UNEDITED VERSION—PLEASE DO NOT CITE WITHOUT PERMISSION

The Role of Conjunctive Representations in Stopping Actions

Atsushi Kikumoto

Department of Cognitive, Linguistics, and Psychological sciences, Brown University and Center for Brain Sciences, RIKEN

Ulrich Mayr

Department of Psychology University of Oregon

Corresponding Author: Ulrich Mayr, University of Oregon, Department of Psychology, 1227 University of Oregon, 97403 Eugene, OR,

Abbreviated Title: Conjunctive Representations and Action Regulation

Keywords: conjunctive representations, response inhibition, EEG

Acknowledgements: This research was supported by NIA grant R01 AG037564-01A1, and by NSF grant NSF grant 1734264.

Abstract

Action selection appears to rely on conjunctive representations that nonlinearly integrate task-relevant features (Kikumoto & Mayr, 2020). We test here the flip-side of this hypothesis that such representations are also intricately involved during attempts to stop an action—a key aspect of action regulation. We tracked both conjunctive representations and those of constituent rule, stimulus, or response features through trial-by-trial representational similarity analysis of the EEG signal in a combined, rule-selection and stop-signal paradigm. Across two experiments with student participants (N = 57), we found (a) that the strength of decoded conjunctive representations prior to the stop-signal uniquely predicted trial-by-trial stopping success (Exp. 1) and (b) that these representations were selectively suppressed following the onset of the stop-signal (Exp. 1 and 2). We conclude that conjunctive representations are key to successful action execution and therefore need to be suppressed when an intended action is no longer appropriate.

Statement of Relevance

A key aspect of self-control is the ability to stop on ongoing action. While considerable progress has been made in understanding the stopping process itself, much less is known about which exact representations are targeted by this process. Based on the idea that actions rely on conjunctive representations that bind action-relevant features together, we hypothesized that it is exactly these conjunctive representations the stopping process is up against. Through EEG-decoding techniques, we tracked conjunctive and basic feature representations and showed that indeed the stronger the conjunctive representations, the harder it is to stop an intended action. Conversely, when the stopping process was successful, the strength of conjunctive representations was selectively reduced. These results further our knowledge about action regulation by showing that conjunctive representations are a necessary precursor for carrying out actions successfully and for that reason also need to be the target of self-regulatory stopping attempts.

The Role of Conjunctive Representations in Stopping Actions

Research on cognitive control has made considerable progress in determining the set of processes involved in stopping an intended or ongoing action, a key aspect of self-regulation (Aron, Robbins, & Poldrack, 2014; Logan & Cowan, 1984; Swann et al., 2009; Verbruggen et al., 2019; Wessel, 2020). Yet, much less is known about the exact representations that are targeted by the presumed stopping processes.

Even simple goal-directed actions, such as kicking a soccer ball to a teammate, rely on various features—the location of the ball, the presence of opponent players and teammates, as well as on abstract rules (e.g., "kick softly when the grass is wet"). In principle, for stopping of an action it may be sufficient to simply suppress the low-level motor representations that directly control the kicking movement (Coxon, Stinear, & Byblow, 2006; Duque, Greenhouse, Labruna, & Ivry, 2017; Greenhouse, Sias, Labruna, & Ivry, 2015; Labruna et al., 2014), while leaving other task-relevant representations intact.

An alternative hypothesis can be derived from theories and recent findings about the representations that serve as the main drivers of successful action. For example, event-file theory (Hommel, 2019; Hommel, Müsseler, Aschersleben, & Prinz, 2001; Schumacher & Hazeltine, 2016) posits that an action becomes executable only once all task-relevant features are integrated within a *conjunctive representation*, also referred to as event file. Moreover, recent research in non-human primates indicates that populations of neurons with nonlinear response properties integrate various task aspects in a conjunctive manner and are critical for successful actions (Parthasarathy et al., 2017; Rigotti et al., 2013; Stokes et al., 2013).

If conjunctive representations are necessary, and maybe even sufficient precursors of goal-directed behavior, it follows that the path towards stopping a given action may also need to run through these representations. Thus, when in the above soccer scenario, an opponent defender suddenly blocks the goal, not just the codes that directly control the kicking

movement that need to be canceled, but the entire integrated representation of all involved action-relevant features.

So far, it has been difficult to characterize the representational targets of the stopping process because this would require tracking the fate of various relevant representations concurrently, as actions are selected and regulated. However, recently Kikumoto and Mayr (2020) demonstrated that time-resolved representational similarity analysis (RSA, Kriegeskorte, Mur, & Bandettini, 2008) of the EEG signal can be used to monitor both conjunctive and constituent feature representations during rule-based action selection in humans, and on the level of individual trials. These analyses indeed revealed conjunctive representations that integrated action rules and specific sensory/motor settings throughout the entire selection period. Moreover, the strength of conjunctions was a robust and unique predictor of trial-to-trial variability in RTs—as one would expect if conjunctive representations are necessary and sufficient conditions for action execution.

For a direct test of the hypothesis that the pathway to canceling an action leads through the corresponding conjunctive representation, we combined here a rule-based action selection task (Fig. 1ab, Kikumoto & Mayr, 2020; Mayr & Bryck, 2005) with an occasional stop signal. In Exp. 1, the stop-signal timing was adjusted via an adaptive tracking procedure. This allowed us to test the prediction that the strength of conjunctions prior to stopping, inversely predicts stopping success and to provide initial information about which representations are targeted by the stopping process. In Exp. 2, the stop-signal was presented 100 ms after stimulus onset, which is early enough for successfully stopping actions in most trials. Here, our main goal was to clearly characterize the consequences of successful stopping of actions on task-relevant representations and specifically test the prediction that conjunctive representations are selectively suppressed following the stop signal.

Materials and Methods

Participants

A convenience sample of 64 students of the University of Oregon participated after signing an informed consent following a protocol approved by the University of Oregon's Human Subjects Committee in exchange for the compensation of \$10 per hour and additional performance-based incentives. Our target samples size was based on our previous work (Kikumoto & Mayr, 2020), where we had obtained highly robust results with sample sizes around 20 participants. For Exp. 1, the target sample size was increased to 36 because the design included an additional comparison between failed and successful stop trials. Exp. 2 followed more closely the design of Kikumoto & Mayr, 2020, but included stop-trials for 33% of trials, and therefore we chose a target sample size of 24 participants. In both experiments, we oversampled participants to account for removal of participants with excessive amount of EEG artifacts (i.e., more than 35% of trials; see *EEG recordings and preprocessing* for detail). As a result, we retained 36 out of 38 participants for Exp. 1, and 24 out of 26 participants for Exp. 2. In Exp. 1, three additional participants were excluded because their stopping accuracy was above 75% on stop trials.

