al* 2003, Alpert *et al* 2014), independent component analysis (Bigdely-Shamlo *et al* 2008, Blankertz *et al* 2011, Kumar and Sahin 2013), or the xDAWN algorithm, which maximizes the SNR between target and non-target stimuli classes (Rivet *et al* 2009, Rivet and Souloumiac 2013, Cecotti *et al* 2014). In the case of image triage where the intention is to classify single-trial ERPs, spatial filters are used to enhance the SNR and exploit spatial redundancy (e.g. Parra *et al* (2005)). Yu *et al* (2011) went a step further by utilizing a methodology that considers spatial and temporal features to ensure augmented

Table 4. Parameter and recommendations for RSVP-based BCIs.

Parameter	Surveillance	RSVP-speller	Face recognition	Categorization/medical
Stimuli number	>5000	>5000	2000	>4000
% targets	~5–10	≤5	~10	10–25
Stimulus presentation duration (ms)	100–200	500	500	100–200
Target examples	Helipads, planes, vehicles, people, etc	Letters	Faces	Animals, mammograms, etc
ERP component	P300	P300	N170	P300
Feature extraction	XDAWN	—	XDAWN	BCSP/XDAWN
Classifier	BLDA, SVM, LP, SP	RDA/ SWLDA	SVM	BLDA

single-trial detection accuracy (Yu *et al* 2011). A bilinear common spatial pattern (BCSP) was suggested to outperform common spatial pattern (CSP) filters (composite and common spatial pattern filters) (Yu *et al* 2011). It should be noted however that CSP spatial filters were not designed to classify ERP but to classify oscillatory EEG activity. CSPs, indeed, ignore the EEG time course—i.e. the ERP—and are thus suboptimal for RSVP-BCI. We would recommend using spatial filters dedicated to ERP classification, such as xDAWN, which were used successfully in many RSVP-BCI. Spatial filtering is normally only performed on high-density EEG data, which might be impractical in certain real-life applications (Parra *et al* 2005). High-density EEG data has been reported to increase accuracy (Ušćumlić *et al* 2013). Table 4 shows the most common method used for different application types.

Face recognition applications differ from other applications as face images evoke different ERPs, in addition to the P300. Faces typically evoke a N170 component that changes between targets and non-targets (Maurer *et al* 2008, Luo *et al* 2010). The vertex positive potential is also associated with face recognition (Zhang *et al* 2012). The midfrontal FN400 and later parietal FP600 components have been associated with familiarity and recollection, respectively (MacKenzie and Donaldson 2007). Specifically, the amplitude of FP600 (a positive deflection >500 ms post-stimulus) was found to significantly correlate with the extent of face familiarity (Touryan *et al* 2011). The use of spatial filters that utilize spatial and temporal features may act as an advantage over conventional spatial filters that only exploit spatial redundancy, e.g. Yu *et al* (2011). However, spatial filters can only be performed on high-density EEG data, which might be impractical in certain real-life applications (Parra *et al* 2005).

5.7.3. Classification. This review found many different classification methods had been used in the acknowledged studies, however some conclusions can be drawn. Linear classifiers are most populous within RSVP-based BCIs. Often EEG can contain information that enables classification of the stimuli correctly even when a participant's behavioral response is incorrect (Sajda *et al* 2003, Bigdely-Shamlo *et al* 2008). The two most commonly used classifiers were linear discriminant analysis (LDA) and SVM, or variations of the two, such as Bayesian LDA (BLDA) and radial basis function SVM, respectively. Parra *et al* (2008) presented an RSVP framework that projects the EEG data matrix bi-linearly onto temporal

and spatial axes (Parra *et al* 2008). This framework is versatile upon implementation, for example, it has been applied to classify target natural scenes and satellite missile images (Gerson *et al* 2006, Sajda *et al* 2010). Contrastingly, Alpert *et al* (2014) presented a two-step linear classifier, which achieved classification accuracy suited to real-world applications (Sajda *et al* 2014). Whilst Sajda *et al* (2010) proposed a two-step system utilizing computer vision and EEG subsequently to optimize the classification (Sajda *et al* 2010). The performance of an ensemble LDA classifier diminished when eight centro-parietal EEG channels were utilized as opposed to the full 41 EEG channels (Ušćumlić *et al* 2013). Contrastingly, Healy and Smeaton (2011) claimed that consideration of additional channels might introduce noise as opposed to advancing categorical information, as indicated by results from one study participant.

