

# Stability-Flexibility Dilemma in Cognitive Control: A Dynamical System Perspective

Sebastian Musslick<sup>1,\*</sup>, Anastasia Bizyaeva<sup>2</sup>, Shamay Agaron<sup>1</sup>, Naomi Leonard<sup>2</sup>, and Jonathan D. Cohen<sup>1</sup>

<sup>1</sup>Princeton Neuroscience Institute, Princeton University, Princeton, NJ 08544, USA.

<sup>2</sup>Department of Mechanical and Aerospace Engineering, Princeton University, Princeton, NJ 08544, USA.

\*Corresponding Author: [musslick@princeton.edu](mailto:musslick@princeton.edu)

## Abstract

Constraints on control-dependent processing have become a fundamental concept in general theories of cognition that explain human behavior in terms of rational adaptations to these constraints. However, theories miss a rationale for why such constraints would exist in the first place. Recent work suggests that constraints on the allocation of control facilitate flexible task switching at the expense of the stability needed to support goal-directed behavior in face of distraction. Here, we formulate this problem in a dynamical system, in which control signals are represented as attractors and in which constraints on control allocation limit the depth of these attractors. We derive formal expressions of the stability-flexibility tradeoff, showing that constraints on control allocation improve cognitive flexibility but impair cognitive stability. Finally, we provide evidence that human participants adapt higher constraints on the allocation of control as the demand for flexibility increases but that participants deviate from optimal constraints.

**Keywords:** cognitive control; task switching; stability-flexibility tradeoff; bounded rationality; capacity constraints

## Introduction

Numerous theories of cognition are grounded in the assumption that there are fundamental constraints on the allocation of cognitive control (Anderson, 2013; Kurzban, Duckworth, Kable, & Myers, 2013; Shenhav, Botvinick, & Cohen, 2013). Theories that assume such limitations have been successful in explaining how humans rationally allocate control under such constraints (Lieder, Shenhav, Musslick, & Griffiths, 2018; Musslick, Shenhav, Botvinick, & Cohen, 2015; Shenhav et al., 2013). However, they do not provide a rationale for *why* such limitations would exist in the first place.

A recent line of work attempts to explain the limitations of control allocation in terms of fundamental computational dilemmas in neural processing systems. For instance, Musslick et al. (2017) suggest that neural architectures are subject to a tradeoff between learning efficiency that is promoted through the use of shared task representations (Bengio, Courville, & Vincent, 2013; Caruana, 1997), on the one hand, and multitasking capability that is achieved through the separation of task representations, on the other hand (Allport, 1980; Musslick et al., 2016; Meyer & Kieras, 1997; Navon & Gopher, 1979; Salvucci & Taatgen, 2008; Feng, Schwemmer, Gershman, & Cohen, 2014). From this perspective, limitations in multitasking may reflect a preference of the neural system to learn tasks more quickly (Musslick et al., 2017; Sagiv, Musslick, Niv, & Cohen, 2018).

One way to circumvent limitations in concurrent multitasking is to execute multiple tasks in series, through flexible switching between tasks (Salvucci, Taatgen, & Borst, 2009; Fischer & Plessow, 2015). The serial execution of tasks, however, gives rise to another tradeoff known as the stability-flexibility dilemma: allocating more control to a task results in greater activation of its neural representation but also in greater persistence of this activity upon switching to a new task, yielding switch costs (Ueltzhöffer, Armbruster-Genç, & Fiebach, 2015; Goschke, 2000). By considering the problem in terms of the parameterization of a nonlinear dynamical system, in which control signals are represented as attractors, Musslick, Jang Jun, Shvartsman, Shenhav, and Cohen (2018) showed that constraints on control allocation can promote cognitive flexibility at the expense of cognitive stability. Their simulations suggest that higher constraints on control allocation are optimal in environments with higher demand for task switches. While the simulations provide a computational rationale for constraints on control, a formal analysis of the problem is lacking. It also remains to be tested whether humans adapt their constraints on control in response to demands for flexibility.