Stimuli, Tasks and Procedure

Participants were randomly cued on a trial-by-trial basis to execute one of the three possible actions rules, resulting in a rule switch rate of p=.66. (Fig. 1a, Mayr & Bryck, 2005). Based on the cued rule, participants responded to the location of a circle (1.32° in radius) that randomly appeared in the corner of a white frame (6.6° in one side) by selecting one of four response keys that were arranged in a 2 x 2 matrix (4, 5, 1, and 2 on the number pad). For example, the vertical rule mapped the left-top dot to the bottom-left response as a correct response. Two different cue words (e.g., "vertical" or "updown") were used for each rule.

In 33.3% of trials, the stop-signal (i.e., a yellow frame; Fig. 1a) indicated to participants that the planned action had to be cancelled. Stop-trials were counted as successful when

participants did not make any responses within 800 ms time-window following the stop-signal onset. In Exp. 1, the interval between the stimulus and stop-signal onset (i.e., the stop-signal delay or SSD) was adjusted using an adaptive tracking method based on participants' trial-to-trial stopping success. Specifically, individuals' SSDs varied between 0 ms to 800 ms counting from the onset of the stimulus and starting with 100 ms at the beginning of session.

Correct/incorrect stop trials increased/decreased SSDs by the step size that was randomly selected from 11.8 ms, 23.5 ms, or 35.3 ms for each trial. In Exp. 2, the stop-signal appeared 100 ms after the stimulus onset. Go trials lasted until the response was executed; stop trials lasted either until the 800 ms response window expired, or until a response (i.e., failed stopping) was recorded.

There were two practice blocks and 250 and 200 experimental blocks for Exp. 1 and 2 respectively. Each block lasted 15 seconds, within which participants were instructed to complete as many trials as possible. Trials that were initiated within the 15 second block duration were allowed to complete. The average number of go-trials and stop-trials were 1576 (SD=162) and 773 (SD=75) for Exp. 1, and 1378 (SD=91) and 685 (SD=33) for Exp. 2. Throughout the experimental session, participants were reminded to respond as accurately and fast as possible and refrain from waiting for the stop signal. In Exp. 1, participants were instructed that the adaptive tracking procedure would make it easier to stop on some trials and more challenging on others. Participants were given a performance-based incentive for trials with RTs on go-trials faster than the 75th percentile of correct responses in the preceding blocks when 1) the overall accuracy in go-trials was above 90 percent and 2) there were more than 5 completed trials in a given block. While performing the task, participants were asked to rest the index finger in the center of the four response keys at the start of each trial (i.e., no lateralization of response sides). At the end of each trial, feedback (a green fixation cross for correct and a red cross for incorrect trials) was presented based on the accuracy of responses in go-trials or on correct stopping in stop-trials. At the end of each block, the number of

completed trials, the number of correct responses in go/stop-trials, and the amount of earned incentives based on the speed of responses in go-trials, were presented as a feedback. All stimuli were created in Matlab (Mathworks) using the Psychophysics Toolbox (Brainard, 1997) and were presented on a 17-inch CRT monitor (refresh rate: 60 *Hz*) at a viewing distance of 100 cm.

Stop-signal Reaction Time

In Exp. 1, we computed individuals' stop-signal reaction time (SSRT), according to the integration method as specified by Verbruggen et al. (2019). First, for each quantile bin of SSDs (Fig. 2b), the mean SSDs and the proportion of successful stop trials (p(respond|signal)) were calculated. Then, the matching go RTs were defined in each SSD bin by taking the nth RT in the rank ordered go-trial RTs (including all go-trials), where n is defined by multiplying the number of RTs in the distribution by the probability of responding, p(respond|signal) or unsuccessful stopping, for each SSD bin. Within each SSD bin, SSRT was calculated by subtracting the corresponding SSD from the matching go RT, then scores from 6 SSD bins were averaged within individuals to obtain a single metric of SSRT for each individual. For all participants, failed-stop RTs were faster than correct go RTs.

EEG recordings and preprocessing

Scalp EEG activities were recorded from 20 tin electrodes on an elastic cap (Electro-Caps) using the International 10/20 system. The 10/20 sites F3, Fz, F4, T3, C3, CZ, C4, T4, P3, PZ, P4, T5, T6, O1, and O2 were used along with five nonstandard sites: OL halfway between T5 and O1; OR halfway between T6 and O2; PO3 halfway between P3 and OL; PO4 halfway between P4 and OR; and POz halfway between PO3 and PO4. Electrodes placed ~1cm to the left and right of the external canthi of each eye recorded horizontal electrooculogram (EOG) to measure horizontal saccades. To detect blinks, vertical EOG was recorded from an electrode placed beneath the left eye and reference to the left mastoid. The left-mastoid was used as reference for all recording sites, and data were re-referenced off-line

to reflect the signal at the right-mastoid. The scalp EEG and EOG were amplified with an SA Instrumentation amplifier with a bandpass of 0.01–80 Hz, and signals were digitized at 250 Hz in LabView 6.1 running on a PC. EEG data was first segmented by 18.5 second intervals to include all trials within a block. After time-frequency decomposition was performed, these epochs were further segmented into trial-to-trial epochs (the time interval of -600 to 800 ms relative to the onset of the stimulus for both experiments). These trial-to-trial epochs including blinks (>80 μ v, window size = 200 ms, window step = 50 ms), large eye movements (>1°, window size = 200 ms, window step = 10 ms), blocking of signals (range = -0.01 μ v to 0.01 μ v) were excluded from subsequent analyses.

Time-Frequency Analysis

We used the time-frequency decomposed EEG signal for decoding of representations (e.g., Foster, Sutterer, Serences, Vogel, & Awh, 2017; Kikumoto & Mayr, 2018). Temporal-spectral profiles of single-trial EEG data were obtained via complex wavelet analysis (Cohen, 2014) by applying time-frequency analysis to preprocessed EEG data epoched for each block (>18 seconds to exclude the edge artifacts). The power spectrum was convolved with a series of complex Morlet wavelets ($e^{i2\pi ft}e^{-t^2/(2\sigma^2)}$), where t is time, f is frequency increased from 1 to 35 Hz in 35 logarithmically spaced steps, and σ defines the width of each frequency band, set according to $n/2\pi f$, where n increased from 3 to 10. The logarithmic scaling was used to keep the width across frequency band approximately equal, and the incremental number of wavelet cycles was used to balance temporal and frequency precision as a function of frequency of the wavelet. After convolution was performed in the frequency-domain, we took an inverse of the Fourier transform, resulting in complex signals in the time-domain. A frequency band-specific estimate at each time point was defined as the squared magnitude of the convolved signal $Z(real[z(t)]^2 + imag[z(t)]^2)$ for instantaneous power.

