



## UvA-DARE (Digital Academic Repository)

### A consensus guide to capturing the ability to inhibit actions and impulsive behaviors in the stop-signal task

Verbruggen, F.; Aron, A.R.; Band, G.P.H.; Beste, C.; Bissett, P.G.; Brockett, A.T.; Brown, J.W.; Chamberlain, S.R.; Chambers, C.D.; Colonius, H.; Colzato, L.S.; Corneil, B.D.; Coxon, J.P.; Dupuis, A.; Eagle, D.M.; Garavan, H.; Greenhouse, I.; Heathcote, A.; Huster, R.J.; Jahfari, S.; Kenemans, J.L.; Leunissen, I.; Li, C.-S.R.; Logan, G.D.; Matzke, D.; Morein-Zamir, S.; Murthy, A.; Paré, M.; Poldrack, R.A.; Ridderinkhof, K.R.; Robbins, T.W.; Roesch, M.; Rubia, K.; Schachar, R.J.; Schall, J.D.; Stock, A.-K.; Swann, N.C.; Thakkar, K.N.; van der Molen, M.W.; Vermeylen, L.; Vink, M.; Wessel, J.R.; Whelan, R.; Zandbelt, B.B.; Boehler, C.N.

**DOI**

[10.7554/eLife.46323](https://doi.org/10.7554/eLife.46323)

**Publication date**

2019

**Document Version**

Final published version

**Published in**

eLife

**License**

CC BY

[Link to publication](#)

**Citation for published version (APA):**

Verbruggen, F., Aron, A. R., Band, G. P. H., Beste, C., Bissett, P. G., Brockett, A. T., Brown, J. W., Chamberlain, S. R., Chambers, C. D., Colonius, H., Colzato, L. S., Corneil, B. D., Coxon, J. P., Dupuis, A., Eagle, D. M., Garavan, H., Greenhouse, I., Heathcote, A., Huster, R. J., ... Boehler, C. N. (2019). A consensus guide to capturing the ability to inhibit actions and impulsive behaviors in the stop-signal task. *eLife*, *8*, [e46323].  
<https://doi.org/10.7554/eLife.46323>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

# A consensus guide to capturing the ability to inhibit actions and impulsive behaviors in the stop-signal task

Frederick Verbruggen<sup>1\*</sup>, Adam R Aron<sup>2</sup>, Guido PH Band<sup>3</sup>, Christian Beste<sup>4</sup>, Patrick G Bissett<sup>5</sup>, Adam T Brockett<sup>6</sup>, Joshua W Brown<sup>7</sup>, Samuel R Chamberlain<sup>8</sup>, Christopher D Chambers<sup>9</sup>, Hans Colonius<sup>10</sup>, Lorenza S Colzato<sup>3</sup>, Brian D Corneil<sup>11</sup>, James P Coxon<sup>12</sup>, Annie Dupuis<sup>13</sup>, Dawn M Eagle<sup>8</sup>, Hugh Garavan<sup>14</sup>, Ian Greenhouse<sup>15</sup>, Andrew Heathcote<sup>16</sup>, René J Huster<sup>17</sup>, Sara Jahfari<sup>18</sup>, J Leon Kenemans<sup>19</sup>, Inge Leunissen<sup>20</sup>, Chiang-Shan R Li<sup>21</sup>, Gordon D Logan<sup>22</sup>, Dora Matzke<sup>23</sup>, Sharon Morein-Zamir<sup>24</sup>, Aditya Murthy<sup>25</sup>, Martin Paré<sup>26</sup>, Russell A Poldrack<sup>5</sup>, K Richard Ridderinkhof<sup>23</sup>, Trevor W Robbins<sup>8</sup>, Matthew Roesch<sup>6</sup>, Katya Rubia<sup>27</sup>, Russell J Schachar<sup>13</sup>, Jeffrey D Schall<sup>22</sup>, Ann-Kathrin Stock<sup>4</sup>, Nicole C Swann<sup>15</sup>, Katharine N Thakkar<sup>28</sup>, Maurits W van der Molen<sup>23</sup>, Luc Vermeulen<sup>1</sup>, Matthijs Vink<sup>19</sup>, Jan R Wessel<sup>29</sup>, Robert Whelan<sup>30</sup>, Bram B Zandbelt<sup>31</sup>, C Nico Boehler<sup>1</sup>

<sup>1</sup>Experimental Psychology, Ghent University, Ghent, Belgium; <sup>2</sup>University of California, San Diego, San Diego, United States; <sup>3</sup>Leiden University, Leiden, Netherlands; <sup>4</sup>Dresden University of Technology, Dresden, Germany; <sup>5</sup>Stanford University, Stanford, United States; <sup>6</sup>University of Maryland, College Park, United States; <sup>7</sup>Indiana University, Bloomington, United States; <sup>8</sup>University of Cambridge, Cambridge, United Kingdom; <sup>9</sup>Cardiff University, Cardiff, United Kingdom; <sup>10</sup>Oldenburg University, Oldenburg, Germany; <sup>11</sup>University of Western Ontario, London, Canada; <sup>12</sup>Monash University, Clayton, Australia; <sup>13</sup>University of Toronto, Toronto, Canada; <sup>14</sup>University of Vermont, Burlington, United States; <sup>15</sup>University of Oregon, Eugene, United States; <sup>16</sup>University of Tasmania, Hobart, Australia; <sup>17</sup>University of Oslo, Oslo, Norway; <sup>18</sup>Spinoza Centre Amsterdam, Amsterdam, Netherlands; <sup>19</sup>Utrecht University, Utrecht, Netherlands; <sup>20</sup>KU Leuven, Leuven, Belgium; <sup>21</sup>Yale University, New Haven, United States; <sup>22</sup>Vanderbilt University, Nashville, United States; <sup>23</sup>University of Amsterdam, Amsterdam, Netherlands; <sup>24</sup>Anglia Ruskin University, Cambridge, United Kingdom; <sup>25</sup>Indian Institute of Science, Bangalore, India; <sup>26</sup>Queen's University, Kingston, Canada; <sup>27</sup>King's College London, London, United Kingdom; <sup>28</sup>Michigan State University, East Lansing, United States; <sup>29</sup>University of Iowa, Iowa City, United States; <sup>30</sup>Trinity College Dublin, Dublin, Ireland; <sup>31</sup>Donders Institute, Nijmegen, Netherlands

\*For correspondence:  
frederick.verbruggen@ugent.be

Competing interest: See  
page 11

Funding: See page 11

Received: 22 February 2019

Accepted: 09 April 2019

Published: 29 April 2019

Reviewing editor: David Badre,  
Brown University, United States

© Copyright Verbruggen et al.  
This article is distributed under  
the terms of the [Creative  
Commons Attribution License](#),  
which permits unrestricted use  
and redistribution provided that  
the original author and source are  
credited.

**Abstract** Response inhibition is essential for navigating everyday life. Its derailment is considered integral to numerous neurological and psychiatric disorders, and more generally, to a wide range of behavioral and health problems. Response-inhibition efficiency furthermore correlates with treatment outcome in some of these conditions. The stop-signal task is an essential tool to determine how quickly response inhibition is implemented. Despite its apparent simplicity, there are many features (ranging from task design to data analysis) that vary across studies in ways that can easily compromise the validity of the obtained results. Our goal is to facilitate a more accurate use of the stop-signal task. To this end, we provide 12 easy-to-implement consensus

recommendations and point out the problems that can arise when they are not followed. Furthermore, we provide user-friendly open-source resources intended to inform statistical-power considerations, facilitate the correct implementation of the task, and assist in proper data analysis. DOI: <https://doi.org/10.7554/eLife.46323.001>

---

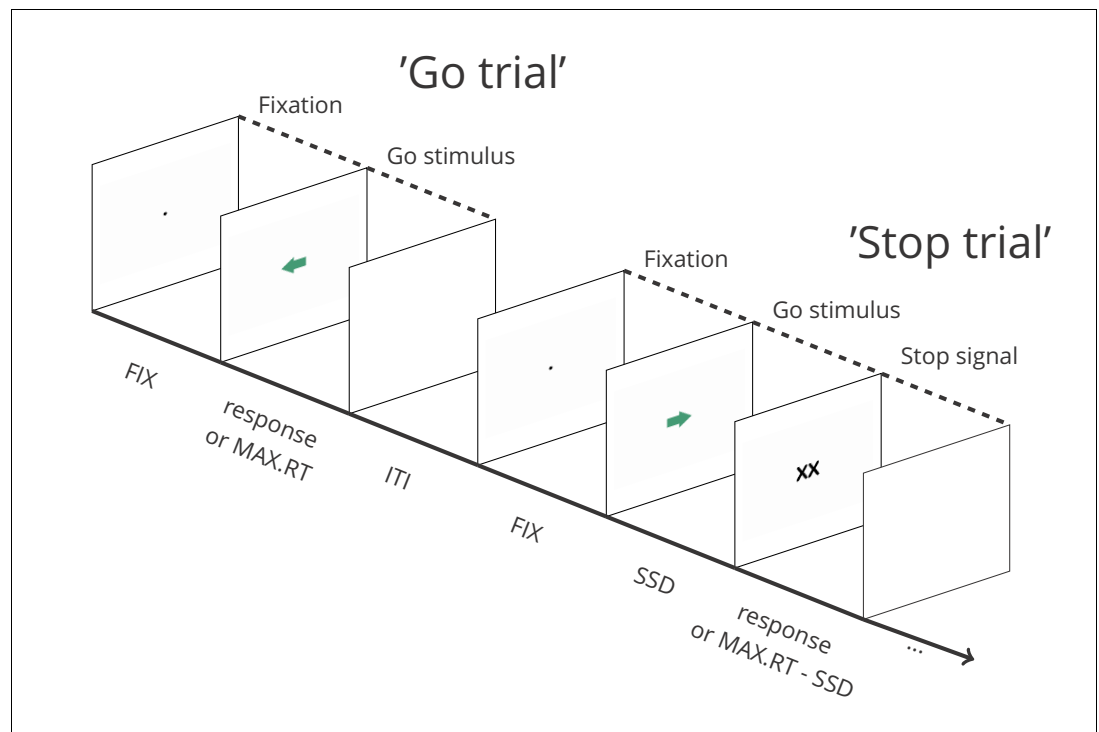
## Introduction

The ability to suppress unwanted or inappropriate actions and impulses ('response inhibition') is a crucial component of flexible and goal-directed behavior. The stop-signal task (*Lappin and Eriksen, 1966; Logan and Cowan, 1984; Vince, 1948*) is an essential tool for studying response inhibition in neuroscience, psychiatry, and psychology (among several other disciplines; see Appendix 1), and is used across various human (e.g. clinical vs. non-clinical, different age groups) and non-human (primates, rodents, etc.) populations. In this task, participants typically perform a go task (e.g. press left when an arrow pointing to the left appears, and right when an arrow pointing to the right appears), but on a minority of the trials, a stop signal (e.g. a cross replacing the arrow) appears after a variable stop-signal delay (SSD), instructing participants to suppress the imminent go response (*Figure 1*). Unlike the latency of go responses, response-inhibition latency cannot be observed directly (as successful response inhibition results in the absence of an observable response). The stop-signal task is unique in allowing the estimation of this covert latency (stop-signal reaction time or SSRT; *Box 1*). Research using the task has revealed links between inhibitory-control capacities and a wide range of behavioral and impulse-control problems in everyday life, including attention-deficit/hyperactivity disorder, substance abuse, eating disorders, and obsessive-compulsive behaviors (for meta-analyses, see e.g. *Bartholdy et al., 2016; Lipszyc and Schachar, 2010; Smith et al., 2014*).

Today, the stop-signal field is flourishing like never before (see Appendix 1). There is a risk, however, that the task falls victim to its own success, if it is used without sufficient regard for a number of important factors that jointly determine its validity. Currently, there is considerable heterogeneity in how stop-signal studies are designed and executed, how the SSRT is estimated, and how results of stop-signal studies are reported. This is highly problematic. First, what might seem like small design details can have an immense impact on the nature of the stop process and the task. The heterogeneity in designs also complicates between-study comparisons, and some combinations of design and analysis features are incompatible. Second, SSRT estimates are unreliable when inappropriate estimation methods are used or when the underlying race-model assumptions are (seriously) violated (see *Box 1* for a discussion of the race model). This can lead to artefactual and plainly incorrect results. Third, the validity of SSRT can be checked only if researchers report all relevant methodological information and data.