In this work, we analyze the model by Musslick et al. (2018) from a dynamical system perspective and derive formal definitions for cognitive stability and cognitive flexibility. We then prove that higher gains of a network's activation function (equivalent to inverse temperature, and thought to reflect the effects of neuromodulatory neurotransmitters such as dopamine and norepinephrine; Servan-Schreiber, Printz, & Cohen, 1990; Liljenström, 2003; Cools, 2015) can balance this tradeoff towards cognitive stability at the cost of cognitive flexibility. To assess whether human participants adjust their constraints on control as a function of flexibility demands, we fit the model to participants who performed a task switching experiment with different rates of switching. We specifically test the hypothesis that the behavior of participants in highly flexible environments can be best described by a lower gain, reflecting higher constraints on control allocation. Finally, we use computational simulations to investigate whether participants adapt to the stability-flexibility dilemma in a rational manner, by comparing fitted constraints on control against optimal constraints on control.

## Recurrent Neural Network Model

We analyze the stability-flexibility tradeoff in a recurrent neural network model described by Musslick et al. (2018). The model consists of a control module that simulates control configurations as activities of processing units. The pattern of activity associated with each control configuration evolves in an attractor landscape over the course of trials. Within each trial, the processing units bias an evidence accumulation process in the decision module that integrates information about the stimulus and generates a response.

### Control Module

We simulate the amount of control allocated to a task as the activity of a corresponding processing unit in a recurrent neural network. Here, we consider environments with two tasks, and therefore two processing units, indexed by  $i, j \in \{1, 2\}$ . The activity of each unit is determined by its net input

$$net_i(t) = w_{i,i}act_i(t) + w_{i,j}act_j(t) + I_i \quad (1)$$

which is a linear combination of the unit's own activity  $act_i(t)$  multiplied by the self-recurrent weight  $w_{i,i}$ , the activity  $act_j(t)$  of the other unit  $j \in 1, 2, j \neq i$ , multiplied by an inhibitory weight  $w_{i,j}$ , and an external input  $I_i$  (i.e., an "instruction") provided to the unit (see Figure 1A). The activities for both processing units evolve across trials according to

$$\frac{dact_i(t)}{dt} = -act_i(t) + \frac{1}{1 + e^{-g \cdot net_i(t)}} \quad (2)$$

where the  $g$  is the slope of a sigmoid activation function<sup>1</sup>. The sigmoid activation function constrains the activity of both units to lie between 0 and 1. The gain of the activation function  $g$  regulates the distance between the two control attractors, with lower gain leading to a lower activation of the currently relevant control unit and slightly higher activation of its competitor (see Figure 1B-C). From this perspective, lower gains impose higher constraints on the amount of control that can be allocated to a task but facilitate switches between tasks. Below, we provide a formal analysis of the stability-flexibility dilemma as a function of gain.

### Decision Module

We simulate the decision process using the drift diffusion model (DDM, Ratcliff, 1978). On each trial, the decision module integrates information along two stimulus dimensions  $S_1$  and  $S_2$  of a single stimulus to determine a response. Each dimension (e.g., color or motion of a moving dot stimulus) can take one of two values (e.g., red or blue; up or down), each of which is associated with one of two responses (e.g. pressing left or right button). Each of the two tasks requires mapping the current value of one of the two stimulus dimensions to its corresponding response, while ignoring the other dimension. Since both tasks involve the same

<sup>1</sup>The non-linear dynamical system presented in this work is formally equivalent to the discrete time model by Musslick et al. (2018) for a rate constant of 1.