For the surveillance application, SVM achieved the highest percentage accuracies (Huang *et al* 2011, Weiden *et al* 2012). For the RSVP-speller application, the most common method of classification used was regularized discriminant analysis (RDA). RDA achieved an AUC performance of 0.948–0.973 (Orhan *et al* 2011). Step-wise LDA (SWLDA) was also used in RSVP-speller applications with high AUC performance and accuracies (0.82, 0.84, 86%, 89%) (Hope *et al* 2013). In face recognition applications, the best AUC performance was produced using an SVM classifier (Cai *et al* 2013). Within this review, only one medical application was identified (Hope *et al* 2013) and the researchers had achieved high accuracy using a Fisher discriminant analysis. BLDA classifiers were also used, achieving high levels of accuracy (79%). The SWFP algorithm outperformed the HDCA algorithm by 10% in categorization applications. Touryan *et al* (2011) demonstrated that EEG classification methods applied to categorization procedures can be adapted to rapid face recognition procedures (Touryan *et al* 2011). Window sizes post stimulus onset of 128, 256 and 512 ms were fed into the classifiers. AUC values (average AUC = 0.945) were reported for the customized PCA models utilized to describe the changes in ERPs seen between familiar (famous and personal) and novel faces displayed for 500 ms at a time. It is the customized version of these models, i.e. the models developed for each participant using only that participant's data, which were shown to improve classification performance through the acknowledgment of discrete variability in the windowed ERP components.

Many of the BCI algorithms presented in tables 1 and 2 are linear, enabling simple/fast training with resilience to overfitting often caused by noise, implying suitability to single-trial EEG data classification. Nonetheless, linear methods can limit feature extraction and classification, and non-linear methods, e.g. neural networks, are more versatile in modeling data of greater variability, also implying suitability to single-trial EEG data classification (Erdogmus *et al* 2006, Huang *et al* 2006, Lotte *et al* 2007). The use of neural networks, in particular deep neural network for the RSVP-based BCI framework, represents an attractive venture, and has shown promise over standard linear methods (Manor *et al* 2016, Huang *et al* 2017). A convolution neural network was shown to outperform a two-step linear classifier using the same dataset (Sajda *et al* 2014, Manor and Geva 2015).

The majority of studies reviewed investigated the effectiveness of classifiers in identifying single-trial EEG correlates for target stimuli presented through an RSVP-type paradigm. However, the spatial filtering technique, as well as the type of classifier used, has an impact on proficiency in detecting EEG of single trials (Bigdely-Shamlo *et al* 2008, Cecotti *et al* 2014). For example, independent component analysis reportedly identifies and divides multiple classes of non-brain response artifacts associated with eye and head movements, which would be useful for EEG de-noising during real-world applications when operators are mobile (Bigdely-Shamlo *et al* 2008).

Additionally Cecotti *et al* (2014) evaluated three classifiers using three different spatial filtering methods, so all in all twelve techniques were compared for three different RSVP paradigms. Marathe *et al* (2015) utilized an active learning technique in a bid to reduce the training samples required to calibrate the classifier. Active learning is a partially supervised iterative learning technique reducing the amount of labeled data required for training. Recalibration depends on parameters such as human attentiveness, physical surroundings or task-specific factors. Looking at the real world applicability of RSVP-based BCI systems, Marathe *et al* (2015) built upon work addressing the issue of the thorough recalibration required for real-time BCI system optimization.

There is growing interest in the use of transfer learning (TL) for calibration reduction or suppression to encourage the real-world applicability of BCIs (Wang *et al* 2011). With TL, the EEG data or classifiers from a given domain are transformed in order to be applied to another domain, hence transferring data/classifiers from one domain to another, possibly increasing the amount of data for the target domain (Wang *et al* 2011). For RSVP-BCI, this typically consists in combining EEG data or classifiers from different participants, in order to classify EEG data from another participant, for which very little or even no calibration EEG data is available. An unsupervised transfer method, namely spectral transfer with information geometry (STIG), ranked and collated unlabeled predictions from a group of information geometry classifiers, which was established through training on individual participants (Waytowich *et al* 2016). Waytowich *et al* (2016) showed that STIG can be used for single-trial detection in ERP-based BCIs, eliminating the requirement for taxing data collection

for training. With access to limited data, STIG outperformed alternative zero-calibration and calibration reduction algorithms (Waytowich *et al* 2016). Within the BCI community conventional TL approaches still necessitate training for each condition, however methodologies have been applied to eradicate the need for subject-specific data calibration, where large-scale data is leveraged from other participants (Wei *et al* 2016). This demarcates the potential for single-trial classification via unsupervised TL and user-independent BCI technology deployment.

5.8. Suggested parameters

The parameters reviewed here have been selected because they have an effect on one or all of the following aspects of the RSVP paradigm: task complexity, stimulus complexity, stimulus saliency or information transmission. Performance within RSVP-based BCIs is measured as the participant's ability to correctly identify oddball images in a sequence. RSVP-based BCIs use two different measurements of performance such as accuracy (percentage of targets that are correctly identified using EEG) and ROC curves. 10% of papers assessed in this review did not report at least one out of these performance measures (ROC/percentage accuracy). The accuracies of the different studies need to be put in context, as all the reviewed parameters and other observed parameters i.e. number of trials and participants will influence study accuracy. In table 4 parameter recommendations are provided for designing RSVP-based BCIs within the different application types and these have been discussed thoroughly throughout section 5. In particular, table 4 suggests the parameters to use for each application, according to those leading to the best detection performances (accuracy or AUC) in studies comparatively. If no formal comparisons between parameters were available for a specific application or parameter, the most popular parameter values that yield good performances are mentioned.