Here, we aim to address these issues by consensus. After an extensive consultation round, the authors of the present paper agreed on 12 recommendations that should safeguard and further improve the overall quality of future stop-signal research. The recommendations are based on previous methodological studies or, where further empirical support was required, on novel simulations (which are reported in Appendices 2–3). A full overview of the stop-signal literature is beyond the scope of this study (but see e.g. *Aron, 2011; Bari and Robbins, 2013; Chambers et al., 2009; Schall et al., 2017; Verbruggen and Logan, 2017*, for comprehensive overviews of the clinical, neuroscience, and cognitive stop-signal domains; see also the meta-analytic reviews mentioned above).

Below, we provide a concise description of the recommendations. We briefly introduce all important concepts in the main manuscript and the boxes. Appendix 4 provides an additional systematic overview of these concepts and their common alternative terms. Moreover, this article is accompanied by novel open-source resources that can be used to execute a stop-signal task and analyze the resulting data, in an easy-to-use way that complies with our present recommendations (<https://osf.io/rmqaw/>). The source code of the simulations (Appendices 2–3) is also provided, and can be used in the planning stage (e.g. to determine the required sample size under varying conditions, or acceptable levels of go omissions and RT distribution skew).



**Figure 1.** Depiction of the sequence of events in a stop-signal task (see <https://osf.io/rmqaw/> for open-source software to execute the task). In this example, participants respond to the direction of green arrows (by pressing the corresponding arrow key) in the go task. On one fourth of the trials, the arrow is replaced by 'XX' after a variable stop-signal delay (FIX = fixation duration; SSD = stop signal delay; MAX.RT = maximum reaction time; ITI = intertrial interval).

DOI: <https://doi.org/10.7554/eLife.46323.002>

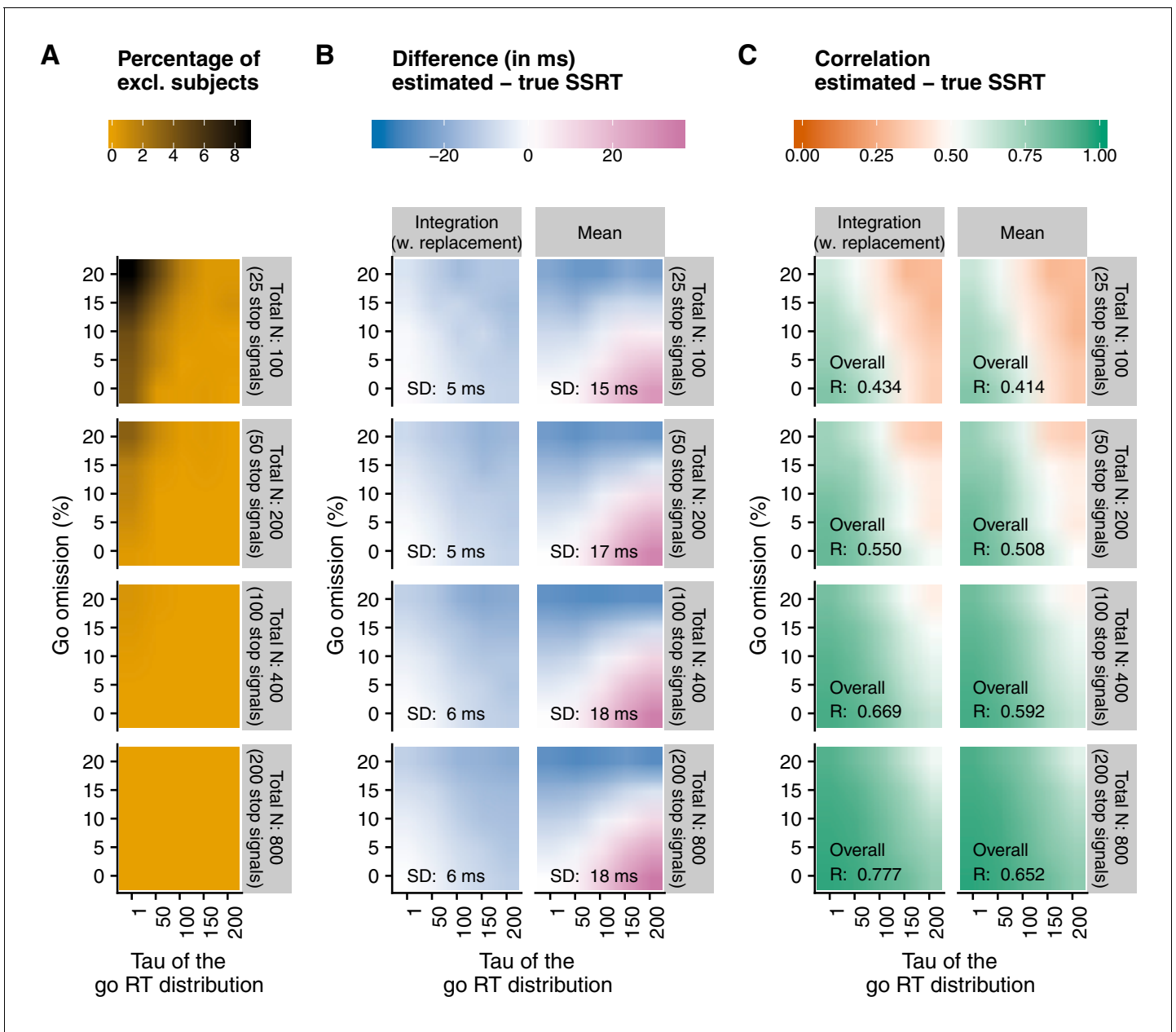
## Results and discussion

The following recommendations are for stop-signal users who are primarily interested in obtaining a reliable SSRT estimate under standard situations. The stop-signal task (or one of its variants) can also be used to study various aspects of executive control (e.g. performance monitoring, strategic adjustments, or learning) and their interactions, for which the design might have to be adjusted. However, researchers should be aware that this will come with specific challenges (e.g. *Bissett and Logan, 2014*; *Nelson et al., 2010*; *Verbruggen et al., 2013*; *Verbruggen and Logan, 2015*).

### How to design stop-signal experiments

#### Recommendation 1: Use an appropriate go task

Standard two-choice reaction time tasks (e.g. in which participants have to discriminate between left and right arrows) are recommended for most purposes and populations. When very simple go tasks are used, the go stimulus and the stop signal will closely overlap in time (because the SSD has to be very short to still allow for the possibility to inhibit a response), leading to violations of the race model as stop-signal presentation might interfere with encoding of the go stimulus. Substantially increasing the difficulty of the go task (e.g. by making the discrimination much harder) might also influence the stop process (e.g. the underlying latency distribution or the probability that the stop process is triggered). Thus, very simple and very difficult go tasks should be avoided unless the researcher has theoretical or methodological reasons for using them (for example, simple detection tasks have been used in animal studies. To avoid responses before the go stimulus is presented or close overlap between the presentation of go stimulus and stop signal, the intertrial interval can be drawn from a random exponential distribution. This will make the occurrence of the go stimulus unpredictable, discouraging anticipatory responses). While two-choice tasks are the most common, we note that the 'anticipatory response' variant of the stop-signal task (in which participants have to



**Figure 2.** Main results of the simulations reported in Appendix 2. Here, we show a comparison of the integration method (with replacement of go omissions) and the mean method, as a function of percentage of go omissions, skew of the RT distribution ( $\tau_{go}$ ), and number of trials. Appendix 2 provides a full overview of all methods. (A) The number of excluded ‘participants’ (RT on unsuccessful stop trials > RT on go trials). As this check was performed before SSRTs were estimated (see Recommendation 7), the number was the same for both estimation methods. (B) The average difference between the estimated and true SSRT (positive values = overestimation; negative values = underestimation). SD = standard deviation of the difference scores (per panel). (C) Correlation between the estimated and true SSRT (higher values = more reliable estimate). Overall R = correlation when collapsed across percentage of go omissions and  $\tau_{go}$ . Please note that the overall correlation does not necessarily correspond to the average of individual correlations.

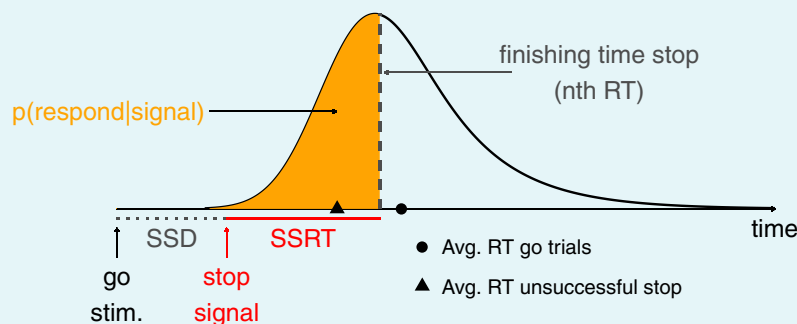
DOI: <https://doi.org/10.7554/eLife.46323.008>

press a key when a moving indicator reaches a stationary target) also holds promise (e.g. *Leunissen et al., 2017*).

## Box 1. The independent race model

Here, we provide a brief discussion of the independent race model, without the specifics of the underlying mathematical basis. However, we recommend that stop-signal users read the original modelling papers (e.g. [Logan and Cowan, 1984](#)) to fully understand the task and the main behavioral measures, and to learn more about variants of the race model (e.g. [Boucher et al., 2007](#); [Colonius and Diederich, 2018](#); [Logan et al., 2014](#); [Logan et al., 2015](#)).

Response inhibition in the stop-signal task can be conceptualized as an independent race between a 'go runner', triggered by the presentation of a go stimulus, and a 'stop runner', triggered by the presentation of a stop signal ([Logan and Cowan, 1984](#)). When the 'stop runner' finishes before the 'go runner', response inhibition is successful and no response is emitted (*successful stop trial*); but when the 'go runner' finishes before the 'stop runner', response inhibition is unsuccessful and the response is emitted (*unsuccessful stop trial*). The independent race model mathematically relates (a) the latencies (RT) of responses on unsuccessful stop trials; (b) RTs on go trials; and (c) the probability of responding on stop trials [ $p(\text{respond}|\text{signal})$ ] as a function of stop-signal delay (yielding 'inhibition functions'). Importantly, the independent race model provides methods for estimating the covert latency of the stop process (stop-signal reaction time; SSRT). These estimation methods are described in Materials and methods.



**Box 1—figure 1.** The independent race between go and stop.

DOI: <https://doi.org/10.7554/eLife.46323.004>

DOI: <https://doi.org/10.7554/eLife.46323.003>

### Recommendation 2: Use a salient stop signal

SSRT is the overall latency of a chain of processes involved in stopping a response, including the detection of the stop signal. Unless researchers are specifically interested in such perceptual or attentional processes, salient, easily detectable stop signals should be used (when auditory stop signals are used, these should not be too loud either, as very loud (i.e. >80 dB) auditory stimuli may produce a startle reflex). Salient stop signals will reduce the relative contribution of perceptual (afferent) processes to the SSRT, and the probability that within- or between-group differences can be attributed to them. Salient stop signals might also reduce the probability of a 'trigger failures' on stop trials (see [Box 2](#)).

### Recommendation 3: Present stop signals on a minority of trials

When participants strategically wait for a stop signal to occur, the nature of the stop-signal process and task change (complicating the comparison between conditions or groups; e.g. SSRT group differences might be caused by differential slowing or strategic adjustments). Importantly, SSRT estimates will also become less reliable when participants wait for the stop signal to occur

(Verbruggen *et al.*, 2013, see also **Figure 2** and Appendix 2). Such waiting strategies can be discouraged by reducing the overall probability of a stop signal. For standard stop-signal studies, 25% stop signals is recommended. When researchers prefer a higher percentage of stop signals, additional measures to minimize slowing are required (see Recommendation 5).