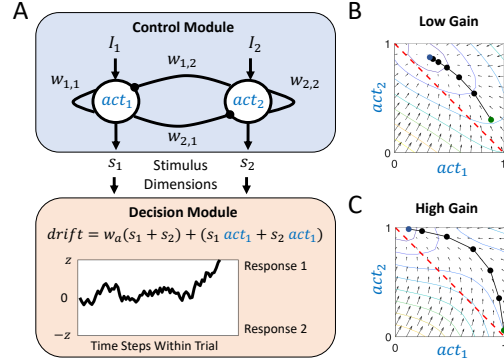


Figure 1: Model architecture. (A) The dynamics of the control module (blue) unfold over the course of trials and are determined by external input signals  $I_1, I_2$ , recurrent connectivity  $w_{1,1}, w_{2,2}$  for each unit, as well as mutual inhibition  $w_{1,2}, w_{2,1}$  between units. The activity of each control unit biases the processing of a corresponding stimulus dimension on a given trial. On each trial, the decision module accumulates evidence for both stimulus dimensions towards one of two responses until a threshold is reached. (B-C) Activation trajectories for models with a (B) low and (C) high gain are shown as a series of connected black dots, evolving from the control attractor for task 1 (green) to the control attractor for task 2 (blue). Contour lines and arrows indicate the energy and shape of the attractor landscape after a task switch from task 1 to task 2. Attractors for both tasks lie approximately on the anti-diagonal of the state space ( $act_{dif}$ ) shown in red.

pair of responses, stimuli can be congruent (stimulus values in both dimensions associated with the same response) or incongruent (associated with different responses). The drift of the DDM integration process is determined by the combined stimulus information from each dimension, weighted by input received from the control module (as described below), and evidence is accumulated over time until one of two response thresholds is reached. The drift rate is decomposed into an automatic and controlled component:

$$drift = \underbrace{w_a(S_1 + S_2)}_{\text{automatic}} + \underbrace{act_1 S_1 + act_2 S_2}_{\text{controlled}} \quad (3)$$

where the automatic component is weighted by  $w_a$  and reflects automatic processing of each stimulus dimension that is unaffected by control. The absolute magnitude of  $S_1, S_2$  depends on the strength of the association of each stimulus with a given response and its sign depends on the response (e.g.  $S_1 < 0$  if the associated response is to press the left button,  $S_1 > 0$  if the associated response is to press the right button). Thus, for congruent trials  $S_1$  and  $S_2$  have the same sign, and the opposite sign for incongruent (conflict) trials. The controlled component of the drift rate is the sum of the two stimulus values, each weighted by the activation of the corresponding control unit. Thus, each unit in the control module biases processing towards one of the stimulus dimensions. As a result, progressively greater activation of a control

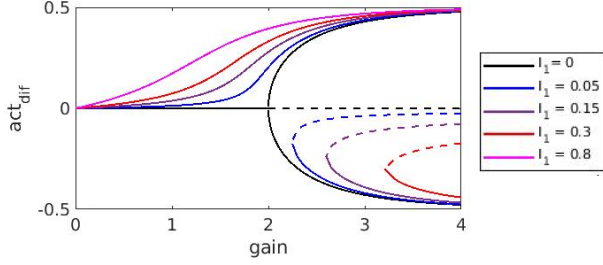


Figure 2: Difference in the activation of the two processing units ( $act_{dif}$ ) for various magnitudes of  $I_1$  in steady-state solutions to Equation (2), with  $I_2 = 0$ . Solid lines are stable attractors and dashed lines are unstable solutions. The black curve is the symmetric no-input system of Equation (6), for which with low values of gain the only attractor is the neutral state  $act_{dif} = 0$ . For values of gain greater than 2, two nonzero attractors emerge in a pitchfork bifurcation. Nonzero input to one of the processing units breaks the symmetry, splitting up the symmetric pitchfork into a continuous branch for the corresponding task and a cusp. The breakup is referred to as imperfect bifurcation (Golubitsky & Schaeffer, 1985).

unit improves performance – speeds responses and improves accuracy – for the corresponding task. Distributions of reaction times (RTs) and error rates for a given parameterization of drift rate at a given trial are derived from an analytical solution to the DDM (Navarro & Fuss, 2009).

## Formal Analysis

Previous simulation work suggests that lower values of gain facilitate switches between tasks but limit how much control can be allocated to any given task (Musslick et al., 2018). Building on work by Franci, Srivastava, and Leonard (2015); Gray, Franci, Srivastava, and Leonard (2018), we derive a formal analysis of this tradeoff as a function of gain.