Representational Similarity Analysis

The decoding analysis in the current study follows closely our previously established methods (Kikumoto & Mayr, 2020). In order to assess the strength of each action feature and conjunction on the level of individual trials and time points, we used a two-step procedure. An initial linear decoding step yielded similarity information that served in the second step as input to a representational similarity analysis. For the initial step, we used a penalized linear discriminant analysis using the caret package in R (Kuhn, 2008) to discriminate between all 12 possible action constellations. At every time sample point, the instantaneous power of rhythmic EEG activity was averaged within the predefined ranges of frequency values (1-3 Hz for the delta-band, 4-7 Hz for the theta-band, 8-12 Hz for the alpha-band, 13-30 Hz for the beta-band, 31-35 Hz for the gamma-band), generating 100 features (5 frequency-bands X 20 electrodes) to train decoders. Within individuals and within each frequency-band, these data points were ztransformed across electrodes at every time sample to remove the effects that uniformly influenced all electrodes. We used a k-fold repeated, cross-validation procedure to evaluate the decoding results (Mosteller & Tukey, 1968), by randomly partitioning single-trial EEG data into four independent folds. All trials except incorrect go-trials were used as the training sets in both experiments. The number of observations of each action constellation was kept equal within and across folds by dropping excess trials randomly. Three folds served as a training set and the remaining fold was used as a test set; this step was repeated until each fold had served as a test set. Each cross-validation cycle was repeated eight times, in which each step generated a new set of randomized folds. Resulting classification probabilities (i.e., evidence estimated for each case of S-R mapping) were averaged across all cross-validated results with the best-tuned hyperparameter to regularize the coefficients. This decoding step yielded for each time point and trial a "confusion-vector" of classification probabilities for both the correct and all possible incorrect classifications (Fig. 1c).

As the second step, we applied time-resolved RSAs to each confusion profile in order to determine the underlying similarity structure. Specifically, we regressed the confusion vector onto model vectors as predictors, which were derived from a set of representational similarity model matrices (Fig. 1c). Each model matrix uniquely represents a potential, underlying representation (e.g., rules, stimuli, responses and conjunctions). For example, the rule model predicts that the decoder would only discriminate instances of different rules, but fail to discriminate instances of the same rule. To estimate the unique variance explained by competing models, we regressed all model vectors simultaneously, resulting in coefficients for each of the four model vectors. These coefficients (i.e., their corresponding *t*-values) allowed us to relate the dynamics of action representations to trial-to-trial variability in behavior during go- and stop-trials (see Multilevel Modeling section for details). For all RSAs, we logit-transformed classification probabilities and further included subject-specific vectors that contained *z*-scored, average RTs and stopping accuracy as nuisance predictors to reduce potential biases in decoding due to idiosyncratic differences in performance among action constellations (see also Fig.S2 and S3 in the Supplemental Material).

We excluded *t*-values that exceeded 5 SDs from means for each sample point, which excluded less than 1% of the entire samples in both experiments. Resulting *t*-values were averaged within 20 ms non-overlapping time samples.

In both experiments, decoders were trained with the stimulus-aligned EEG signal. In Exp. 1, we further computed RSA scores that were re-epoched in reference to the onset of the stop-signal (the right column of Fig. 3 and 4). Matching go-trial results were calculated with the SSDs that would have been used if the stop-signal appeared in those trials.

Estimating Timing of Stop-induced Suppression

In Exp. 2, we used nonparametric permutation tests with a single-threshold method to identify the earliest time sample at which statistically significant differences between go-trials and stop-trials emerged. Specifically, for each action feature, we computed permutation

distributions of the maximum statistic for every sample point from the stop-signal onset (fixed at 100 ms after the stimulus onset) to the end of 800 ms of the hold period. First, we obtained RSA results by decoding data with randomly shuffled condition labels (i.e., of action constellations). We then performed a series of *t*-tests, testing the differences in RSA scores between go- and stop-trials, for every sample against the null level (i.e., 0 for *t*-values). Out of the series of *t*-test results, we retained the maximum *t*-value. We repeated this process 10000 times by randomly drawing samples from all possible permutations of labels, thereby generating the permutation distributions of the maximum statistics. This approach allowed us to identify statistically significant, individual time points by comparing scores from the correct labels to the critical threshold, which was defined as the 95th (i.e., alpha = .05) of the largest member of maximum statistics in the permutation distribution of the corresponding variable.

Multilevel Modeling

In Exp. 1, to analyze predictors of trial-by-trial variability in stopping success, we used multilevel logistic regression models. Specifically, we estimated for stop trials a model predicting stopping success on a given trial using the RSA-derived t-values for basic action features (i.e., rule, stimulus, and response) and the conjunction as predictors. In addition, we also included each trial's log-transformed SSD as a covariate to account for the possibility that SSDs affect both action representations and stopping success as a third-variable. For statistical tests, we used EEG data averaged over a-priori selected, symmetric time intervals, namely a pre-stop-signal (-200 to 0 ms) and a post-stop-signal period (0 ms to 200 ms), relative to the onset of the stop-signal in each trial. Both time intervals clearly precede the average SSRT across individuals (M = 272 ms). We also performed additional control analyses, where we excluded trials with early responses (i.e., responses occurred after the stimulus onset and before the stop-signal in unsuccessful-stop trials) and where we included decoded representations from both pre- and post-stop-signal phases simultaneously (Table 3). In addition, to visualize changes in predictability of stopping success, we separately performed

a series of logistic regression analyses by fitting models at each sample point in reference to the onset of the stimulus and the stop-signal (Fig. 5). In order to replicate the results by (Kikumoto & Mayr, 2020) about how action representations contributed to action selection in go-trials, we also report for both experiments results from multilevel models to assess which action representations predict trial-to-trial RTs on go trials (Table 2). Here, RTs were log-transformed and trials with response errors were excluded.

Open Practice Statement

Neither of the experiments reported in this article was formally preregistered. All data and analysis scripts will be posted on OSF for the final manuscript.