#### Recommendation 4: Use the tracking procedure to obtain a broad range of stop-signal delays

If participants can predict when a stop signal will occur within a trial, they might also wait for it. Therefore, a broad range of SSDs is required. The stop-signal delay can be continuously adjusted via a standard adaptive tracking procedure: SSD increases after each successful stop, and decreases after each unsuccessful stop; this converges on a probability of responding [ $p(\text{respond}|\text{signal})$ ]  $\approx 0.50$ . Many studies adjust SSD in steps of 50 ms (which corresponds to three screen ‘refreshes’ for 60 Hz monitors). When step size is too small (for example 16 ms) the tracking may not converge in short experiments, whereas it may not be sensitive enough if step size is too large. Importantly, SSD should decrease after *all* responses on unsuccessful stop trials; this includes premature responses on unsuccessful stop trials (i.e. responses executed before the stop signal was presented) and choice errors on unsuccessful stop trials (e.g. when a left go response would have been executed on the stop trial depicted in **Figure 1**, even though the arrow was pointing to the right).

An adaptive tracking procedure typically results in a sufficiently varied set of SSD values. An additional advantage of the tracking procedure is that fewer stop trials are required to obtain a reliable SSRT estimate (Band *et al.*, 2003). Thus, the tracking procedure is recommended for standard applications.

#### Recommendation 5: Instruct participants not to wait and include block-based feedback

In human studies, task instructions should also be used to discourage waiting. At the very least, participants should be told that ‘*[they] should respond as quickly as possible to the go stimulus and not wait for the stop signal to occur*’ (or something along these lines). To adults, the tracking procedure (if used) can also be explained to further discourage a waiting strategy (i.e. inform participants that the probability of an unsuccessful stop trial will approximate 0.50, and that SSD will increase if they gradually slow their responses).

Inclusion of a practice block in which adherence to instructions is carefully monitored is recommended. In certain populations, such as young children, it might furthermore be advisable to start with a practice block without stop signals to emphasize the importance of the go component of the task.

Between blocks, participants should also be reminded about the instructions. Ideally, this is combined with block-based feedback, informing participants about their mean RT on go trials, number of go omissions (with a reminder that this should be 0), and  $p(\text{respond}|\text{signal})$  (with a reminder that this should be close to .50). The feedback could even include an explicit measure of response slowing.

#### Recommendation 6: Include sufficient trials

The number of stop trials varies widely between studies. Our novel simulation results (see **Figure 2** and Appendix 2) indicate that reliable and unbiased SSRT group-level estimates can be obtained with 50 stop trials (with 25% stop signals in an experiment, this amounts to 200 trials in total. Usually, this corresponds to an experiment of 7–10 min including breaks), but only under ‘optimal’ or very specific circumstances (e.g. when the probability of go omissions is low and the go-RT distribution is not strongly skewed). Lower trial numbers (here we tested 25 stop trials) rarely produced reliable SSRT estimates (and the number of excluded subjects was much higher; see **Figure 2**). Thus, as a general rule of thumb, we recommend to have at least 50 stop trials for standard group-level comparisons. However, it should again be stressed that this may not suffice to obtain reliable individual estimates (which are required for e.g. individual-differences research or diagnostic purposes).

Thus, our simulations reported in Appendix 2 suggest that reliability increases with number of trials. However, in some clinical populations, adding trials may not always be possible (e.g. when patients cannot concentrate for a sufficiently long period of time), and might even be

## Box 2. Failures to trigger the stop process

The race model assumes that the go runner is triggered by the presentation of the go stimulus, and the stop runner by the presentation of the stop signal. However, go omissions (i.e. go trials without a response) are often observed in stop-signal studies. Our preferred SSRT method compensates for such go omissions (see Materials and methods). However, turning to the stopping process, studies using fixed SSDs have found that  $p(\text{respond}|\text{signal})$  at very short delays (including  $\text{SSD} = 0$  ms, when go and stop are presented together) is not always zero; this finding indicates that the stop runner may also not be triggered on all stop trials ('trigger failures').

The non-parametric estimation methods described in Materials and methods (see also Appendix 2) will overestimate SSRT when trigger failures are present on stop trials (**Band et al., 2003**). Unfortunately, these estimation methods cannot determine the presence or absence of trigger failures on stop trials. In order to diagnose in how far trigger failures are present in their data, researchers can include extra stop signals that occur at the same time of the go stimulus (i.e.  $\text{SSD} = 0$ , or shortly thereafter). Note that this number of zero-SSD trials should be sufficiently high to detect (subtle) within- or between-group differences in trigger failures. Furthermore,  $p(\text{respond}|\text{signal})$  should be reported separately for these short-SSD trials, and these trials should not be included when calculating mean SSD or estimating SSRT (see Recommendation one for a discussion of problems that arise when SSDs are very short. Note that the (neural) mechanisms involved in stopping might also partly differ when  $\text{SSD} = 0$ ; see for example **Swick et al., 2011**). Alternatively, researchers can use a parametric method to estimate SSRT. Such methods describe the whole SSRT distribution (unlike the non-parametric methods that estimate summary measures, such as the mean stop latency). Recent variants of such parametric methods also provide an estimate of the probability of trigger failures on stop trials (for the most recent version and specialized software, see **Matzke et al., 2019**).

DOI: <https://doi.org/10.7554/eLife.46323.005>

counterproductive (as strong fluctuations over time can induce extra noise). Our simulations reported in Appendix 3 show that for standard group-level comparisons, researchers can compensate for lower trial numbers by increasing sample size. Above all, we strongly encourage researchers to make informed decisions about number of trials and participants, aiming for sufficiently powered studies. The accompanying open-source simulation code can be used for this purpose.

### When and how to estimate SSRT

Recommendation 7: Do not estimate the SSRT when the assumptions of the race model are violated

SSRTs can be estimated based on the independent race model, which assumes an independent race between a go and a stop runner (**Box 1**). When this independence assumption is (seriously) violated, SSRT estimates become unreliable (**Band et al., 2003**). Therefore, the assumption should be checked. This can be done by comparing the mean RT on unsuccessful stop trials with the mean RT on go trials. Note that this comparison should include all trials with a response (including choice errors and premature responses), and it should be done for each participant and condition separately. SSRT should not be estimated when RT on unsuccessful stop trials is numerically longer than RT on go trials (see also, **Appendix 2—table 1**). More formal and in-depth tests of the race model can be performed (e.g. examining probability of responding and RT on unsuccessful stop trials as a function of delay); however, a large number of stop trials is required for such tests to be meaningful and reliable.



### Box 3. Check-lists for reporting stop-signal studies

The description of every stop-signal study should include the following information:

- Stimuli and materials
  - Properties of the go stimuli, responses, and their mapping
  - Properties of the stop signal
  - Equipment used for testing
- The procedure
  - The number of blocks (including practice blocks)
  - The number of go and stop trials per block
  - Detailed description of the randomization (e.g. is the order of go and stop trials fully randomized or pseudo-randomized?)
  - Detailed description of the tracking procedure (including start value, step size, minimum and maximum value) or the range and proportion of fixed stop-signal delays.
  - Timing of all events. This can include intertrial intervals, fixation intervals (if applicable), stimulus-presentation times, maximum response latency (and whether a trial is terminated when a response is executed or not), feedback duration (in case immediate feedback is presented), etc.
  - A summary of the instructions given to the participant, and any feedback-related information (full instructions can be reported in Supplementary Materials).
  - Information about training procedures (e.g. in case of animal studies)
- The analyses
  - Which trials were included when analyzing go and stop performance
  - Which SSRT estimation method was used (see Materials and methods), providing additional details on the exact approach (e.g. whether or not go omissions were replaced; how go and stop trials with a choice errors—e.g. left response for right arrows—were handled; how the nth quantile was estimated; etc.)
  - Which statistical tests were used for inferential statistics

Stop-signal studies should also report the following descriptive statistics for each group and condition separately (see Appendix 4 for a description of all labels):

- Probability of go omissions (no response)
- Probability of choice errors on go trials
- RT on go trials (mean or median). We recommend to report intra-subject variability as well (especially for clinical studies).
- Probability of responding on a stop trial (for each SSD when fixed delays are used)
- Average stop-signal delay (when the tracking procedure is used); depending on the set-up, it is advisable to report (and use) the 'real' SSDs (e.g. for visual stimuli, the requested SSD may not always correspond to the real SSD due to screen constraints).
- Stop-signal reaction time
- RT of go responses on unsuccessful stop trials

DOI: <https://doi.org/10.7554/eLife.46323.006>

### Recommendation 8: If using a non-parametric approach, estimate SSRT using the integration method (with replacement of go omissions)

Different SSRT estimation methods have been proposed (see Materials and methods). When the tracking procedure is used, the 'mean estimation' method is still the most popular (presumably because it is very easy to use). However, the mean method is strongly influenced by the right tail (skew) of the go RT distribution (see Appendix 2 for examples), as well as by go omissions (i.e. go trials on which no response is executed). The simulations reported in Appendix 2 and summarized in **Figure 2** indicate that the integration method (which replaces go omissions with the maximum RT in order to compensate for the lacking response) is generally less biased and more reliable than the mean method when combined with the tracking procedure. Unlike the mean method, the integration method also does not assume that  $p(\text{respond}|\text{signal})$  is exactly 0.50 (an assumption that is often not met in empirical data). Therefore, we recommend the use of the integration method (with replacement of go omissions) when non-parametric estimation methods are used. We provide software and the source code for this estimation method (and all other recommended measures; Recommendation 12).

Please note that some parametric SSRT estimation methods are less biased than even the best non-parametric methods and avoid other problems that can beset them (see **Box 2**); however, they can be harder for less technically adept researchers to use, and they may require more trials (see **Matzke et al., 2018**, for a discussion).

### Recommendation 9: Refrain from estimating SSRT when the probability of responding on stop trials deviates substantially from 0.50 or when the probability of omissions on go trials is high

Even though the preferred integration method (with replacement of go omissions) is less influenced by deviations in  $p(\text{respond}|\text{signal})$  and go omissions than other methods, it is not completely immune to them either (**Figure 2** and Appendix 2). Previous work suggests that SSRT estimates are most reliable (**Band et al., 2003**) when probability of responding on a stop trial is relatively close to 0.50. Therefore, we recommend that researchers refrain from estimating individual SSRTs when  $p(\text{respond}|\text{signal})$  is lower than 0.25 or higher than 0.75 (**Congdon et al., 2012**). Reliability of the estimates is also influenced by go performance. As the probability of a go omission increases, SSRT estimates also become less reliable. **Figure 2** and the resources described in Appendix 3 can be used to determine an acceptable level of go omissions at a study level. Importantly, researchers should decide on these cut-offs or exclusion criteria before data collection has started.

## How to report stop-signal experiments

### Recommendation 10: Report the methods in enough detail

To allow proper evaluation and replication of the study findings, and to facilitate follow-up studies, researchers should carefully describe the stimuli, materials, and procedures used in the study, and provide a detailed overview of the performed analyses (including a precise description of how SSRT was estimated). This information can be presented in Supplementary Materials in case of journal restrictions. **Box 3** provides a check-list that can be used by authors and reviewers. We also encourage researchers to share their software and materials (e.g. the actual stimuli).

### Recommendation 11: Report possible exclusions in enough detail

As outlined above, researchers should refrain from estimating SSRT when the independence assumptions are seriously violated or when sub-optimal task performance might otherwise compromise the reliability of the estimates. The number of participants for whom SSRT was not estimated should be clearly mentioned. Ideally, dependent variables which are directly observed (see Recommendation 12) are separately reported for the participants that are not included in the SSRT analyses. Researchers should also clearly mention any other exclusion criteria (e.g. outliers based on distributional analyses, acceptable levels of go omissions, etc.), and whether those were set a-priori (analytic plans can be preregistered on a public repository, such as the [Open Science Framework](#); **Nosek et al., 2018**).