Applying BCI systems commercially and outside the lab in real-world scenarios will ideally require the system to be robust during the execution of tasks of increasing difficulty. Section 5 summarized the five applications areas that have been studied to the greatest extent in the context of RSVP-based BCIs. Specifically, this section tackles intra-application comparisons of various aspects of the papers that met the inclusion/exclusion criteria. A few of the papers found in this review carried out more than one study in different application types. The most common type of application found was surveillance applications, followed by RSVP-speller applications and categorization applications; after this were face recognition and lastly medical applications. Although there is a relatively limited number of studies, the design parameters and the focal points of different applications vary widely.

6. Discussion and conclusions

With the increasing intensity in RSVP-based BCI research there is a need for further standardization of experimental protocols, to compare and contrast development of the

different applications described in this review. This will aid the realization of a platform that researchers can use to develop RSVP paradigms and compare their results and determine the optimal RSVP-based BCI paradigm for their application type. This paper presents a review of the available research, the defining elements of the research and a categorization approach that will facilitate coordination efforts among researchers in the field. Research has revealed that using a combination of RSVP with BCI technology allows the detection of targets at an expedited rate without detriment to accuracy.

Understanding the neural correlates of visual information processing can create symbiotic interaction between human and machine through BCIs. Further development of RSVP-based BCIs will depend on both basic and applied research. Within the last five years there have been advancements in how studies are reported, and a sufficient body of evidence exists in support of the development and application of RSVP BCIs. However, there is a need for the research to be developed further, and standardized protocols applied, so that comparative studies can be done for progressive research. Many ERP reviews have been carried out; however, this paper focuses on RSVP visual search tasks with high variability in targets and the parameters used. This paper gives guidelines on which parameters impact performance but also on which parameters should be reported so that studies can be compared. It is important that the design aspects shown in tables 1 and 2 are reported and described within each research study. It has been shown that RSVP-based BCIs can be used in processing target images in multiple application types with a low-target probability, but consistency of reporting method renders it difficult to truly compare one paradigm to another or one parameter setup to another.

There has been profuse reporting of percentage accuracy and area under the ROC curve values, nonetheless there is room for more studies to utilize this unofficial standardization across RSVP-based BCI research.

To maximize reliability to pre-existing literature in terms of keeping one feature that contributes to cognitive load constant, it is recommended that studies utilizing more than one category type as targets to conduct the same study with just one target category in the first instance.

For all applications, it is of course necessary to choose an epoch for single trial ERP classification corresponding to the temporal evolution of the most robust ERP components that are, on the whole, pre-established in the literature as associated with the specified task at hand, i.e. target stimuli identification due to their infrequency, recognizability, relevancy or contents. However, whether the duration of stimuli presentation must extend beyond the latency between ERP component appearances relative to stimuli presentation is questionable.

This review found a single medical application. More research in applying the RSVP-based BCI paradigm to high throughput screening within medicine is highly encouraged upon the basis that similarly complex imagery has been categorized relatively successfully in other applications, e.g. side scan sonar imagery of mines or aircraft amongst birds eye view of maps in surveillance (Bigdely-Shamlo *et al*

2008, Barngrover *et al* 2016). The medical application of RSVP-based BCIs has immense potential in diagnostics and prognostics through recognition and tracking of established disease biomarkers, and accelerating high throughput health image screening.

Studies utilized varying image sizes, visual angles and participant distance from the screen. Researchers are encouraged to report visual angle as it accounts for both image size and distance of the participant from the screen. A potential way to facilitate uniformity of these variables is to utilize a head mounted display (HMD) or virtual reality (VR) headset such as an oculus rift (Foerster *et al* 2016). The rapid visual information processing capacity is heavily dependent on visual parameters and use of an HMD headset would enable standardization of viewing distance, room lighting and visual angle (Foerster *et al* 2016). Use of a VR headset could distort electrode positions; nonetheless this affect could be easily mitigated. BCIs employing motion-onset visual evoked potentials (mVEP) have been utilized with VR headsets in neurogaming, and shown to be feasible (Beveridge *et al* 2016). The mVEP responses were evaluated in relation to mobile, complex and varying graphics within game distractors (Beveridge *et al* 2016). Foerster *et al* (2016) used the VR device oculus rift for neuropsychological assessment of visual processing capabilities. This VR device is head-mounted and covers the entire visual field, thereby shielding and standardizing the visual stimulation, and therefore may improve test-retest reliability. Compared to a CRT screen performances, visual processing speed, threshold of conscious perception and capacity of visual working memory did not differ significantly using the VR headset. VR headsets may therefore be applicable for standardized and reliable assessment and diagnosis of elementary cognitive functions in laboratory and clinical settings and maximize the opportunity to compare visual processing components between individuals and institutions and to establish statistical norm distributions. Recently, a new VR-EEG combined headset with electrodes embedded in occipital areas for ERP detection has been reported for neurogaming (www.neurable.com). RVSP-based BCI paradigms may therefore benefit from the head mounted visual displays however a vision obscuring headset may not be appropriate in some contexts as it could limit the ability of the users, e.g. a person with disabilities, to communicate with their peers and environment. Such a headset may prevent the expressive or receptive use of non-verbal communication skills, such as eye movement and facial expressions, that are vital for users with non-verbal communication skills.