For unit weights  $w_{1,2} = w_{2,1} = -1$  and  $w_{1,1} = w_{2,2} = 1$ , the attractors for both tasks are observed to lie near the antidiagonal in the activation space (see red dashed line in Figure 1B-C). We examine the dynamics of the system in a rotated frame of reference such that the attractors lie near the vertical axis. We introduce translated and rotated variables

$$\begin{pmatrix} act_{avg} \\ act_{dif} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} act_1 - 1/2 \\ act_2 - 1/2 \end{pmatrix} \quad (4)$$

where  $act_{avg}$  corresponds to the average of the two (shifted) activity states of the processing units, and  $act_{dif}$  is the average difference between the two activity states. Here,  $act_{dif}$  can be considered a proxy for cognitive stability, indexing how much control is allocated to one task versus the other. We can get an intuition for the dynamics of the system by first considering the symmetric case, in which the control module receives no input to either task processing unit  $I_1 = I_2 = 0$ .

The dynamical equations (2) with zero input decouple in

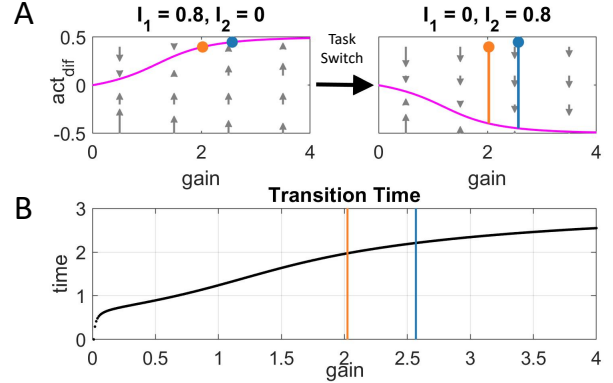


Figure 3: Relationship between  $act_{dif}$ , convergence time and gain. (A) Configuration of the system before and after a task switch. (B) Convergence time as a function of gain. Vertical lines mark examples for gain parameters that were fitted to participants' performance in environments with low (blue) and high (orange) rates of task switches.

the new variables:

$$\frac{d}{dt} act_{avg} = -act_{avg} \quad (5)$$

$$\frac{d}{dt} act_{dif} = -act_{dif} + \frac{1}{2} \tanh(g \cdot act_{dif}) \quad (6)$$

The attractors of the system are the stable steady-state solutions of (5), (6). Note that  $act_{avg}$  decays to zero and the no-input system always settles on the antidiagonal  $act_1 + act_2 = 1$ . According to the dynamics of  $act_{dif}$ , the available attractors vary with the value of the gain parameter (Figure 2).

With nonzero input, the dynamics on the diagonal and antidiagonal directions do not completely decouple. The contribution in the  $act_{avg}$  direction results in the system settling near the antidiagonal with a small offset rather than directly on the antidiagonal (Figure 1B-C). However, analogously to the symmetric no-input case, the dominant dynamical behavior is in the  $act_{dif}$  direction, shown in Figure 2. From this we can recover intuition for the tradeoff between cognitive stability and flexibility. The relationship between network gain  $g$  on the relevant domain ( $0 < act_{dif} < 0.5$  when  $I_1 \neq 0$  or  $-0.5 < act_{dif} < 0$  when  $I_2 \neq 0$ ) and  $act_{dif}$  is defined by

$$g = f(act_{dif}) = \frac{\tanh^{-1}(2act_{dif})}{act_{dif}} + E(|act_{dif}|, I) \quad (7)$$

where the first term is the explicit solution for steady state gain of the no-input system (6) and  $E(act_{dif})$  is the deviation from the symmetric case, which is a monotonically decaying function of the magnitude of  $act_{dif}$ . We can approximate this deviation with a decaying exponential fit

$$E(|act_{dif}|, I) \approx 1.4e^{-5I^{-1.1} \cdot |act_{dif}|} + 0.6 \quad (8)$$

where  $I$  is the magnitude of input. Since (7) is locally invertible on the given domain, we can express cognitive stability as a function of the network gain,  $act_{dif}(g) = f^{-1}(act_{dif})$ .