## Recommendation 12: Report all relevant behavioral data

Researchers should report all relevant descriptive statistics that are required to evaluate the findings of their stop-signal study (see **Box 3** for a check-list). These should be reported for each group or condition separately. As noted above (Recommendation 7), additional checks of the independent race model can be reported when the number of stop trials is sufficiently high. Finally, we encourage researchers to share their anonymized raw (single-trial) data when possible (in accordance with the FAIR data guidelines; *Wilkinson et al., 2016*).

## Conclusion

Response inhibition and impulse control are central topics in various fields of research, including neuroscience, psychiatry, psychology, neurology, pharmacology, and behavioral sciences, and the stop-signal task has become an essential tool in their study. If properly used, the task can reveal unique information about the underlying neuro-cognitive control mechanisms. By providing clear recommendations, and open-source resources, this paper aims to further increase the quality of research in the response-inhibition and impulse-control domain and to significantly accelerate its progress across the various important domains in which it is routinely applied.

## Materials and methods

The independent race model (**Box 1**) provides two common ‘non-parametric’ methods for estimating SSRT: the integration method and the mean method. Both methods have been used in slightly different flavors in combination with the SSD tracking procedure (see Recommendation 4). Here, we discuss the two most typical estimation variants, which we further scrutinized in our simulations (Appendix 2). We refer the reader to Appendices 2 and 3 for a detailed description of the simulations.

### Integration method (with replacement of go omissions)

In the integration method, the point at which the stop process finishes (**Box 1**) is estimated by ‘integrating’ the RT distribution and finding the point at which the integral equals  $p(\text{respond}|\text{signal})$ . The finishing time of the stop process corresponds to the  $n$ th RT, with  $n = \text{the number of RTs in the RT distribution of go trials multiplied by } p(\text{respond}|\text{signal})$ . When combined with the tracking procedure, overall  $p(\text{respond}|\text{signal})$  is used. For example, when there are 200 go trials, and overall  $p(\text{respond}|\text{signal})$  is 0.45, then the  $n$ th RT is the 90th fastest go RT. SSRT can then be estimated by subtracting mean SSD from the  $n$ th RT. To determine the  $n$ th RT, all go trials with a response are included (*including go trials with a choice error and go trials with a premature response*). Importantly, go omissions (i.e. go trials on which the participant did not respond before the response deadline) are assigned the maximum RT in order to compensate for the lacking response. Premature responses on unsuccessful stop trials (i.e. responses executed before the stop signal is presented) should also be included when calculating  $p(\text{respond}|\text{signal})$  and mean SSD (as noted in Recommendation 4, SSD should also be adjusted after such trials). This version of the integration method produces the most reliable and least biased non-parametric SSRT estimates (Appendix 2).

### The mean method

The mean method uses the mean of the inhibition function (which describes the relationship between  $p(\text{respond}|\text{signal})$  and SSD). Ideally, this mean corresponds to the average SSD obtained with the tracking procedure when  $p(\text{respond}|\text{signal}) = 0.50$  (and often this is taken as a given despite some variation). In other words, the mean method assumes that the mean RT equals  $\text{SSRT} + \text{mean SSD}$ , so SSRT can be estimated easily by subtracting mean SSD from mean RT on go trials when the tracking procedure is used. The ease of use has made this the most popular estimation method. However, our simulations show that this simple version of the mean method is biased and generally less reliable than the integration method with replacement of go omissions.

## Acknowledgements

This work was mainly supported by an ERC Consolidator grant awarded to FV (European Union's Horizon 2020 research and innovation programme, grant agreement No 769595).

## Additional information

### Competing interests

Nicole C Swann: Reviewing editor, *eLife*. Adam R Aron: Reviewing editor, *eLife*. Christian Beste: has received payment for consulting and speaker's honoraria from GlaxoSmithKline, Novartis, Genzyme, and Teva. He has recent research grants with Novartis and Genzyme. Samuel R Chamberlain: consults for Shire, Ieso Digital Health, Cambridge Cognition, and Promentis. Dr Chamberlain's research is funded by Wellcome Trust (110049/Z/15/Z). Trevor W Robbins: consults for Cambridge Cognition, Mundipharma and Unilever. He receives royalties from Cambridge Cognition (CANTAB) and has recent research grants with Shionogi and SmallPharma. Katya Rubia: has received speaker's honoraria and grants for other projects from Eli Lilly and Shire. Russell J Schachar: has consulted to Highland Therapeutics, Eli Lilly and Co., and Purdue Pharma. He has commercial interest in a cognitive rehabilitation software company, eHave. The other authors declare that no competing interests exist.

### Funding








Funder	Grant reference number	Author
H2020 European Research Council	769595	Frederick Verbruggen

The funders had no role in study design, data collection and interpretation, or the decision to submit the work for publication.

### Author contributions

Frederick Verbruggen, Conceptualization, Resources, Data curation, Software, Formal analysis, Supervision, Funding acquisition, Validation, Investigation, Visualization, Methodology, Writing—original draft, Project administration, Writing—review and editing; Adam R Aron, Christian Beste, Patrick G Bissett, Adam T Brockett, Joshua W Brown, Samuel R Chamberlain, Christopher D Chambers, Hans Colonius, Lorenza S Colzato, Brian D Corneil, James P Coxon, Annie Dupuis, Dawn M Eagle, Hugh Garavan, Ian Greenhouse, René J Huster, Sara Jahfari, J Leon Kenemans, Inge Leunissen, Chiang-Shan R Li, Dora Matzke, Sharon Morein-Zamir, Aditya Murthy, Martin Paré, Russell A Poldrack, K Richard Ridderinkhof, Trevor W Robbins, Matthew Roesch, Katya Rubia, Russell J Schachar, Jeffrey D Schall, Ann-Kathrin Stock, Nicole C Swann, Katharine N Thakkar, Maurits W van der Molen, Matthijs Vink, Jan R Wessel, Robert Whelan, Bram B Zandbelt, Conceptualization, Writing—review and editing; Guido PH Band, Andrew Heathcote, Gordon D Logan, Conceptualization, Methodology, Writing—review and editing; Luc Vermeylen, Conceptualization, Resources, Software, Writing—review and editing; C Nico Boehler, Conceptualization, Resources, Software, Formal analysis, Validation, Investigation, Visualization, Methodology, Writing—original draft, Writing—review and editing

### Author ORCIDs

Frederick Verbruggen  <https://orcid.org/0000-0002-7958-0719>  
 Adam T Brockett  <http://orcid.org/0000-0001-7712-5053>  
 Hans Colonius  <http://orcid.org/0000-0002-9733-6939>  
 Brian D Corneil  <http://orcid.org/0000-0002-4702-7089>  
 James P Coxon  <http://orcid.org/0000-0003-2351-8489>  
 Ian Greenhouse  <http://orcid.org/0000-0003-1467-739X>  
 Sara Jahfari  <http://orcid.org/0000-0002-1979-589X>  
 Russell A Poldrack  <http://orcid.org/0000-0001-6755-0259>

Matthew Roesch  <https://orcid.org/0000-0003-2854-6593>

Nicole C Swann  <https://orcid.org/0000-0003-2463-5134>

Jan R Wessel  <http://orcid.org/0000-0002-7298-6601>

C Nico Boehler  <http://orcid.org/0000-0001-5963-2780>

### Decision letter and Author response

Decision letter <https://doi.org/10.7554/eLife.46323.026>

Author response <https://doi.org/10.7554/eLife.46323.027>

## Additional files

### Supplementary files

- Transparent reporting form

DOI: <https://doi.org/10.7554/eLife.46323.007>

### Data availability

The code used for the simulations and all simulated data can be found on Open Science Framework (<https://osf.io/rmqaw/>).

The following dataset was generated:

Author(s)	Year	Dataset title	Dataset URL	Database and Identifier
Verbruggen F	2019	Race model simulations to determine estimation bias and reliability of SSRT estimates	<a href="https://dx.doi.org/10.17605/OSF.IO/JWSF9">https://dx.doi.org/10.17605/OSF.IO/JWSF9</a>	Open Science Framework, 10.17605/OSF.IO/JWSF9

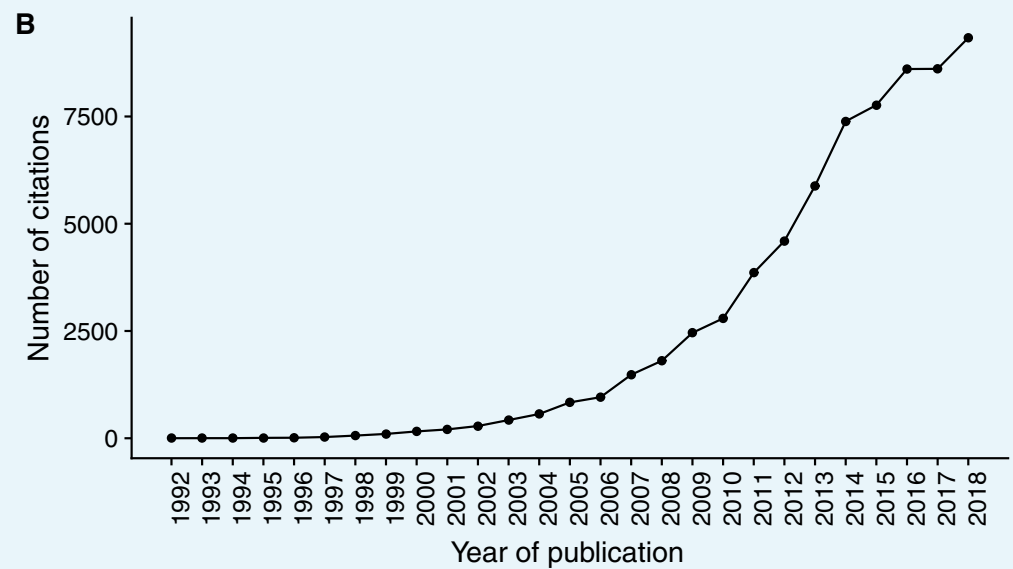
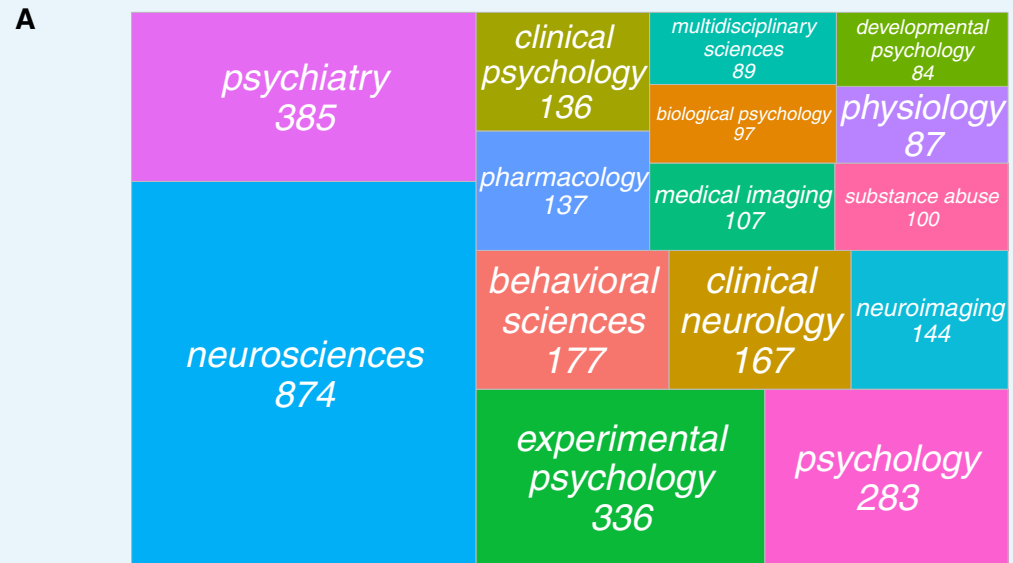
## References