Advancements towards RSVP of targets during moving sequences have shown promising results, although it is more difficult to study movie clips since the stimulus start event is not as clear. A remaining challenge in this area is for researchers to design signal processing tools that can deal with imprecise stimulus beginning/end (Cecotti 2015). However, an advantage of moving mode is that the target stimulus remains on the screen for longer than with static mode, allowing participants the opportunity to confirm a target stimulus. Moving stimuli studies to date have been limited to surveillance applications so there is a need for further investigation in this area.

Just over half the papers used the button press mode in conjunction with one of the other modes, as not all of the studies are concerned with comparing EEG responses to motor responses. It is important to develop a scale in order to rank the difficulty of tasks. This will enable the comparison of paradigms that are at the same level. The key outcomes of this study are shown in table 4, provided as suggested guidelines. These are suggested parameters that may be useful to researchers when designing RSVP-based BCI paradigms within the different application types. From this review, we can conclude that using these parameters will enable more consistent performance for the different application types and will enable improved comparison with new studies.

In acknowledgment of the need for standardization of parameters for RSVP-based BCI protocols, Cecotti *et al* (2011) raise an interesting proposal, stating that other parameters could be automatically prescribed in accordance with the chosen target likelihood, such as the optimal ISI length, classifiers and spatial filters (Cecotti *et al* 2011). Such an infrastructure for parameter choices does not currently exist with studies focusing on the impact of different parameters.

Future studies would benefit from engaging with iterative changes in design parameters. This would allow for a comparative study of the different design parameters and enable the identification of parameters that most affect the experimental paradigm. A study involving increasing the rate of presentation until classification starts to deteriorate significantly for various types of stimulus categories may indicate the maximum possible speed of RSVP-BCI. Additionally, a future development for RSVP-based BCIs might be to use real life imagery with numerous distractor stimuli amongst the target stimuli. This is a more difficult task but it would enhance paradigm relatability to real-life applications. Hybridizing RSVP BCIs with other BCI paradigms has also started to receive more attention (Kumar and Sahin 2013). Users of this system navigate using motor imagery movements (left, right, up and down). Search queries are spelt using the Hex-O-Speller and results retrieved from a web search engine may be fed back to the user using RSVP. This study shows the potential benefits of the RSVP paradigm and how it may be used in order to aid physically impaired users. Eye tracking can be used as an outcome measure to assess and enhance RSVP stimuli and presentation modes. Specifically, using eye-tracking researchers can establish where the participant's gaze is focused during erroneous trials and explore correlations between gaze variability and performance. With the RSVP-based BCI paradigm there is much scope to evaluate different data types/imagery. This is a fast-growing field with a promising future. There are multiple opportunities and a large array of potential RSVP-BCI paradigm setups. Researchers in the field are therefore recommended to consider the literature to date and the comparative framework proposed in this paper.