Further, we can express cognitive flexibility in terms of the time it takes to switch from one task to another, that is, the time it takes for  $act_{dif}$  to pass through zero and switch sign. From simulation we observe that the transition time is a monotonically increasing function of the network gain (see Figure 3). For an input  $I_j = 0.8$  we approximate the transition time with a linear fit  $T(g) \approx 0.8g + 0.6$ . Substituting (7) with  $I = 0.8$  for  $g$ , we obtain an expression for the stability-flexibility tradeoff

$$T(act_{dif}) \approx 0.8 \frac{\tanh^{-1}(2act_{dif})}{act_{dif}} + 1.1e^{-6.4|act_{dif}|} + 1.1 \quad (9)$$

by relating convergence time and  $act_{dif}$ . This analysis supports intuitions from prior computational work, showing that a higher network gain promotes cognitive stability at the expense of cognitive flexibility (Musslick et al., 2018). Moreover, the formal results described in this section offer a quantitative interpretation of network gain in terms of both  $act_{dif}$ , as well as  $T$ , when fitting the model to human behavior.

## Experiment

Our analysis results suggest that a system should adapt higher constraints on control (lower gains) if the demand for cognitive flexibility increases. To examine whether human participants rationally adapt constraints on control to the flexibility demands of their environment, we conducted a task switching experiment in which the rate of task switches was varied across participants. We then fit the network model to each participant and evaluated the fitted gain against the gain that optimizes the stability-flexibility tradeoff for each participant.

### Method

*Participants.* We recruited 67 participants from Amazon Turk. All participants signed a consent form prior to participation and received \$6 US for participation. The study was approved by the Institutional Review Board of Princeton University. We only included participants with an accuracy above 65% into our analysis, yielding a total of 31 participants in the low switch rate group and 27 participants in the high switch rate group.

*Apparatus and Stimuli.* Stimuli consisted of a web-based random-dot kinematogram (RDK) that we adapted from Rajananda, Lau, and Odegaard (2018). The RDK contained blue and red moving dots, some of which consistently moved in either an upward or a downward direction, and some of which moved in a random direction.

*Task and Procedure.* Participants switched between a color task, in which they had to indicate the color of the majority of the presented dots (red or blue), using the response buttons ‘A’ and ‘L’, respectively, and a motion task in which they had to indicate the direction of coherent motion (up or down), using the same response buttons ‘A’ and ‘L’, respectively. Participants performed each task over a mini-block of four to six trials. Each mini-block was preceded by a task cue (one of two cues for each task to control for cue repetition effects) that

instructed participants which tasks to perform. In some mini-blocks, participants had to repeat the task that they performed in the previous mini-block (task repetition), whereas in other mini-blocks, they had to switch to the other task (task switch). The cue was displayed for 700ms and disappeared for another 600ms. On each trial of a miniblock, the RDK stimulus was shown for 1500ms, followed by an inter-trial interval of 700ms. Participants were asked to indicate the task-relevant response while the stimulus was on the screen. In the beginning of the experiment, we used a staircasing procedure to identify coherence levels (i.e. the percent of dots having the same motion or color) for each participant that standardized performance at around 85% accuracy for both tasks. After training participants to associate the task cues with each task, participants switched between tasks over a sequence of two larger blocks of 66 miniblocks each.

*Design.* Participants were divided into two experimental groups, one that switched tasks between mini-blocks 25% of the time (low switch rate) and one that switched tasks 75% of the time (high switch rate). For each task switching sequence, we counterbalanced seven factors with respect to the first trial of each mini-block: task (color or motion task), task transition (task switch or task repetition), task cue (first or second cue associated with a task), congruency (congruent or incongruent), dot motion (upward or downward), color (mostly blue or red) and correct response (‘A’ or ‘L’ key).