- Aron AR. 2011. From reactive to proactive and selective control: developing a richer model for stopping inappropriate responses. *Biological Psychiatry* **69**:e55–e68. DOI: <https://doi.org/10.1016/j.biopsych.2010.07.024>, PMID: 20932513
- Band GP, van der Molen MW, Logan GD. 2003. Horse-race model simulations of the stop-signal procedure. *Acta Psychologica* **112**:105–142. DOI: [https://doi.org/10.1016/S0001-6918\(02\)00079-3](https://doi.org/10.1016/S0001-6918(02)00079-3), PMID: 12521663
- Bari A, Robbins TW. 2013. Inhibition and impulsivity: behavioral and neural basis of response control. *Progress in Neurobiology* **108**:44–79. DOI: <https://doi.org/10.1016/j.pneurobio.2013.06.005>, PMID: 23856628
- Bartholdy S, Dalton B, O'Daly OG, Campbell IC, Schmidt U. 2016. A systematic review of the relationship between eating, weight and inhibitory control using the stop signal task. *Neuroscience & Biobehavioral Reviews* **64**:35–62. DOI: <https://doi.org/10.1016/j.neubiorev.2016.02.010>, PMID: 26900651
- Bissett PG, Logan GD. 2014. Selective stopping? maybe not. *Journal of Experimental Psychology: General* **143**:455–472. DOI: <https://doi.org/10.1037/a0032122>, PMID: 23477668
- Boucher L, Palmeri TJ, Logan GD, Schall JD. 2007. Inhibitory control in mind and brain: an interactive race model of countermanding saccades. *Psychological Review* **114**:376–397. DOI: <https://doi.org/10.1037/0033-295X.114.2.376>, PMID: 17500631
- Chambers CD, Garavan H, Bellgrove MA. 2009. Insights into the neural basis of response inhibition from cognitive and clinical neuroscience. *Neuroscience & Biobehavioral Reviews* **33**:631–646. DOI: <https://doi.org/10.1016/j.neubiorev.2008.08.016>, PMID: 18835296
- Colonius H, Diederich A. 2018. Paradox resolved: Stop signal race model with negative dependence. *Psychological Review* **125**:1051–1058. DOI: <https://doi.org/10.1037/rev0000127>, PMID: 30272461
- Congdon E, Mumford JA, Cohen JR, Galvan A, Canli T, Poldrack RA. 2012. Measurement and reliability of response inhibition. *Frontiers in Psychology* **3**. DOI: <https://doi.org/10.3389/fpsyg.2012.00037>, PMID: 22363308
- Lappin JS, Eriksen CW. 1966. Use of a delayed signal to stop a visual reaction-time response. *Journal of Experimental Psychology* **72**:805–811. DOI: <https://doi.org/10.1037/h0021266>
- Leunissen I, Zandbelt BB, Potocanac Z, Swinnen SP, Coxon JP. 2017. Reliable estimation of inhibitory efficiency: to anticipate, choose or simply react? *European Journal of Neuroscience* **45**:1512–1523. DOI: <https://doi.org/10.1111/ejn.13590>, PMID: 28449195
- Lipszyc J, Schachar R. 2010. Inhibitory control and psychopathology: a meta-analysis of studies using the stop signal task. *Journal of the International Neuropsychological Society* **16**:1064–1076. DOI: <https://doi.org/10.1017/S1355617710000895>, PMID: 20719043
- Logan GD, Van Zandt T, Verbruggen F, Wagenmakers EJ. 2014. On the ability to inhibit thought and action: general and special theories of an act of control. *Psychological Review* **121**:66–95. DOI: <https://doi.org/10.1037/a0035230>, PMID: 24490789

- Logan GD**, Yamaguchi M, Schall JD, Palmeri TJ. 2015. Inhibitory control in mind and brain 2.0: blocked-input models of saccadic countermanding. *Psychological Review* **122**:115–147. DOI: <https://doi.org/10.1037/a0038893>, PMID: 25706403
- Logan GD**, Cowan WB. 1984. On the ability to inhibit thought and action: A theory of an act of control. *Psychological Review* **91**:295–327. DOI: <https://doi.org/10.1037/0033-295X.91.3.295>
- Matzke D**, Verbruggen F, Logan GD. 2018. The Stop-Signal Paradigm. In: Wixted J. T (Ed). *Stevens' Handbook of Experimental Psychology and Cognitive Neuroscience*. John Wiley & Sons, Inc. DOI: <https://doi.org/10.1002/9781119170174.epcn510>
- Matzke D**, Curley S, Gong CQ, Heathcote A. 2019. Inhibiting responses to difficult choices. *Journal of experimental psychology. General* **148**:124142. DOI: <https://doi.org/10.1037/xge0000525>, PMID: 30596441
- Nelson MJ**, Boucher L, Logan GD, Palmeri TJ, Schall JD. 2010. Nonindependent and nonstationary response times in stopping and stepping saccade tasks. *Attention, Perception & Psychophysics* **72**:1913–1929. DOI: <https://doi.org/10.3758/APP.72.7.1913>, PMID: 20952788
- Nosek BA**, Ebersole CR, DeHaven AC, Mellor DT. 2018. The preregistration revolution. *PNAS* **115**:2600–2606. DOI: <https://doi.org/10.1073/pnas.1708274114>
- R Development Core Team**. 2017. *R: a language and environment for statistical computing*. R Foundation for Statistical Computing. 3.4.2. Vienna, Austria: <http://www.r-project.org/>.
- Rigby RA**, Stasinopoulos DM. 2005. Generalized additive models for location, scale and shape (with discussion). *Journal of the Royal Statistical Society: Series C* **54**:507–554. DOI: <https://doi.org/10.1111/j.1467-9876.2005.00510.x>
- Schall JD**, Palmeri TJ, Logan GD. 2017. Models of inhibitory control. *Philosophical Transactions of the Royal Society B: Biological Sciences* **372**:20160193. DOI: <https://doi.org/10.1098/rstb.2016.0193>
- Smith JL**, Mattick RP, Jamadar SD, Iredale JM. 2014. Deficits in behavioural inhibition in substance abuse and addiction: a meta-analysis. *Drug and Alcohol Dependence* **145**:1–33. DOI: <https://doi.org/10.1016/j.drugalcdep.2014.08.009>, PMID: 25195081
- Swick D**, Ashley V, Turken U. 2011. Are the neural correlates of stopping and not going identical? quantitative meta-analysis of two response inhibition tasks. *NeuroImage* **56**:1655–1665. DOI: <https://doi.org/10.1016/j.neuroimage.2011.02.070>, PMID: 21376819
- Tannock R**, Schachar RJ, Carr RP, Chajczyk D, Logan GD. 1989. Effects of methylphenidate on inhibitory control in hyperactive children. *Journal of Abnormal Child Psychology* **17**:473–491. DOI: <https://doi.org/10.1007/BF00916508>, PMID: 2681316
- Verbruggen F**, Chambers CD, Logan GD. 2013. Fictitious inhibitory differences: how skewness and slowing distort the estimation of stopping latencies. *Psychological Science* **24**. DOI: <https://doi.org/10.1177/0956797612457390>, PMID: 23399493
- Verbruggen F**, Logan GD. 2015. Evidence for capacity sharing when stopping. *Cognition* **142**:81–95. DOI: <https://doi.org/10.1016/j.cognition.2015.05.014>, PMID: 26036922
- Verbruggen F**, Logan GD. 2017. Control in response inhibition. In: Egnér T (Ed). *The Wiley Handbook of Cognitive Control*. Wiley. DOI: <https://doi.org/10.1002/9781118920497.ch6>
- Vince MA**. 1948. The intermittency of control movements and the psychological refractory period1. *British Journal of Psychology. General Section* **38**:149–157. DOI: <https://doi.org/10.1111/j.2044-8295.1948.tb01150.x>
- Wickham H**. 2016. *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer. DOI: <https://doi.org/10.1007/978-0-387-98141-3>
- Wilkinson MD**, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, Blomberg N, Boiten JW, da Silva Santos LB, Bourne PE, Bouwman J, Brookes AJ, Clark T, Crosas M, Dillo I, Dumon O, Edmunds S, Evelo CT, Finkers R, Gonzalez-Beltran A, et al. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data* **3**:160018. DOI: <https://doi.org/10.1038/sdata.2016.18>, PMID: 26978244

# Appendix 1

DOI: <https://doi.org/10.7554/eLife.46323.003>



**Appendix 1—figure 1.** The number of stop-signal publications per research area (Panel A) and the number of articles citing the ‘stop-signal task’ per year (Panel B). Source: Web of Science, 27/01/2019, search term: ‘topic = stop signal task’. The research areas in Panel A are also taken from Web of Science.

DOI: <https://doi.org/10.7554/eLife.46323.010>

## Appendix 2

DOI: <https://doi.org/10.7554/eLife.46323.003>

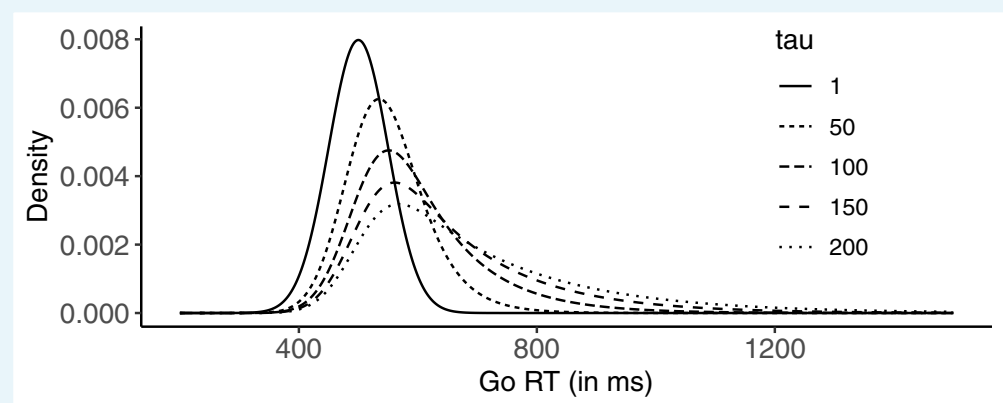
### Race model simulations to determine estimation bias and reliability of SSRT estimates

#### Simulation procedure

To compare different SSRT estimation methods, we ran a set of simulations which simulated performance in the stop-signal task based on assumptions of the independent race model: on stop trials, a response was deemed to be stopped (successful stop) when the RT was larger than SSRT + SSD; a response was deemed to be executed (unsuccessful stop) when RT was smaller than SSRT + SSD. Go and stop were completely independent.

All simulations were done using R (*R Development Core Team, 2017*, version 3.4.2). Latencies of the go and stop runners were sampled from an ex-Gaussian distribution, using the `rexGaus` function (*Rigby and Stasinopoulos, 2005*, version 5.1.2). The ex-Gaussian distribution has a positively skewed unimodal shape and results from a convolution of a normal (Gaussian) distribution and an exponential distribution. It is characterized by three parameters:  $\mu$  (mean of the Gaussian component),  $\sigma$  (SD of Gaussian component), and  $\tau$  (both the mean and SD of the exponential component). The mean of the ex-Gaussian distribution =  $\mu + \tau$ , and variance =  $\sigma^2 + \tau^2$ . Previous simulation studies of the stop-signal task also used ex-Gaussian distributions to model their reaction times (e.g. *Band et al., 2003; Verbruggen et al., 2013; Matzke et al., 2019*).

For each simulated 'participant',  $\mu_{go}$  of the ex-Gaussian go RT distribution was sampled from a normal distribution with mean = 500 (i.e. the population mean) and SD = 50, with the restriction that it was larger than 300 (see *Verbruggen et al., 2013*, for a similar procedure).  $\sigma_{go}$  was fixed at 50, and  $\tau_{go}$  was either 1, 50, 100, 150, and 200 (resulting in increasingly skewed distributions). The RT cut-off was set at 1,500 ms. Thus, go trials with an RT >1,500 ms were considered go omissions. For some simulations, we also inserted extra go omissions, resulting in five 'go omission' conditions: 0% inserted go omissions (although the occasional go omission was still possible when  $\tau_{go}$  was high), 5%, 10%, 15%, or 20%. These go omissions were randomly distributed across go and stop trials. For the 5%, 10%, 15%, and 20% go-omission conditions, we first checked if there were already go omissions due to the random sampling from the ex-Gaussian distribution. If such go omissions occurred 'naturally', fewer 'artificial' omissions were inserted.