Results

Experiment 1

Behavior

Behavioral performance is summarized in Table 1. Most participants (33 out of 36 participants) exhibited p(stop|signal) in the range of .40-.65; individuals with the stopping accuracy higher than 75% were excluded from further analyses. The average RTs in go-trials were longer than the RT in failed stop-trials for all participants. This pattern is consistent with the race model as a basis for estimating individuals' stop-signal reaction time (SSRT, Fig. 2a). Also, the probability of stopping errors covaried with the increase of SSDs, indicating the overall efficacy of the SSD staircase algorithm (Fig. 2b).

Action Representations in Go-trials, Failed Stop-trials, and Successful Stop-trials

Fig. 3 shows the time-course of RSA scores estimated on the level of single trials for each of the basic features (i.e., rules, stimuli, and responses) and the conjunction all trial types. For go-trials, the flow of activated representations was highly consistent with our previous results. Rule information appeared in the pre-stimulus period, stimulus information peaked shortly after the stimulus appeared, followed by the emergence of response information (Hubbard, Kikumoto, & Mayr, 2019; Kikumoto & Mayr, 2020). Importantly,

conjunctive information was present throughout the entire response-selection period. We also replicated the previous finding that trial-to-trial variability in conjunctive representations robustly predicted go-trial RTs (Table 2), over and above other representations of constituent features. In the Supplemental Material, we include additional analyses that probe the degree to which key RSA results are specific to particular frequency bands (Fig.S4, modulated by posterior or global alpha/theta power (Tables S1 and S2), or can be explained through the overfitting of statistical noise (Fig.S5).

The strength of conjunctive representations and late response representations were selectively reduced in successful stop-trials relative to go-trials and failed stop-trials. In contrast, there were no clear effects on rule and stimulus representations. For the conjunctive representations, the divergence between failed and successful stop trials occurred even before the onset of the stop-signal (see individuals' average SSDs in Fig. 3 left column), suggesting stopping was particularly impaired when the conjunctive representations were strong in the early response selection phase. Indeed, when we replotted RSA scores relative to trial-to-trial SSDs (see Fig. 3, right column), differences in successful and failed stop-trials emerged clearly before the average SSRT (M=272 ms, SE=9.23 ms) and even before the stop signal. No other action features showed similar differences in the pre-stop-signal period (t < .13), and post-stop-signal effects on the response representation were apparent only when the conjunction model was excluded, b=-.010, SE=.004, t=2.34.

Before accepting the conclusion that the state of conjunctive representations prior to the onset of the stop-signal determined the success of stopping, we need to consider the fact that trial-to-trial variability in SSDs might covary with both the strength of conjunctions and the probability of successful stopping. Therefore, to rule out SSDs as a potential third-variable explanation, we performed multilevel logistic regressions to predict single-trial stopping failures using decoded action features and SSDs as simultaneous predictors. Fig. 4, shows time-point by time-point results of these analyses, which clearly indicate that pre-stop-signal conjunctions

are a unique predictor of stopping success. Statistical tests of these relationships confirmed that the average state of conjunctions prior to the onset of the stop-signal strongly predicted stopping failures over and above the state of the other feature representations (Table 3, top panel). To ensure that these results are not due to very fast responses that occurred prior to the stop signal, we confirmed that they were robust when eliminating premature responses (Table 3, middle panel). In addition, when entering pre-stop-signal (-200 ms – 0 ms) and post-stop-signal predictors (0 ms – 200 ms) simultaneously, we found that during each phase, conjunctions uniquely predicted stopping success (Table 3, bottom panel).

Arguably, a pre-stop-signal effect of the conjunctive representations may be driven by participants who occasionally used a strategy of waiting for the stop-signal in order to initiate the go process. We used a median split to separate subjects by the proportion of stopping failures. As shown in the Supplemental Material, both pre- and post-stop-signal states of conjunctions predicted stopping failures within both subgroups, even in the above-median group, which should be less likely to use a waiting strategy (Table S3). This indicates that the observed results cannot be attributed to a subset of particularly cautious participants.

The stop-signal-aligned pattern shown in Fig. 3 (right column) provides initial evidence that conjunctions are the representational target of the stopping process. Specifically, in failed stop trials, conjunctive representations were particularly strong at the time the stop-signal arrived, but showed a rapid reduction immediately following the stop-signal (Fig. 3 right column). Such a pattern is consistent with a stop-signal-induced suppression that came too late to influence behavior. We used multi-level regression to test the linear trend in the strength of conjunction during the 320 ms period after the stop-signal (using consecutive, 40ms time-windows). This linear trend was indeed stronger for failed-stop compared to successful-stop or go trials combined, b=.008, SE=.0038, t=2.04. However, this pattern remains somewhat inconclusive as the reduction following the peak at the time of the stop-signal might also be

interpreted as a regression towards the mean level of conjunction strength. In Exp. 2, we will seek more definitive evidence regarding the representational targets of the stopping process.

Experiment 2

Exp. 1 demonstrated that the strength of conjunctive representations predicts of trial-to-stopping success and provided initial evidence that conjunctive representations are predominantly affected by the stopping process. However, while the adaptive calibration of the stopping success rate around 50%, was useful for determining the relationship between conjunctive representations and behavioral performance, it made it more difficult to distinguish between the representational precursors and consequences of stopping success. In order to clearly establish the effects of stopping on different representations, we used a consistent, early stop-signal in Exp. 2, which ensured high, overall stopping success.

Behavior

Behavioral performance in go/stop-trials are summarized in Table 1. The probability of stopping failures—incorrectly executing responses in the presence of the stop-signal (i.e., p(respond|signal))—was low because of the early presentation of stop-signal at the fixed timing (100 ms after the stimulus onset). This allowed us to estimate the time-course of suppression of action representations from a fixed starting point. Note, that the substantially lower RTs for go trials in Exp. 2 compared to Exp. 1, are likely due to the reduced uncertainty about the timing of the stop-signal in Exp. 2.

Action Representations in Go-trials and Stop-trials

As shown in Fig. 5, the pattern of activated representations was highly consistent with our previous results. Also again, conjunctive representations robustly predicted go-trial RTs (Table 2), over and above representations of constituent features.

Our main goal in Exp. 2 was to test the prediction that conjunctive representations are suppressed on stop-trials relative to go-trials. Indeed, we found stopping of actions markedly reduced the strength of conjunctive representations right after the onset of the stop-signal (Fig.