*Data Analysis.* We focused our analysis on the second block of the experiment, assuming that subjects take the first block to adjust to the frequency manipulation of the experiment. We were specifically interested in the performance costs associated with task switches. Prior work suggests that switch costs diminish after the first trial of a mini-block (Rogers & Monsell, 1995). We therefore analyzed reaction times (RTs) and error rates associated with the first trial of a miniblock. Furthermore, RT data was limited to correct trials that were preceded by at least one correct trial. For each group of participants, we computed switch costs as the difference in performance between switch trials and repetition trials for both RTs and error rates. We also computed incongruency costs on task repetitions<sup>2</sup> as the difference in performance between congruent and incongruent trials. Finally, we conducted two-tailed t-tests to assess whether participants in the low switch rate group exhibited different switch costs and different incongruency costs compared to the high switch rate group.

*Model Fitting Procedure.* Before fitting parameters of the model to behavior of human participants, we evaluated how well we can recover these parameters from simulated behavior generated by the model. Motivated by the formal analysis described above, we parameterized the control module with balanced recurrent and inhibitory weights,  $w_{i,i} = 1, w_{i,j} = -1$ ,

<sup>2</sup>Incongruency costs have been shown to interact with task transition (Rogers & Monsell, 1995; Goschke, 2000; Wendt & Kiesel, 2008). To avoid confounding effects of congruency with the frequency of task switches we conditioned incongruency costs on task repetition trials.

Table 1: Fitted model parameters with prior distributions.

Parameter	Prior Distribution	Lower Bound	Upper Bound
$g$	$\text{Gamma}(2.5, 0.75)$	0	4
$z$	$\text{Gamma}(3, 0.02)$	0.01	0.25
$c$	$\text{Gamma}(3, 0.75)$	0.015	0.25
$h$	$\text{Beta}(1.2, 1.2)$	0	1
$w_a$	$\text{Gamma}(16, 0.05)$	0.1	0.5

and computed the activities of both processing units trial-by-trial, by numerically integrating Equation (2) with step size  $h$ . We set the input for the currently relevant task unit to  $I_i = 0.8$  and the input for the task-irrelevant unit to  $I_{j \neq i} = 0$ . The stimulus dimension encoding the color feature was set to  $S_1 = 0.1$  if the majority of the dots was red and set to  $S_1 = -0.1$  if the majority of the dots was blue. Similarly, the stimulus dimension encoding the motion feature was set to  $S_2 = 0.1$  if the dots were moving upward and set to  $S_2 = -0.1$  if dots were moving downward. We fixed the non-decision time of the DDM to  $T_0 = 0.2$  and fit five free parameters with priors shown in Table 1: network gain  $g$ , DDM response threshold  $z$ , DDM noise  $c$ , integration constant  $h$  and automaticity weight  $w_a$ . The number of free parameters was determined based on prior analyses of parameter identifiability, indicating that larger or different sets of free parameters may not be reliably recovered. To assess how well the five parameters can be recovered from the simulated behavior of the model, we first sampled 10 parameter configurations uniformly from the intervals shown in Table 1. We then generated distributions of response times for each trial of the second experiment block and identified parameters that maximized the likelihood of the model’s responses given the data. The identifiability of each parameter was quantified by regressing the true parameter against the fitted parameter across all sampled parameter configurations. We used the same procedure to fit the model to each participant. Finally, we conducted a one-tailed t-test to assess whether fitted gain parameters of the participants in the low switch rate group were higher relative to fitted gain parameters in the high switch rate group.

**Optimality Analysis.** To evaluate whether participants adapt rationally to the stability-flexibility dilemma, we identified the optimal gain that maximizes accuracy across all trials in the experiment, given all other fitted parameters for a given participant<sup>3</sup>. For each participant group, we computed the difference between fitted gains and optimal gains, and performed a two-sided t-test to evaluate whether fitted gain parameters systematically deviate from their optimal gain.

## Results

We found that participants who switched tasks less frequently took more time to switch tasks,  $t(56) = 2.04, p < 0.05$ , but found no significant differences in terms of error rates,  $t(56) = 0.20, p = 0.84$ . Participants showed no significant differences in incongruency costs between the two experi-

mental groups in terms of both RTs,  $t(56) = 0.93, p = 0.35$  and error rates,  $t(56) = 1.30, p = 0.20$ .