**Appendix 2—figure 1.** Examples of ex-Gaussian (RT) distributions used in our simulations. For all distributions,  $\mu_{go} = 500$  ms, and  $\sigma_{go} = 50$  ms.  $\tau_{go}$  was either 1, 50, 100, 150, and 200 (resulting in increasingly skewed distributions). Note that for a given RT cut-off (1,500 ms in the simulations), cut-off-related omissions are rare, but systematically more likely as tau



increases. In addition to such ‘natural’ go omissions, we introduced ‘artificial’ ones in the different go-omission conditions of the simulations (not depicted).

DOI: <https://doi.org/10.7554/eLife.46323.012>

For each simulated ‘participant’,  $\mu_{stop}$  of the ex-Gaussian SSRT distribution was sampled from a normal distribution with mean = 200 (i.e. the population mean) and SD = 20, with the restriction that it was larger than 100.  $\sigma_{stop}$  and  $\tau_{stop}$  were fixed at 20 and 10, respectively. For each ‘participant’, the start value of SSD was 300 ms, and was continuously adjusted using a standard tracking procedure (see main text) in steps of 50 ms. In the present simulations, we did not set a minimum or maximum SSD.

The total number of trials simulated per participant was either 100, 200, 400, or 800, whereas the probability of a stop signal was fixed at .25; thus, the number of stop trials was 25, 50, 100, or 200, respectively. This resulted in 5 (go omission: 0, 5, 10, 15, or 20%)  $\times$  5 ( $\tau_{go}$ : 1, 50, 100, 150, 200)  $\times$  4 (total number of trials: 100, 200, 400, 800) conditions. For each condition, we simulated 1000 participants. Overall, this resulted in 100,000 participants (and 375,000,000 trials).

The code used for the simulations and all simulated data can be found on Open Science Framework (<https://osf.io/rmqaw/>).

## Analyses

We performed three sets of analyses. First, we checked if RT on unsuccessful stop trials was numerically shorter than RT on go trials. Second, we estimated SSRTs using the two estimation methods described in the main manuscript (Materials and methods), and two other methods that have been used in the stop-signal literature. The first additional approach is a variant of the integration method described in the main manuscript. The main difference is the exclusion of go omissions (and sometimes choice errors on unsuccessful stop trials) from the go RT distribution when determining the nth RT. The second additional variant also does not assign go omissions the maximum RT. Rather, this method adjusts  $p(\text{respond}|\text{signal})$  to compensate for go omissions (**Tannock et al., 1989**):

$$p(\text{respond}|\text{signal})_{\text{adjusted}} = 1 - \frac{p(\text{inhibit}|\text{signal}) - p(\text{omission}|\text{go})}{1 - p(\text{omission}|\text{go})}$$

The nth RT is then determined using the adjusted  $p(\text{respond}|\text{signal})$  and the distribution of RTs of all go trials with a response.

Thus, we estimated SSRT using four different methods: (1) integration method with replacement of go omissions; (2) integration method with exclusion of go omissions; (3) integration method with adjustment of  $p(\text{respond}|\text{signal})$ ; and (4) the mean method. For each estimation method and condition (go omission  $\times$   $\tau_{go}$   $\times$  number of trials), we calculated the difference between the estimated SSRT and the actual SSRT; positive values indicate that SSRT is overestimated, whereas negative values indicate that SSRT is underestimated. For each estimation method, we also correlated the true and estimated values across participants; higher values indicate more reliable SSRT estimates.

We investigated all four mentioned estimation approaches in the present appendix. In the main manuscript, we provide a detailed overview focussing on (1) the integration method with replacement of go omissions and (2) the mean method. As described below, the integration method with replacement of go omissions was the least biased and most reliable, but we also show the mean method in the main manuscript to further highlight the issues that arise when this (still popular) method is used.

## Results

All figures were produced using the ggplot2 package (version 3.1.0 **Wickham, 2016**). The number of excluded ‘participants’ (i.e. RT on unsuccessful stop trials > RT on go trials) is presented in **Figure 2** of the main manuscript. Note that these are only apparent violations of

the independent race model, as go and stop were always modelled as independent runners. Instead, the longer RTs on unsuccessful stop trials result from estimation uncertainty associated with estimating mean RTs using scarce data. However, as true SSRT of all participants was known, we could nevertheless compare the SSRT bias for included and excluded participants. As can be seen in the table below, estimates were generally much more biased for 'excluded' participants than for 'included' participants. Again this indicates that extreme data are more likely to occur when the number of trials is low.

**Appendix 2—table 1.** The mean difference between estimated and true SSRT for participants who were included in the main analyses and participants who were excluded (because average RT on unsuccessful stop trials > average RT on go trials). We did this only for  $\tau_{go} = 1$  or 50,  $p(\text{go omission}) = 10, 15, \text{ or } 20$ , and number of trials = 100 (i.e. when the number of excluded participants was high; see Panel A, **Figure 2** of the main manuscript).

Estimation method	Included	Excluded
Integration with replacement of go omissions	-6.4	-35.8
Integration without replacement of go omissions	-19.4	-48.5
Integration with adjusted $p(\text{respond} \text{signal})$	12.5	-17.4
Mean	-16.0	-46.34

DOI: <https://doi.org/10.7554/eLife.46323.013>

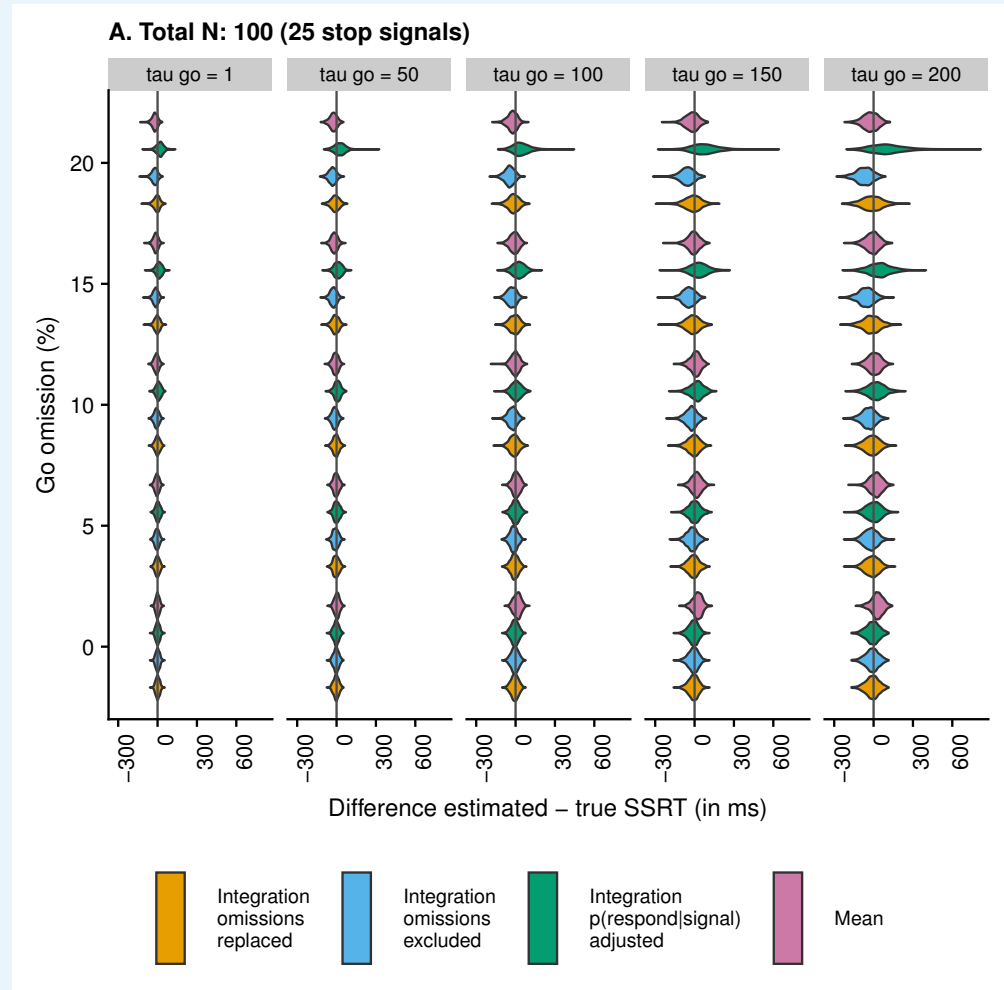
To further compare differences between estimated and true SSRTs for the included participants, we used 'violin plots'. These plots show the distribution and density of SSRT difference values. We created separate plots as a function of the total number of trials (100, 200, 400, and 800), and each plot shows the SSRT difference as a function of estimation method, percentage of go omissions, and  $\tau_{go}$  (i.e. the skew of the RT distribution on go trials; see **Appendix 2—figure 1**). The plots can be found below. The first important thing to note is that the scales differ between subplots. This was done intentionally, as the distribution of difference scores was wider when the number of trials was lower (with fixed scales, it is difficult to detect meaningful differences between estimation methods and conditions for higher trial numbers; i.e. Panels C and D). In other words, low trial numbers will produce more variable and less reliable SSRT estimates.

Second, the violin plots show that SSRT estimates are strongly influenced by an increasing percentage of go omissions. The figures show that the integration method with replacement of go omissions, integration method with exclusion of go omissions, and the mean method all have a tendency to underestimate SSRT as the percentage of go omissions increases; importantly, *this underestimation bias is most pronounced for the integration method with exclusion of go omissions*. By contrast, the integration method which uses the adjusted  $p(\text{respond}|\text{signal})$  will overestimate SSRT when go omissions are present; compared with the other methods, this bias was the strongest in absolute terms.

Consistent with previous work (**Verbruggen et al., 2013**), skew of the RT distribution also strongly influenced the estimates. SSRT estimates were generally more variable as  $\tau_{go}$  increased. When the probability of a go omission was low, the integration methods showed a small underestimation bias for high levels of  $\tau_{go}$ , whereas the mean method showed a clear overestimation bias for high levels of  $\tau_{go}$ . In absolute terms, this overestimation bias for the mean method was more pronounced than the underestimation bias for the integration methods. For higher levels of go omissions, the pattern became more complicated as the various biases started to interact. Therefore, we also correlated the true SSRT with the estimated SSRT to compare the different estimation methods.