ORCID iDs

Stephanie Lees  <https://orcid.org/0000-0003-1036-5639>

References

- Acqualagna L and Blankertz B 2011 A gaze independent spelling based on rapid serial visual presentation *2011 Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society* (IEEE) p 45603
- Acqualagnav L *et al* 2010 A novel brain–computer interface based on the rapid serial visual presentation paradigm *2010 Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society* (IEEE) pp 2686–9
- Alpert G F, Manor R, Spanier A B, Deouel L Y and Geva A B 2014 Spatiotemporal representations of rapid visual target detection: a single-trial EEG classification algorithm *IEEE Trans. Biomed. Eng.* **61** 2290–303
- Barngrover C *et al* 2016 A brain–computer interface (BCI) for the detection of mine-like objects in sidescan sonar imagery *IEEE J. Ocean. Eng.* **41** 124–39
- Beveridge R, Wilson S and Coyle D 2016 3D graphics, virtual reality, and motion-onset visual evoked potentials in neurogaming *Brain–Computer Interfaces: Lab Experiments to Real-World Applications* vol 228, ed D Coyle (Amsterdam: Elsevier) pp 329–53
- Bigdely-Shamlo N *et al* 2008 Brain activity-based image classification from rapid serial visual presentation *IEEE Trans. Neural Syst. Rehabil. Eng.* **16** 432–41
- Blankertz B *et al* 2011 Single-trial analysis and classification of ERP components—a tutorial *NeuroImage* **56** 814–25
- Bohannon A W *et al* 2017 Collaborative image triage with humans and computer vision *2016 IEEE Int. Conf. on Systems, Man, and Cybernetics, SMC 2016—Conf. Proc.* pp 4046–51
- Boksem M A S, Meijman T F and Lorist M M 2005 Effects of mental fatigue on attention: an ERP study *Brain Res. Cogn. Brain Res.* **25** 107–16
- Cai B *et al* 2013 A rapid face recognition BCI system using single-trial ERP *2013 6th Int. IEEE/EMBS Conf. on Neural Engineering* pp 89–92
- Cecotti H 2015 Toward shift invariant detection of event-related potentials in non-invasive brain–computer interface *Pattern Recogn. Lett.* **66** 127–34
- Cecotti H 2016 Single-trial detection with magnetoencephalography during a dual-rapid serial visual presentation task *IEEE Trans. Biomed. Eng.* **63** 220–7
- Cecotti H *et al* 2011a Impact of target probability on single-trial EEG target detection in a difficult rapid serial visual presentation task *Conf. Proc., Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conf.* pp 6381–4
- Cecotti H *et al* 2011b Multimodal target detection using single trial evoked EEG responses in single and dual-tasks *Conf. Proc., Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conf.* pp 6311–4
- Cecotti H *et al* 2012a Multiclass classification of single-trial evoked EEG responses *2012 Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society* (IEEE) pp 1719–22
- Cecotti H, Eckstein M P and Giesbrecht B 2012b Effects of performing two visual tasks on single-trial detection of event-related potentials *Conf. Proc., Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conf.* pp 1723–6
- Cecotti H, Eckstein M P and Giesbrecht B 2014 Single-trial classification of event-related potentials in rapid serial visual presentation tasks using supervised spatial filtering *IEEE Trans. Neural Netw. Learn. Syst.* **25** 2030–42
- Chahine G and Krekelberg B 2009 Cortical contributions to saccadic suppression *PLoS One* **4** e6900