5). Not surprisingly, the response representation was also suppressed, whereas we found no effect on the rule representation and only a late, and small effect for the stimulus representation. Yet, when only the basic constituent features (i.e., rules, stimuli and responses) were used in the RSA model, the suppression effect was substantially increased for the rule representation, highlighting the importance of including the conjunction model (see also, Kikumoto & Mayr, 2020). Suppression of the conjunction occurred at the same time, or even slightly before suppression of the response representation. This suggests that the reduction of conjunctive information is not just an aftereffect of response suppression. Rather, it supports the notion that the conjunctive representation is a direct target of the stopping activity.

Discussion

Even simple, goal-directed actions rely on various aspects of the task environment. Both psychological and neural-level theories assume that successful action requires the integration of all relevant action features within a conjunctive representation (Hommel et al., 2001; Rigotti et al., 2013). Using a convenience sample of student participants, we tested the novel hypothesis that because such representations are critical for action selection, they are also intricately involved when a planned or initiated action needs to be stopped. Consistent with this hypothesis, we found that the strength of conjunctive representations at the time the stopping process is initiated, inversely predicted stopping success (Fig. 3 and 4) and that conjunctive representations were a main target of the stopping process (Fig. 3 and 5).

In principle, stopping of actions could require suppression of the entire set of task-relevant representations. Alternatively, only low-level, "motor" representations might be targeted. Instead, our results are most consistent with the hypothesis that the conjunctive representation is the primary target of suppression, followed by the response representation (Fig. 3, 4, and 5). It is an open question whether conjunctive and response representations are separately targeted, or whether the deactivation of response representations is a downstream

consequence of the suppressed conjunctive representations. Representations of the rule or the stimulus remained intact, or showed very minor suppression, and only after completion of the stopping process (Fig. 3 and 5).

Conjunctive representations that integrate stimulus, response, and rule information are by definition situated on more central level than representations that directly control motor output. The fact that conjunctive representations were targeted by the stopping process, is consistent with results indicating that inhibition of simple motor responses and inhibition of more "central", thoughts or memories are handled by a shared process (Anderson, 2004; Guo, Schmitz, Mur, Ferreira, & Anderson, 2018). For example, the same right lateral prefrontal area that is typically involved in stopping of motor responses, was also critical in suppressing thoughts, leading to longer-term negative effects on their accessibility.

The just-mentioned observation of aftereffects of suppression raises interesting questions about the functional benefits of conjunctive representations. By adding contextual specificity to abstract feature codes, conjunctive representations should help constrain inhibition to the currently relevant context and thereby mitigate suppression aftereffects on memory representations (Anderson & Green, 2001). For example, in the action context studied here, a potential benefit of selectively suppressing conjunctions is that, once the reason for stopping has been removed, actions can be easily reimplemented on the basis of the intact rule and stimulus representations. Thus, conjunctive representations as suppression targets may increase flexibility by keeping representations of constituent features available for further use.

Our finding that the strength of conjunctive representations inversely predicts subsequent stopping success (Table 3), is also important for event-file theory (Hommel, 2004; Hommel et al., 2001). A central claim of this theory is that conjunctive representations are a key causal factor behind successful actions. However, this claim has never been explicitly

tested. Short of a difficult to achieve, experimental manipulation of conjunctive representation strength, our results come as close as possible to establishing such a test.

One might argue that the finding that conjunction strength predicts stopping failure is little more than a direct consequence of the observation that conjunction strength predicts response speed (see Table 2; Fig. 3 and 4 in Kikumoto & Mayr, 2020). After all, by the logic of the race model of stopping, faster go responses are per definition more difficult to stop (Logan & Cowan, 1984). However, our results show that basic feature representations also have some predictive power, making them potential targets for the stopping process (see Table 3). Therefore, the finding that conjunctive representations emerged both as a strong early predictor of stopping success (Fig. 4), and also as the earliest suppression target (Fig. 5), is a novel result. Even more importantly, the finding that strong conjunctions predict faster responses is open to an alternative explanation, namely in terms of an unspecific relationship between decoding quality and efficiency of goal-directed behavior (e.g., due to fluctuations in neural noise as a third variable). However, as the present results show, strong conjunctions make actions both faster (Table 2) and more difficult to stop. Thus, the predictive power of conjunctive representations is highly specific to the efficiency of the corresponding action, not to goal-directed behavior in general—which in the current case would have included successful stopping.

The RSA analyses presented here allow us to say with certainty that at least two different task-relevant features were integrated within conjunctions, but the task spaces we used did not allow disambiguating between conjunctions that include binary combinations of rules, stimuli, or responses, or the combination between all three aspects. However, previous work had also used an expanded task space that allowed disambiguating between different types of conjunctions (Exp.2 in Kikumoto & Mayr, 2020). In terms of functional characteristics, the rule-specific conjunctions from the extended task space, had behaved in a highly similar manner to the conjunctions derived from the more ambiguous situations. In addition, we also

show here that when the conjunction predictor was dropped from the RSA analyses, the coefficients for the rule predictor absorbed much of the conjunction effect—indicating a contribution of rule information to the decoded conjunctions (Fig.5 inlets). Therefore, the observed conjunctions likely reflect an integration between both the rule, and stimulus/response features.

Our results by themselves do not identify the underlying mechanisms that modulate the state of the conjunctive representation prior to the stop signal. One possibility is that conjunction strength depends on endogenous fluctuations of attention towards the go-action across trials. An emphasis on initiating action may induce strong conjunctions and thereby cause the failures to trigger the stop process altogether (Matzke, Hughes, Badcock, Michie, & Heathcote, 2017; Matzke, Love, & Heathcote, 2017). The fact that conjunctive representations in failed-stop trials, prior to the arrival of stop signal, were even stronger than on go-trials, is consistent with such an attentional fluctuation account. As another, not necessarily exclusive possibility, there is evidence that variations in strategic, proactive inhibition (Aron, 2011) affect the state of conjunctions. On trials in which subjects anticipate stopping, proactive inhibition may keep conjunctive representations from fully developing. Such proactive control processes could be directly tested by cuing the stop probability on a trial-by-trial basis (Chikazoe et al., 2009; Vink, Kaldewaij, Zandbelt, Pas, & du Plessis, 2015; Zandbelt, Bloemendaal, Neggers, Kahn, & Vink, 2013). In any case, our results clearly confirm that the pre-stop-signal state of action representations must be taken into account to fully understand subsequent reactive inhibition and stopping.