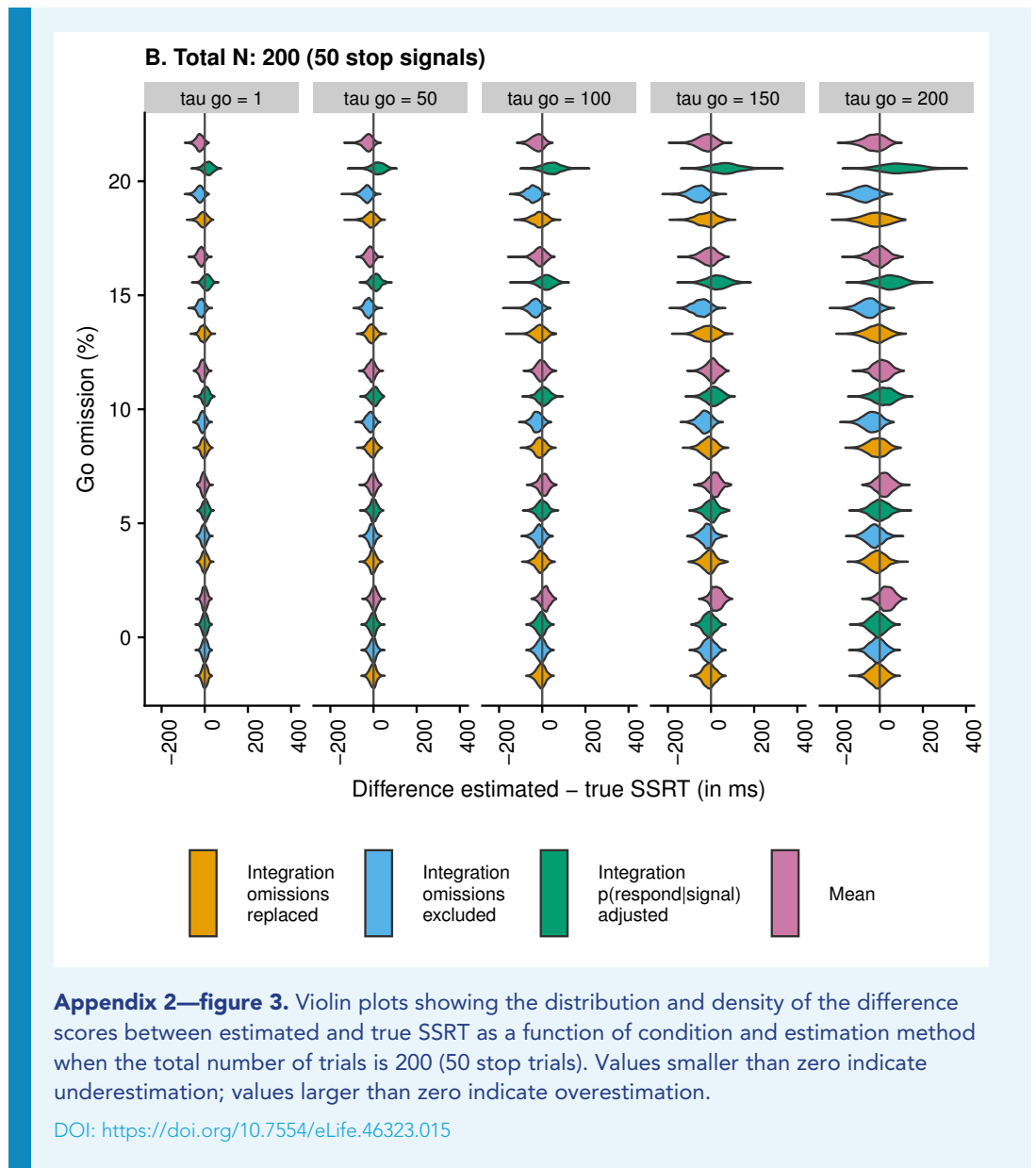
To calculate the correlation between true and estimated SSRT for each method, we collapsed across all combinations of  $\tau_{go}$ , go-omission rate, and number of trials. The correlation (i.e. reliability of the estimate) was highest for the integration method with replacement of go omissions,  $r = 0.57$  (as shown in the violin plots, this was also the least biased method); intermediate for the mean method,  $r = 0.53$ , and the integration method

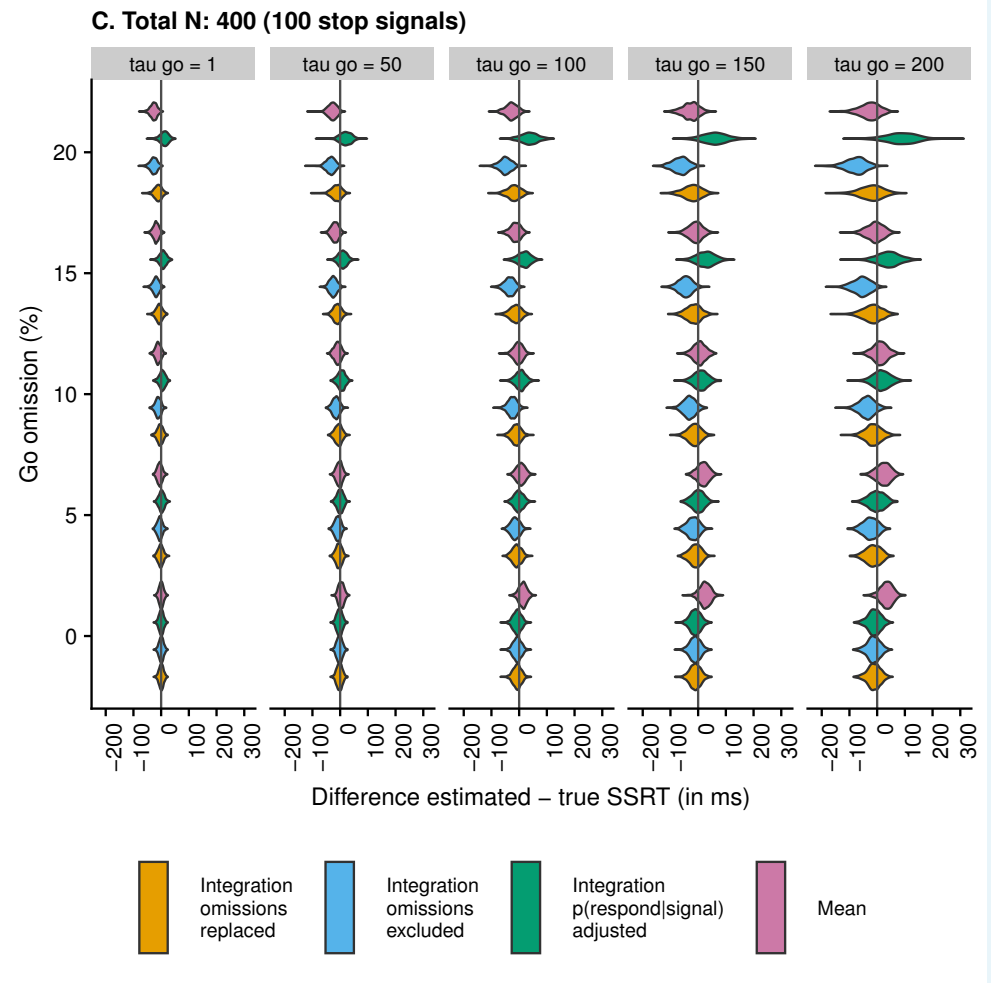
with exclusion of go errors,  $r = 0.51$ ; and lowest for the integration method using adjusted  $p$  (respond|signal),  $r = 0.43$ .



**Appendix 2—figure 2.** Violin plots showing the distribution and density of the difference scores between estimated and true SSRT as a function of condition and estimation method when the total number of trials is 100 (25 stop trials). Values smaller than zero indicate underestimation; values larger than zero indicate overestimation.

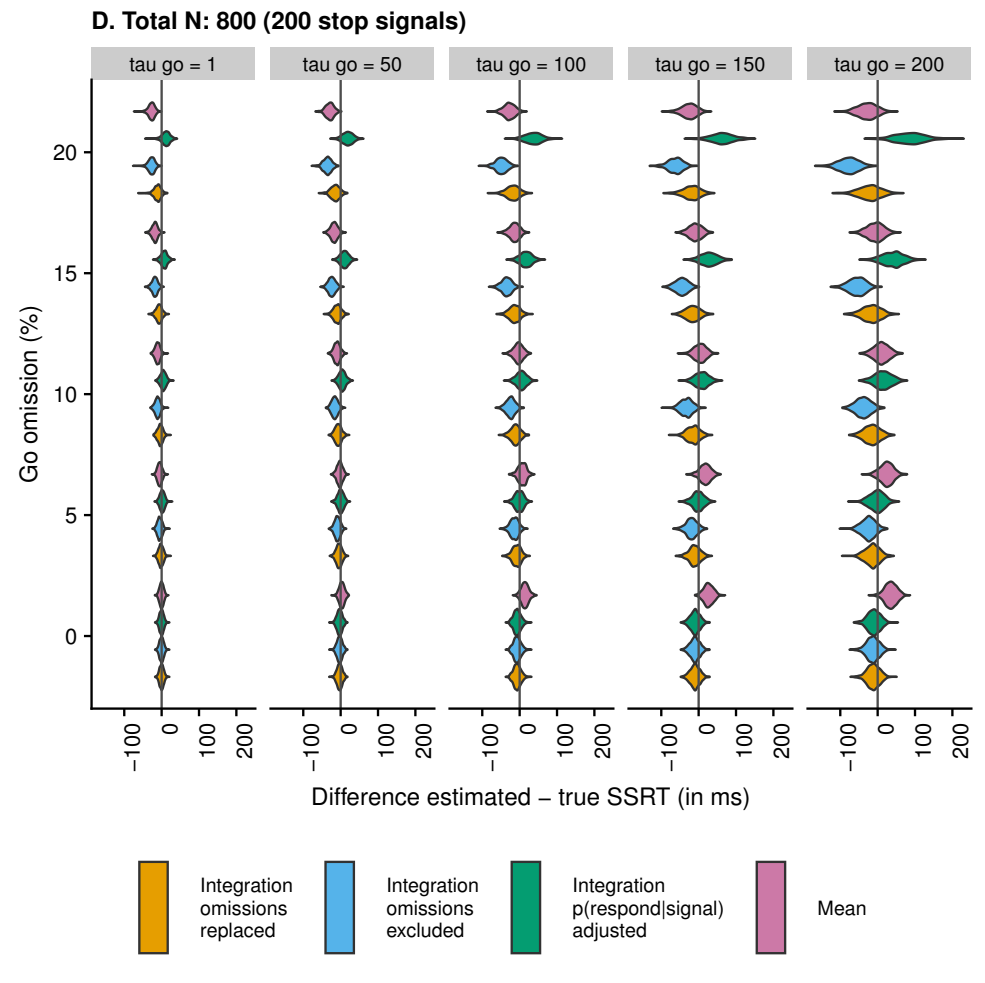
DOI: <https://doi.org/10.7554/eLife.46323.014>





**Appendix 2—figure 4.** Violin plots showing the distribution and density of the difference scores between estimated and true SSRT as a function of condition and estimation method when the total number of trials is 400 (100 stop trials). Values smaller than zero indicate underestimation; values larger than zero indicate overestimation.

DOI: <https://doi.org/10.7554/eLife.46323.016>



**Appendix 2—figure 5.** Violin plots showing the distribution and density of the difference scores between estimated and true SSRT as a function of condition and estimation method when the total number of trials is 800 (200 stop trials). Values smaller than zero indicate underestimation; values larger than zero indicate overestimation.

DOI: <https://doi.org/10.7554/eLife.46323.017>

## Appendix 3

DOI: <https://doi.org/10.7554/eLife.46323.003>

### Race model simulations to determine achieved power

#### Simulation procedure

To determine how different parameters affected the power to detect SSRT differences, we simulated 'experiments'. We used the same general procedure as described in Appendix 2. In the example described below, we used a simple between-groups design with a control group and an experimental group.

For each simulated 'participant' of the 'control group',  $\mu_{go}$  of the ex-Gaussian go RT distribution was sampled from a normal distribution with mean = 500 (i.e. the population mean) and SD = 100, with the restriction that it was larger than 300.  $\sigma_{go}$  and  $\tau_{go}$  were both fixed at 50, and the percentage of (artificially inserted) go omissions was 0% (see Appendix 2).  $\mu_{stop}$  of the ex-Gaussian SSRT distribution was also sampled from a normal distribution with mean = 200 (i.e. the population mean) and SD = 40, with the restriction that it was larger than 100.  $\sigma_{stop}$  and  $\tau_{stop}$  were fixed at 20 and 10, respectively. Please note that the SDs for the population means were higher than the values used for the simulations reported in Appendix 2 to allow for extra between-subjects variation in our groups.

For the 'experimental group', the go and stop parameters could vary across 'experiments'.  $\mu_{go}$  was sampled from a normal distribution with population mean = 500, 525, or 575 (SD = 100).  $\sigma_{go}$  was 50, 52.5, or 57.5 (for population mean of  $\mu_{go}$  = 500, 525, and 575, respectively), and  $\tau_{go}$  was either 50, 75, or 125 (also for population mean of  $\mu_{go}$  = 500, 525, and 575, respectively). Remember that the mean of the ex-Gaussian distribution =  $\mu + \tau$  (Appendix 2). Thus, mean go RT of the experimental group was either 550 ms (500 + 50, which is the same as the control group), 600 (525 + 75), or 700 (575 + 125). The percentage of go omissions for the experimental group was either 0% (the same as the experimental group), 5% (for  $\mu_{go}$  = 525) or 10% (for  $\mu_{go}$  = 575).