- Chennu S *et al* 2013 The cost of space independence in P300-BCI spellers *J. Neuroeng. Rehabil.* **10** 82
- Cohen M X 2014 *Analyzing Neural Time Series Data: Theory and Practice* (Cambridge, MA: MIT Press)
- Diamond M R, Ross J and Morrone M C 2000 Extraretinal control of saccadic suppression *J. Neurosci.* **20** 3449–55
- Dias J C and Parra L C 2011 No EEG evidence for subconscious detection during rapid serial visual presentation 2011 *IEEE Signal Process. Medicine and Biology Symp.* (IEEE) pp 1–4
- Erdogmus D, Mathan S and Pavel M 2006 Comparison of linear and nonlinear approaches on single trial ERP detection in rapid serial visual presentation tasks 2006 *IEEE Int. Joint Conf. on Neural Network Proc.* pp 1136–42
- Fawcett T 2006 An introduction to ROC analysis *Pattern Recogn. Lett.* **27** 861–74
- Files B T and Marathe A R 2016 A regression method for estimating performance in a rapid serial visual presentation target detection task *J. Neurosci. Methods* **258** 114–23
- Fine E M and Peli E 1995 Scrolled and rapid serial visual presentation texts are read at similar rates by the visually impaired *J. Opt. Soc. Am. A* **12** 2286–92
- Foerster R M *et al* 2016 Using the virtual reality device oculusrift for neuropsychological assessment of visual processing capabilities *Sci. Rep.* **6** 37016
- Forster K I 1970 Visual perception of rapidly presented word sequences of varying complexity *Perception Psychophys.* **8** 215–21
- Gerson A D, Parra L C and Sajda P 2005 Cortical origins of response time variability during rapid discrimination of visual objects *NeuroImage* **28** 342–53
- Gerson A D, Parra L C and Sajda P 2006 Cortically coupled computer vision for rapid image search *IEEE Trans. Neural Syst. Rehabil. Eng.* **14** 174–9
- Gonsalvez C L and Polich J 2002 P300 amplitude is determined by target-to-target interval *Psychophysiology* **39** 388–96
- Guo C, Voss J L and Paller K A 2005 Electrophysiological correlates of forming memories for faces, names, and face–name associations *Cogn. Brain Res.* **22** 153–64
- Hart B M *et al* 2013 Attention in natural scenes: contrast affects rapid visual processing and fixations alike *Phil. Trans. R. Soc. B* **368** 20130067
- Healy G and Smeaton A F 2011 Optimising the number of channels in EEG-augmented image search *Proc. of HCI 2011—25th BCS Conf. on Human Computer Interaction* (British Computer Society) pp 157–62
- Hild K E *et al* 2011 An ERP-based brain–computer interface for text entry using rapid serial visual presentation and language modeling *ACL HLT 2011—49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of Student Session* pp 38–43
- Hope C *et al* 2013 High throughput screening for mammography using a human–computer interface with rapid serial visual presentation (RSVP) *Proc. SPIE* **867303**
- Huang L *et al* 2017 BHCR: RSVP target retrieval BCI framework coupling with CNN by a Bayesian method *Neurocomputing* **238** 255–68
- Huang Y *et al* 2006 Comparison of linear and nonlinear approaches on single trial ERP detection in rapid serial visual presentation tasks 2006 *IEEE Int. Joint Conf. on Neural Network Proc.* (IEEE) pp 1136–42
- Huang Y *et al* 2007 A fusion approach for image triage using single trial erp detection *Proc. of the 3rd Int. IEEE EMBS Conf. on Neural Engineering* pp 473–6
- Huang Y *et al* 2008 Large-scale image database triage via EEG evoked responses *IEEE Int. Conf. on Acoustics, Speech and Signal Processing* pp 429–32
- Huang Y *et al* 2011 A framework for rapid visual image search using single-trial brain evoked responses *Neurocomputing* **74** 2041–51
- Johnson R 1986 A triarchic model of P300 amplitude *Psychophysiology* **23** 367–84
- Kindermans P-J *et al* 2014 Integrating dynamic stopping, transfer learning and language models in an adaptive zero-training ERP speller *J. Neural Eng.* **11** 035005
- Kranczioch C, Debener S and Engel A K 2003 Event-related potential correlates of the attentional blink phenomenon *Cogn. Brain Res.* **17** 177–87
- Kumar S and Sahin F 2013 Brain computer interface for interactive and intelligent image search and retrieval 2013 *High Capacity Optical Networks and Emerging/Enabling Technologies* (IEEE) pp 136–40
- Lamme V A F and Roelfsema P R 2000 The distinct modes of vision offered by feedforward and recurrent processing *Trends Neurosci.* **23** 571–9
- Leutgeb V, Schäfer A and Schienle A 2009 An event-related potential study on exposure therapy for patients suffering from spider phobia *Biol. Psychol.* **82** 293–300
- Lin Z, Ying Z, Hui G, Li T, Chi Z, Xiaojuan W, Qunjian W and Bin Y 2017 Multi-rapid serial visual presentation framework for EEG-based target detection *BioMed. Res. Int.* **2017** 2049094
- Liu R and Kwon M 2016 Integrating oculomotor and perceptual training to induce a pseudofovea: a model system for studying central vision loss *J. Vis.* **16** 10
- Lotte F *et al* 2007 A review of classification algorithms for EEG-based brain–computer interfaces *J. Neural Eng.* **4** R1–3
- Luck S J 2005 *An Introduction to the Event-Related Potential Technique* (Cambridge, MA: MIT Press)
- Luck S, Woodman G and Vogel E 2000 Event-related potential studies of attention *Trends Cogn. Sci.* **4** 432–40
- Luo W *et al* 2010 Three stages of facial expression processing: ERP study with rapid serial visual presentation *NeuroImage* **49** 1857–67
- MacKenzie G and Donaldson D I 2007 Dissociating recollection from familiarity: Electrophysiological evidence that familiarity for faces is associated with a posterior old/new effect *NeuroImage* **36** 454–63
- Maguire J F and Howe P D L 2016 Failure to detect meaning in RSVP at 27 ms per picture *Atten. Percept. Psychophys.* **78** 1405–13
- Manor R and Geva A B 2015 Convolutional neural network for multi-category rapid serial visual presentation BCI *Front. Comput. Neurosci.* **9** 146
- Manor R, Mishali L and Geva A B 2016 Multimodal neural network for rapid serial visual presentation brain computer interface *Front. Comput. Neurosci.* **10** 130
- Marathe A *et al* 2015a Improved neural signal classification in a rapid serial visual presentation task using active learning *IEEE Trans. Neural Syst. Rehabil. Eng.* **4320** 1–11
- Marathe A R *et al* 2014a Confidence metrics improve human–autonomy integration *Proc. of the 2014 ACM/IEEE Int. Conf. on Human-Robot Interaction* pp 240–1
- Marathe A R *et al* 2015b The effect of target and non-target similarity on neural classification performance: a boost from confidence *Front. Neurosci.* **9** 1–11
- Marathe A R, Ries A J and McDowell K 2014b Sliding HDCA: single-trial eeg classification to overcome and quantify temporal variability *IEEE Trans. Neural Syst. Rehabil. Eng.* **22** 201–11
- Mathan S *et al* 2008 Rapid image analysis using neural signals *Proc. of the 26th Annual CHI Conf. Extended Abstracts on Human Factors in Computing Systems* (New York: ACM Press) p 3309