Our EEG-based decoding results also provide no precise information about the neural-anatomical location of conjunctive representations (see Supplementary Information). However, recently there has been increasing evidence from research with non-human primates about the high prevalence of neurons in parietal/frontal areas that show very similar properties as the conjunctive representations we report on (Fusi, Miller, & Rigotti, 2016; Parthasarathy et al.,

2017; Rigotti et al., 2013). Specifically, these so-called mixed-selectivity neurons integrate various task-features in a nonlinear and diverse manner and, just as the EEG-decoded conjunctive representations, are uniquely predictive of successful action selection. It would be important to establish the degree to which the conjunctive representations examined here reflect activity of mixed-selectivity neurons. One way to test this hypothesis is to look for equivalent functional and computational properties of both conjunctive and nonlinearly mixed-selectivity representations in both human and animal models (Badre, Bhandari, Keglovits, & Kikumoto, 2021).

In conclusion, we provide for the first time detailed results about which precise representations are suppressed when trying to stop an intended action. Specifically, we show that conjunctive representations are the primary target of action inhibition, exactly because they are a key driver of successful action.

References

- Anderson, M. C. (2004). Neural Systems Underlying the Suppression of Unwanted Memories. *Science*, 303(5655), 232-235. doi:10.1126/science.1089504
- Anderson, M. C., & Green, C. (2001). Suppressing unwanted memories by executive control. *Nature*, *410*(6826), 366-369. doi:10.1038/35066572
- Aron, A. R. (2011). From reactive to proactive and selective control: developing a richer model for stopping inappropriate responses. *Biological Psychiatry*, 69(12), e55-68. doi:10.1016/j.biopsych.2010.07.024
- Aron, A. R., Robbins, T. W., & Poldrack, R. A. (2014). Inhibition and the right inferior frontal cortex; one decade on. *Trends in Cognitive Sciences*. *18*(4), 177-185.
- Badre, D., Bhandari, A., Keglovits, H., & Kikumoto, A. (2021). The dimensionality of neural representations for control. *Current Opinion in Behavioral Sciences*, 38, 20-28.
- Brainard, D. H. (1997). The psychophysics toolbox. Spatial Vision, 10(4), 433-436.
- Chikazoe, J., Jimura, K., Hirose, S., Yamashita, K.-i., Miyashita, Y., & Konishi, S. (2009). Preparation to inhibit a response complements response inhibition during performance of a stop-signal task. *Journal of Neuroscience*, *29*(50), 15870-15877.
- Cohen, M. X. (2014). Analyzing Neural Time Series Data: Theory and Practice: MIT Press.
- Coxon, J. P., Stinear, C. M., & Byblow, W. D. (2006). Intracortical inhibition during volitional inhibition of prepared action. *Journal of Neurophysiology*, *95*(6), 3371-3383.
- Duque, J., Greenhouse, I., Labruna, L., & Ivry, R. B. (2017). Physiological Markers of Motor Inhibition during Human Behavior. *Trends in Neurosciences*, *40*(4), 219-236. doi:10.1016/j.tins.2017.02.006
- Foster, J. J., Sutterer, D. W., Serences, J. T., Vogel, E. K., & Awh, E. (2017). Alpha-band oscillations enable spatially and temporally resolved tracking of covert spatial attention. *Psychological Science*, *28*(7), 929-941.
- Fusi, S., Miller, E. K., & Rigotti, M. (2016). Why neurons mix: high dimensionality for higher cognition. *Current Opinion in Neurobiology*, *37*, 66-74.
- Greenhouse, I., Sias, A., Labruna, L., & Ivry, R. B. (2015). Nonspecific Inhibition of the Motor System during Response Preparation. *Journal of Neuroscience*, *35*(30), 10675-10684.
- Guo, Y., Schmitz, T. W., Mur, M., Ferreira, C. S., & Anderson, M. C. (2018). A supramodal role of the basal ganglia in memory and motor inhibition: Meta-analytic evidence. *Neuropsychologia*, *108*, 117-134. doi:10.1016/j.neuropsychologia.2017.11.033
- Hommel, B. (2004). Event files: Feature binding in and across perception and action. *Trends in Cognitive Sciences*, 8(11), 494-500.
- Hommel, B. (2019). Theory of Event Coding (TEC) V2. 0: Representing and controlling perception and action. *Attention, Perception, & Psychophysics*, 1-16.
- Hommel, B., Müsseler, J., Aschersleben, G., & Prinz, W. (2001). The theory of event coding (TEC): A framework for perception and action planning. *Behavioral and Brain Sciences*, 24(5), 849-878.
- Hubbard, J., Kikumoto, A., & Mayr, U. (2019). EEG decoding reveals the strength and temporal dynamics of goal-relevant representations. *Scientific reports*, *9*(1), 1-11.
- Kikumoto, A., & Mayr, U. (2018). Decoding hierarchical control of sequential behavior in oscillatory EEG activity. *eLife*, 7, e38550.
- Kikumoto, A., & Mayr, U. (2020). Conjunctive representations that integrate stimuli, responses, and rules are critical for action selection. *Proceedings of the National Academy of Sciences*, 201922166. doi:10.1073/pnas.1922166117
- Kriegeskorte, N., Mur, M., & Bandettini, P. (2008). Representational similarity analysis connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2, 4. doi:10.3389/neuro.06.004.2008
- Kuhn, M. (2008). Caret package. Journal of Statistical Software, 28(5).