**Appendix 3—table 1.** Parameters of the go distribution for the control group and the three experimental conditions. SSRT of all experimental groups differed from SSRT in the control group (see below).

Parameters of go distribution	Control	Experimental 1	Experimental 2	Experimental 3
$\mu_{go}$	500	500	525	575
$\sigma_{go}$	50	50	52.5	57.5
$\tau_{go}$	50	50	75	125
go omission	0	0	5	10

DOI: <https://doi.org/10.7554/eLife.46323.019>

$\mu_{stop}$  of the 'experimental-group' SSRT distribution was sampled from a normal distribution with mean = 210 or 215 (SD = 40).  $\sigma_{stop}$  was 21 or 21.5 (for  $\mu_{stop}$  = 210 and 215, respectively), and  $\tau_{stop}$  was either 15 or 20 (for  $\mu_{stop}$  = 210 and 215, respectively). Thus, mean SSRT of the experimental group was either 225 ms (210 + 15, corresponding to a medium effect size; Cohen's  $d \approx 0.50$ – $0.55$ ). Note that the exact value could differ slightly between simulations as random samples were taken) or 235 (215 + 20, corresponding to a large effect size; Cohen's  $d \approx 0.85$ – $0.90$ ). SSRT varied independently from the go parameters (i.e.  $\mu_{go}$  +  $\tau_{go}$ , and % go omissions).

The total number of trials per experiment was either 100 (25 stop trials), 200 (50 stop trials) or 400 (100 stop trials). Other simulation parameters were the same as those described in Appendix 2. Overall, this resulted in 18 different combinations: 3 (go difference between control and experimental; see **Appendix 3—table 1** above) x 2 (mean SSRT difference

between control and experimental: 15 or 30) x 3 (total number of trials: 100, 200 or 400). For each parameter combination, we simulated 5000 'pairs' of subjects.

The code and results of the simulations are available via the Open Science Framework (<https://osf.io/rmqaw/>); stop-signal users can adjust the scripts (e.g. by changing parameters or even the design) to determine the required sample size given some consideration about the expected results. Importantly, the present simulation code provides access to a wide set of parameters (i.e. go omission, parameters of the go distribution, and parameters of the SSRT distribution) that could differ across groups or conditions.

## Analyses

SSRTs were estimated using the integration method with replacement of go omissions (i.e. the method that came out on top in the other set of simulations). Once the SSRTs were estimated, we randomly sampled 'pairs' to create the two groups for each 'experiment'. For the 'medium' SSRT difference (i.e. 210 vs. 225 ms), group size was either 32, 64, 96, 128, 160, or 192 (the total number of participants per experiment was twice the group size). For the 'large' SSRT difference (i.e. 210 vs. 235 ms), group size was either 16, 32, 48, 64, 80, or 96 (the total number of participants per experiment was twice the group size). For each sample size and parameter combination (see above), we repeated this procedure 1000 times (or 1000 experiments).

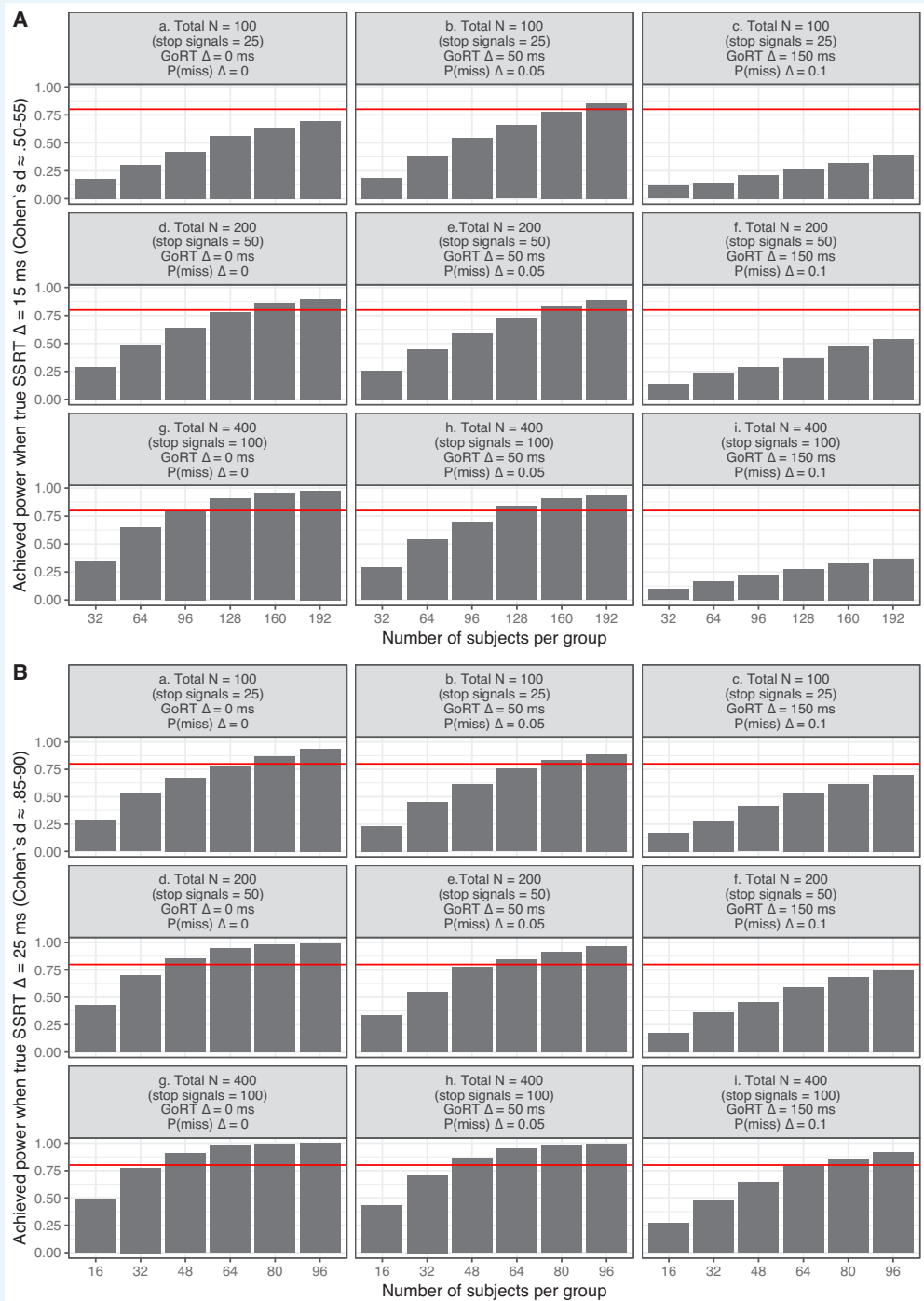
For each experiment, we subsequently compared the estimated SSRTs of the control and experiment groups with an independent-samples t-test (assuming unequal variances). Then we determined for each sample size x parameter combination the proportion of t-tests that were significant (with  $\alpha = 0.05$ ).

## Results

The figure below plots achieved power as a function of sample size (per group), experimental vs. control group difference in true SSRT, and group differences in go performance. Note that if true and estimated SSRTs would exactly match (i.e. estimations reliability = 1), approximately 58 participants per group would be required to detect a medium-sized true SSRT difference with power = 0.80 (i.e. when Cohen's  $d \approx 0.525$ ), and 22 participants per group for a large-sized true SSRT difference (Cohen's  $d \approx 0.875$ ).

Inspection of the figure clearly reveals that achieved power generally increases when sample size and number of trials increase. Obviously achieved power is also strongly dependent on effect size (Panel A vs. B). Interestingly, the figure also shows that the ability to detect SSRT differences is reduced when go performance of the groups differ substantially (see second and third columns of Panel A). As noted in the main manuscript and Appendix 2, even the integration method (with replacement of go omissions) is not immune to changes in the go performance. More specifically, SSRT will be underestimated when the RT distribution is skewed (note that all other approaches produce even stronger biases). In this example, the underestimation bias will reduce the observed SSRT difference (as the underestimation bias is stronger for the experimental group than for the control group). Again, this highlights the need to encourage consistent fast responding (reducing the right-end tail of the distribution).





**Appendix 3—figure 1.** Achieved power for an independent two-groups design as function of differences in go omission, go distribution, SSRT distribution, and the number of trials in the ‘experiments’.

DOI: <https://doi.org/10.7554/eLife.46323.020>

## Appendix 4

DOI: <https://doi.org/10.7554/eLife.46323.003>

### Overview of the main labels and common alternatives

Label	Description	Common alternative labels
Stop-signal task	A task used to measure response inhibition in the lab. Consists of a go component (e.g. a two-choice discrimination task) and a stop component (suppressing the response when an extra signal appears).	Stop-signal reaction time task, stop-signal paradigm, countermanding task
Go trial	On these trials (usually the majority), participants respond to the go stimulus as quickly and accurately as possible (e.g. left arrow = left key, right arrow = right key).	No-signal trial, no-stop-signal trial
Stop trial	On these trials (usually the minority), an extra signal is presented after a variable delay, instructing participants to stop their response to the go stimulus.	Stop-signal trial, signal trial
Successful stop trial	On these stop trials, the participants successfully stopped (inhibited) their go response.	Stop-success trial, signal-inhibit trial, canceled trial
Unsuccessful stop trial	On these stop trials, the participants could not inhibit their go response; hence, they responded despite the (stop-signal) instruction not to do so.	Stop-failure trial, signal-respond trial, noncanceled trial, stop error
Go omission	Go trials without a go response.	Go-omission error, misses, missed responses
Choice errors on go trials	Incorrect response on a go trial (e.g. the go stimulus required a left response but a right response was executed).	(Go) errors, incorrect (go or no-signal) trials
Premature response on a go trial	A response executed before the presentation of the go stimulus on a go trial. This can happen when go-stimulus presentation is highly predictable in time (and stimulus identity is not relevant to the go task; e.g. in a simple detection task) or when participants are 'impulsive'. Note that response latencies will be negative on such trials.	
$p(\text{respond} \text{signal})$	Probability of responding on a stop trial. Non-parametric estimation methods (Materials and Methods) use $p(\text{respond} \text{signal})$ to determine SSRT.	$P(\text{respond})$ , response rate, $p(\text{inhibit})=1 - p(\text{respond} \text{signal})$
Choice errors on unsuccessful stop trials	Unsuccessful stop trials on which the incorrect go response was executed (e.g. the go stimulus required a left response but a right response was executed).	Incorrect signal-respond trials
Premature responses on unsuccessful stop trials	This is a special case of unsuccessful stop trials, referring to go responses executed <i>after</i> the presentation of the go stimulus but <i>before</i> the presentation of the stop signal. In some studies, this label is also used for go responses executed before the presentation of the go stimulus on stop trials (see description premature responses on go trials).	Premature signal-respond

*continued on next page*

*continued*

<b>Label</b>	<b>Description</b>	<b>Common alternative labels</b>
Trigger failures on stop trials	Failures to launch the stop process or 'runner' on stop trials (see <b>Box 2</b> for further discussion).	
Reaction time (RT) on go trials	How long does it take to respond to the stimulus on go trials? This corresponds to the finishing time of the go runner in the independent race model.	Go RT, go latency, no-signal RT
Stop-signal delay (SSD)	The delay between the presentation of the go stimulus and the stop signal	Stimulus-onset asynchrony (SOA)
Stop-signal reaction time (SSRT)	How long does it take to stop a response? SSD + SSRT correspond to the finishing time of the stop runner in the independent race model.	Stop latency
RT on unsuccessful stop trials	Reaction time of the go response on unsuccessful stop trials	Signal-respond RT, SR-RT (note that this abbreviation is highly similar to the abbreviation for stop-signal reaction time, which can cause confusion)

Note: The different types of unsuccessful stop trials (including choice errors and premature responses) are usually collapsed when calculating  $p(\text{respond}|\text{signal})$ , estimating SSRT, or tracking SSD.