- Matran-Fernandez A and Poli R 2014 Collaborative brain–computer interfaces for target localisation in rapid serial visual presentation *6th Computer Science and Electronic Engineering Conf. CEEC 2014—Conf. Proc.* pp 127–32
- Maurer U, Rossion B and McCandliss B D 2008 Category specificity in early perception: face and word n170 responses differ in both lateralization and habituation properties *Front. Hum. Neurosci.* **2** 18
- McCarthy G and Donchin E 1981 A metric for thought: a comparison of P300 latency and reaction time *Science* **211** 77–80
- Ming D *et al* 2010 Time-locked and phase-locked features of P300 event-related potentials (ERPs) for brain–computer interface speller *Biomed. Signal Process. Control* **5** 243–51
- Mohedano E *et al* 2014 Object segmentation in images using EEG signals *Proc. of the ACM Int. Conf. on Multimedia* (New York: ACM Press) pp 417–26
- Mohedano E *et al* 2015 Exploring EEG for object detection and retrieval *Proc. 5th ACM Int. Conf. on Multimedia Retrieval (ICMR 2015) (Shanghai, China, 23–26 June 2015)* pp 591–4
- Oken B S *et al* 2014 Brain–computer interface with language model–electroencephalography fusion for locked-in syndrome *Neurorehabil. Neural Repair* **28** 387–94
- Oliva A 2005 Gist of the scene *Encyclopedia of Neurobiology of Attention* (San Diego, CA: Elsevier)
- Orhan U *et al* 2011 Fusion with language models improves spelling accuracy for ERP-based brain computer interface spellers *Conf. Proc. Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conf.* pp 5774–7
- Orhan U *et al* 2012a Improved accuracy using recursive bayesian estimation based language model fusion in ERP-based BCI typing systems *Conf. Proc., Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conf.* pp 2497–500
- Orhan U *et al* 2012b RSVP keyboard: an EEG based typing interface *ICASSP, IEEE Int. Conf. on Acoustics, Speech and Signal Processing—Proc.* pp 645–8
- Orhan U *et al* 2013 Offline analysis of context contribution to ERP-based typing BCI performance *J. Neural Eng.* **10** 066003
- Paller K A and Wagner A D 2002 Observing the transformation of experience into memory *Trends Cogn. Sci.* **6** 93–102
- Parra L *et al* 2008 Spatiotemporal linear decoding of brain state *IEEE Signal Process. Mag.* **25** 107–15
- Parra L C *et al* 2005 Recipes for the linear analysis of EEG *NeuroImage* **28** 326–41
- Polich J and Donchin E 1988 ‘P300 and the word frequency effect’ *Electroencephalogr. Clin. Neurophysiol.* **70** 33–45
- Poolman P, Frank R M, Luu P, Pederson S M and Tucker D M 2008 A single-trial analytic framework for eeg analysis and its application to target detection and classification *NeuroImage* **42** 787–98
- Potter M C *et al* 2002 Recognition memory for briefly presented pictures: the time course of rapid forgetting *J. Exp. Psychol. Hum. Percept. Perform.* **28** 1163–75
- Potter M C *et al* 2014 Detecting meaning in RSVP at 13 ms per picture *Atten. Percept. Psychophys.* **76** 270–9
- Primativo S *et al* 2016 Perceptual and cognitive factors imposing ‘speed limits’ on reading rate: a study with the rapid serial visual presentation *PLOS One* **11** e0153786
- Raymond J E, Shapiro K L and Arnell K M 1992 Temporary suppression of visual processing in an RSVP task: an attentional blink? *J. Exp. Psychol. Hum. Percept. Perform.* **18** 459–60
- Rensink R A 2000 When good observers go bad: change blindness, inattention blindness, and visual experience *Change* **6** 458–68
- Rivet B and Souloumiac A 2013 Optimal linear spatial filters for event-related potentials based on a spatio-temporal model: asymptotical performance analysis *Signal Process.* **93** 387–98
- Rivet B *et al* 2009 xDAWN algorithm to enhance evoked potentials: application to brain–computer interface *IEEE Trans. Biomed. Eng.* **56** 2035–43
- Rosenthal D *et al* 2014 Evoked neural responses to events in video *IEEE J. Sel. Top. Signal Process.* **8** 358–65
- Rousset G A, Thorpe S J and Fabre-Thorpe M 2004 How parallel is visual processing in the ventral pathway? *Trends Cogn. Sci.* **8** 363–70
- Rugg M D *et al* 1998 Dissociation of the neural correlates of implicit and explicit memory *Nature* **392** 595–8
- Saavedra C and Bougrain L 2012 Processing stages of visual stimuli and event-related potentials *NeuroComp/KEOps’12 Workshop*
- Sajda P *et al* 2010 ‘In a blink of an eye and a switch of a transistor: cortically coupled computer vision’ *Proc. IEEE* **98** 462–78
- Sajda P *et al* 2014 Evoked neural responses to events in video *IEEE J. Sel. Top. Signal Process.* **8** 358–65
- Sajda P, Gerson A and Parra L 2003 High-throughput image search via single-trial event detection in a rapid serial visual presentation task *Proc. 1st Int. IEEE EMBS Conf. on Neural Engineering (Capri Island, Italy, 20–22 March 2003)* pp 7–10
- Sasane S and Schwabe L 2012 Decoding of EEG activity from object views : active detection versus passive visual tasks *Brain Informatics. BI 2012. Lecture Notes in Computer Science*, vol 7670, ed F M Zanzotto, S Tsumoto, N Taatgen, Y Yao (Berlin: Springer) pp 277–87
- Serre T, Oliva A and Poggio T 2007 A feedforward architecture accounts for rapid categorization *Proc. Natl Acad. Sci. USA* **104** 6424–9
- Simons D J and Levin D T 1997 Change blindness *Trends Cogn. Sci.* **1** 261–7
- Smith M E 1993 Neurophysiological manifestations of recollective experience during recognition memory judgments *J. Cogn. Neurosci.* **5** 1–13
- Stoica A *et al* 2013 Multi-brain fusion and applications to intelligence analysis *Proc. SPIE* **87560N**
- Touryan J *et al* 2011 Real-time measurement of face recognition in rapid serial visual presentation *Front. Psychol.* **2** 1–8
- Ušćumlić M, Chavarriaga R and Millán J D R 2013 An iterative framework for EEG-based image search: robust retrieval with weak classifiers *PloS One* **8** e72018
- Vidal J J 1973 Toward direct brain–computer communication *Ann. Rev. Biophys. Bioeng.* **157**–80
- Wang J *et al* 2009 Brain state decoding for rapid image retrieval *Proc. of the 17th ACM Int. Conf. on Multimedia* (New York: ACM Press) p 945
- Wang Y and Jung T-P 2011 A collaborative brain–computer interface for improving human performance *PLoS One* **6** e20422
- Ward R, Duncan J and Shapiro K 1997 Effects of similarity, difficulty, and nontarget presentation on the time course of visual attention *Percept. Psychophys.* **59** 593–600
- Waytowich N R *et al* 2016 Spectral transfer learning using information geometry for a user-independent brain–computer interface *Front. Neurosci.* **10** 430
- Wei C-S, Lin Y-P, Wang Y-T, Lin C-T and Jung T-P 2016 Transfer learning with large-scale data in brain–computer interfaces *38th Annual Int. Conf. of the IEEE Engineering in Medicine and Biology Society (EMBC) (Orlando, FL, USA, 16–20 August 2016)* pp 4666–9
- Weiden M, Khosla D and Keegan M 2012 Electroencephalographic detection of visual saliency of motion towards a practical brain–computer interface for video analysis *Proc. of the 14th ACM Int. Conf. on Multimodal Interaction* (New York: ACM Press) p 601