- Labruna, L., Lebon, F., Duque, J., Klein, P.-A., Cazares, C., & Ivry, R. B. (2014). Generic inhibition of the selected movement and constrained inhibition of nonselected movements during response preparation. *Journal of Cognitive Neuroscience*, *26*(2), 269-278.
- Logan, G. D., & Cowan, W. B. (1984). On the ability to inhibit thought and action: A theory of an act of control. *Psychological Review*, *91*(3), 295.
- Matzke, D., Hughes, M., Badcock, J. C., Michie, P., & Heathcote, A. (2017). Failures of cognitive control or attention? The case of stop-signal deficits in schizophrenia. *Attention, Perception, & Psychophysics, 79*(4), 1078-1086.
- Matzke, D., Love, J., & Heathcote, A. (2017). A Bayesian approach for estimating the probability of trigger failures in the stop-signal paradigm. *Behavior Research Methods*, 49(1), 267-281. doi:10.3758/s13428-015-0695-8
- Mayr, U., & Bryck, R. L. (2005). Sticky rules: integration between abstract rules and specific actions. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 31*(2), 337-350. doi:10.1037/0278-7393.31.2.337
- Mosteller, F., & Tukey, J. W. (1968). Handbook of Social Psychology. 2, 80-203.
- Parthasarathy, A., Herikstad, R., Bong, J. H., Medina, F. S., Libedinsky, C., & Yen, S.-C. (2017). Mixed selectivity morphs population codes in prefrontal cortex. *Nature Neuroscience*, *20*(12), 1770-1779.
- Rigotti, M., Barak, O., Warden, M. R., Wang, X.-J., Daw, N. D., Miller, E. K., & Fusi, S. (2013). The importance of mixed selectivity in complex cognitive tasks. *Nature*, *497*(7451), 585.
- Schumacher, E. H., & Hazeltine, E. (2016). Hierarchical task representation: Task files and response selection. *Current Directions in Psychological Science*, *25*(6), 449-454.
- Snijders, T. A., & Bosker, R. J. (2011). *Multilevel analysis:* An introduction to basic and advanced multilevel modeling: Sage.
- Stokes, M. G., Kusunoki, M., Sigala, N., Nili, H., Gaffan, D., & Duncan, J. (2013). Dynamic coding for cognitive control in prefrontal cortex. *Neuron*, *78*(2), 364-375. doi:10.1016/j.neuron.2013.01.039
- Swann, N., Tandon, N., Canolty, R., Ellmore, T. M., McEvoy, L. K., Dreyer, S., . . . Aron, A. R. (2009). Intracranial EEG reveals a time-and frequency-specific role for the right inferior frontal gyrus and primary motor cortex in stopping initiated responses. *Journal of Neuroscience*, *29*(40), 12675-12685.
- Verbruggen, F., Aron, A. R., Band, G. P., Beste, C., Bissett, P. G., Brockett, A. T., . . . Boehler, C. N. (2019). A consensus guide to capturing the ability to inhibit actions and impulsive behaviors in the stop-signal task. *eLife*. *8*. doi:10.7554/eLife.46323
- Vink, M., Kaldewaij, R., Zandbelt, B. B., Pas, P., & du Plessis, S. (2015). The role of stop-signal probability and expectation in proactive inhibition. *European Journal of Neuroscience*, *41*(8), 1086-1094.
- Wessel, J. R. (2020). β-bursts reveal the trial-to-trial dynamics of movement initiation and cancellation. *Journal of Neuroscience*, *40*(2), 411-423.
- Zandbelt, B. B., Bloemendaal, M., Neggers, S. F. W., Kahn, R. S., & Vink, M. (2013). Expectations and violations: delineating the neural network of proactive inhibitory control. *Human Brain Mapping*, *34*(9), 2015-2024. doi:10.1002/hbm.22047

Acknowledgements

This research was supported by NIA grant R01 AG037564-01A1, and by NSF grant NSF grant 1734264.

Figures

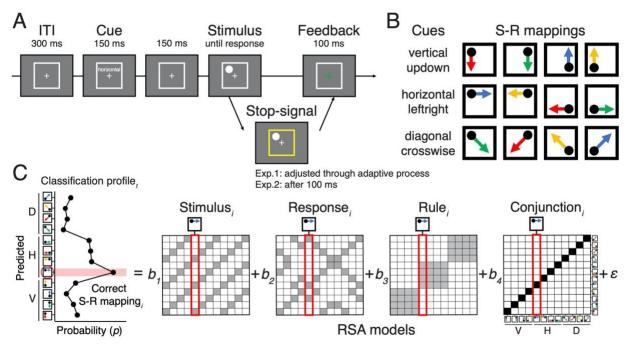


Fig. 1. Task design and analytic procedure.

(A) Sequence of trial events in the combined rule-selection/stop-signal task for both Exp. 1 and 2. (B) Spatial translation rules mapping specific stimuli to responses. Two different cue words were used for each rule to disambiguate between cue and rule-level representations. (C) Schematic steps of the representational similarity analysis. The raw EEG signal was decomposed into activity in specific frequency-bands via time-frequency analysis (see EEG recordings and preprocessing and Time-Frequency Analysis). For each sample time (t), a scalp-distributed pattern of EEG power was used to decode the specific rule/stimulus/response configuration of a given trial, producing a set of classification probabilities for each of the possible configurations. The profile of classification probabilities reflects the similarity structure of the underlying representations, where similar action constellations are more likely to be confused. The figure shows an example classification probabilities for a case where both a unique conjunction and rule information are expressed (peak at the correct S-R mapping, plus confusion to other instances with the same rule). For each trial and timepoint, the profile of classification probabilities was regressed onto model vectors as predictors that reflect the different, possible representations. In each matrix of model vectors, the x-axis corresponds to the correct constellation for the decoder to pick, and the y-axis shows all possible constellation. The shading of squares indicates the theoretically predicted classification probabilities (darker shading means higher probabilities). The coefficients associated with each predictor (i.e., tvalues) reflect the unique variance explained by each of the constituent features and their conjunction.

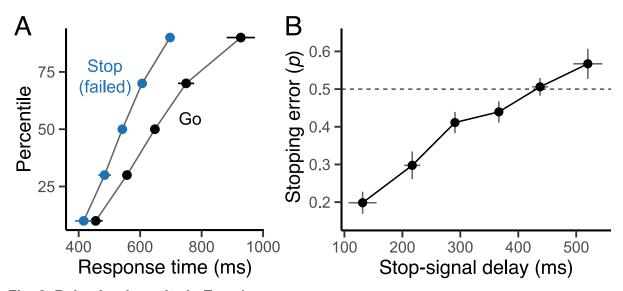


Fig. 2. **Behavioral results in Exp. 1**. (A) Vincentized mean response times (RTs) for go-trials and failed stop-trials. (B) Average rates of stopping failures as a function of stop-signal delays. Error bars specify 95% within-subject confidence intervals.

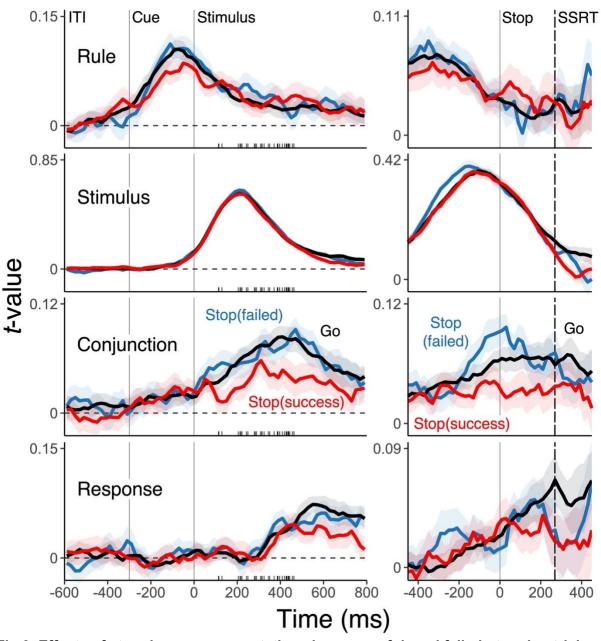


Fig 3. Effects of stopping on representations in successful- and failed-stopping trials compared to go-trials.