- Wolpaw J R and Wolpaw E W 2012 *Brain–Computer Interfaces: Principles and Practice* (Oxford: Oxford University Press)
- Won D-O *et al* 2017 Shifting stimuli for brain computer interface based on rapid serial visual presentation 2017 5th Int. Winter Conf. on Brain–Computer Interface (BCI) (IEEE) pp 40–1
- Yazdani A *et al* 2010 The impact of expertise on brain computer interface based salient image retrieval 2010 Annual Int. Conf. of the IEEE Engineering in Medicine and Biology (IEEE) pp 1646–9
- Yovel G and Paller K A 2004 The neural basis of the butcher-on-the-bus phenomenon: when a face seems familiar but is not remembered *NeuroImage* **21** 789–800
- Yu K *et al* 2011 A bilinear feature extraction method for rapid serial visual presentation triage (<https://doi.org/10.1109/IWBE.2011.6079025>)
- Yu K *et al* 2014 The analytic bilinear discrimination of single-trial EEG signals in rapid image triage *PloS One* **9** e100097
- Zander T O *et al* 2010 Enhancing human-computer interaction with input from active and passive brain–computer interfaces *Brain–Computer Interfaces* ed D S Tan and A Nijholt (London: Springer) pp 181–99
- Zander T O and Kothe C 2011 Towards passive brain–computer interfaces: applying brain–computer interface technology to human–machine systems in general *J. Neural Eng.* **8** 025005
- Zhang Y *et al* 2012 A novel BCI based on ERP components sensitive to configural processing of human faces *J. Neural Eng.* **9** 026018