Average, single-trial *t*-values associated with each of the basic features (rule, stimulus, and response) and their conjunction derived from the RSA, separately for go-trials (black), successful stop-trials (red), and failed stop-trials (blue). The left panels show the results aligned with the stimulus onset, the right panels aligned with the stop-signal onset. Shaded regions specify the 95% within-subject confidence intervals. Tick marks on the x-axis of the stimulus-aligned panels mark individuals' average stop-signal delays.

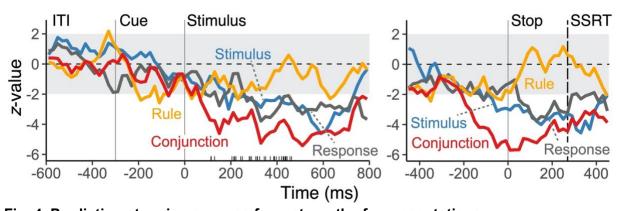


Fig. 4. Predicting stopping success from strength of representations.Time-course of *z* values from multilevel, logistic regression models predicting the variability in

Time-course of *z* values from multilevel, logistic regression models predicting the variability in trial-to-trial stopping failures in the stop-trials (the "impact" of representations on stopping success), using RSA scores of all features and trial-to-trial SSDs as simultaneous predictors. Negative *z*-value indicates more stopping failures as the strength of decoded representations increase. The left panel shows results aligned to the stimulus onset; in the right panel data are aligned to the stop-signal onset. Tick marks on the x-axis of the stimulus-aligned panels mark individuals' average stop-signal delays.

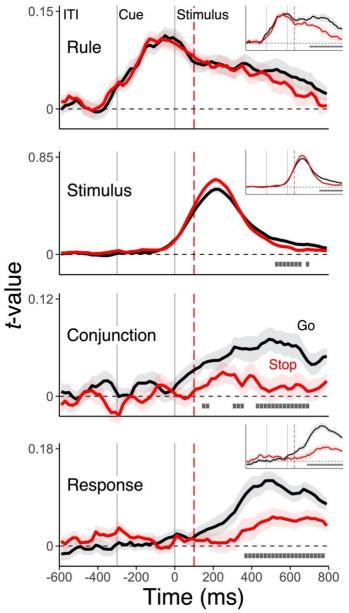


Fig 5. Effects of stopping on representations in the fixed-stop delay. Average, single-trial *t*-values derived from the RSA (see Fig. 1C) for each of the basic features (rule, stimulus, and response) and the conjunction, separately for go-trials (black) and stoptrials (red). Shaded regions specify the 95% within-subject confidence intervals. The vertical, red dashed line marks the onset of the stop-signal at 100 ms after the stimulus onset. Gray squares below lines denote the time points with significant differences between go- and stoptrials, correcting for multiple comparison using a non-parametric permutation test. The inserts for the rule, stimulus, and response features show the same results when the RSA contains only these basic features, but excludes the conjunction as model predictor.

Tables

Table 1. Behavioral performance in go- and stop-trials.

	Exp.1	Exp.2
Go RT (ms)	675 (24.3)	503 (12.8)
Go error (%)	3.55 (.63)	3.24 (.36)
<pre>p(respond signal)</pre>	44.1 (.84)	11.6 (1.81)
Failed stop RT (ms)	552 (15.8)	405 (6.5)
Stop error (%)	2.83 (.40)	2.36 (.54)
SSD (ms)	327 (17.9)	
SSRT (ms)	272 (9.2)	

Note. Stop errors=incorrect, failed stop responses, SSD=stop-signal delay, SSRT=stop-signal RT computed by using the entire RT distributions within each participant.

Table 2. Predicting trial-by-trial RTs in go-trials using the average strength of representations decoded through the RSA analyses during 0-300 ms post-stimulus intervals for each trial.

	Exp.1		Exp.2		
Variable	b (se)	t-value	b (se)	<i>t</i> -value	
Rule	013 (.006)	-2.08	025 (.012)	-2.15	
Stimulus	038 (.009)	-4.12	015 (.009)	-1.59	
Response	032 (.009)	-3.55	016 (.009)	-1.85	
Conjunction	057 (.010)	-5.69	042 (.012)	-3.61	

Note. Coefficients for all decoded variables were included as predictors simultaneously. Negative coefficient imply that stronger representations predict faster RTs. Standardized coefficients are reported, which can serve as effect-size estimates for fixed effects in multilevel models (Snijders & Bosker, 2011).

Table 3. Predicting trial-by-trial stopping accuracy using the strength of decoded representations in Exp. 1.

	•	Pre-Stop-Signal		Post-Stop-Signal	
Model	Variable	b (se)	<i>t</i> -value	b (se)	<i>t</i> -value
SSD control	Rule	077 (.033)	-2.38	012 (.024)	-0.49
	Stimulus	083 (.033)	-2.50	060 (.020)	-3.11
	Response	081 (.034)	-2.41	062 (.020)	-3.11
	Conjunction	193 (.041)	-4.70	151 (.029)	-5.23
Exclude early					
responses	Rule	050 (.020)	-2.43	019 (.025)	-0.76
	Stimulus	036 (.017)	-2.19	057 (.020)	-2.89
	Response	039 (.021)	-1.86	054 (.020)	-2.70
	Conjunction	136 (.029)	-4.78	155 (.029)	-5.31
Pre/post stop signal					
simultaneous	Rule	084 (.036)	-2.31	.016 (.044)	0.36
	Stimulus	035 (.036)	-0.93	104 (.041)	-2.53
	Response	038 (.040)	-0.97	093 (.041)	-2.30
	Conjunction	126 (.049)	-2.60	194 (.049)	-3.93

Note. The "SSD control" model included trial-to-trial stop-signal delays as the fixed and random effect as covariate. In the "exclude early responses" model, all premature responses that occurred prior to the onset of the stop-signal were removed. Whereas these models were fitted separately for pre-stop-signal and post-stop-signal predictors, in the "pre/post stop-signal simultaneous" model, both pre-stop-signal and post-stop-signal predictors were included simultaneously. Pre-stop-signal interval: -200–0 ms; post-stop-signal interval: 0-200 ms. Standardized coefficients are reported, which can serve as effect-size estimates for fixed effects in multi-level models (Snijders & Bosker, 2011).