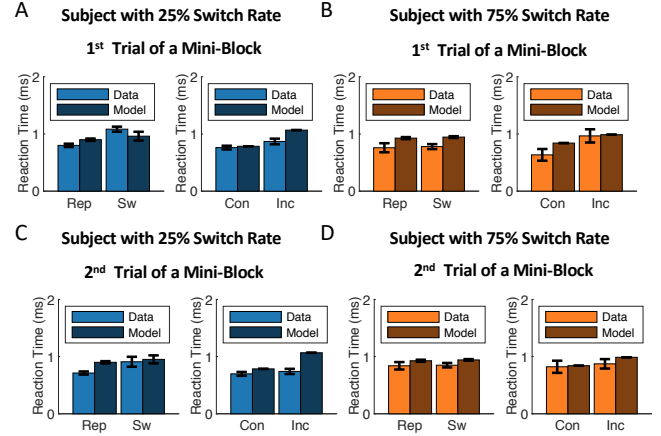


Figure 4: Examples of participant RTs and RTs generated by the fitted model. RTs are shown as a function of task transition (repetition, switch) and response congruency (congruent, incongruent) for the first (A, B) and second (C, D) trial of a mini-block. Data is shown for one participant from the low switch rate group (A, C) and one participant from the high switch rate group (B, D). Dark bars indicate average RTs generated by the fitted model. Error bars indicate the standard error of the mean across trials.

Overall, we were able to recover parameters from behavior generated by the model. The true parameter value significantly predicted the value estimated by the fitting procedure for network gain,  $b = 0.97, t(9) = 4.44, p < 0.01$ , DDM response threshold  $z$ ,  $b = 1.09, t(9) = 14.06, p < 0.001$ , DDM noise  $c$ ,  $b = 0.69, t(9) = 9.03, p < 0.001$ , integration constant  $h$ ,  $b = 0.87, t(9) = 5.02, p < 0.01$ , and automaticity weight  $w_a$ ,  $b = 0.64, t(9) = 3.06, p < 0.05$ . Figure 4 shows the behavior of two participants along with the behavior generated by the fitted model. In line with the prediction made by the model, we observed that the fitted gain parameters to behavior of human participants were significantly higher in the low switch rate group relative to the high switch rate group,  $t(56) = 3.61, p < 0.001$  (Figure 5A). Note that Figure 3 depicts formal expressions of cognitive stability (Figure 3A) and cognitive flexibility (Figure 3B) as a function of the average fitted gains for both groups. Interestingly, the fitted gains were significantly lower than the optimal gains, for both groups: low switch rates,  $t(30) = 7.40, p < 0.001$ , and high switch rates,  $t(26) = 4.24, p < 0.001$ , suggesting that, while participants adapt gain in the predicted way, overall they exert more constraint on control allocation (lower gain) than was predicted to be optimal.

## General Discussion and Conclusion

A fundamental characteristic of control-dependent processing are constraints on the allocation of control (Shiffrin & Schneider, 1977; Posner & Snyder, 1975). Recent work sug-

<sup>3</sup>We chose to maximize accuracy over maximizing reward rate as the duration of each trial was independent of response time. However, we obtained identical results when optimizing for reward rate.

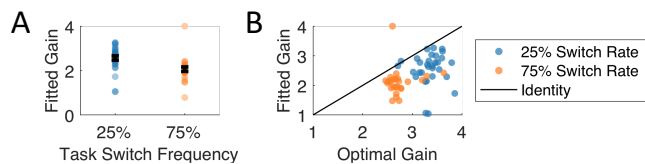


Figure 5: Model fitting results. (A) Fitted gains are shown for participants with a low (blue) and high (orange) switch rates. Each circle corresponds to the fitted gain of a participant. Vertical lines indicate the standard error of the mean fitted gain, centered around the mean. Mean gains for each group are also shown in Figure 3. (B) Fitted gain for each participant is plotted against the optimal gain that maximizes the overall accuracy of the model for a given participant.

gests that these limitations may origin from shared representations (Feng et al., 2014; Musslick et al., 2016; Salvucci & Taatgen, 2008), as well as persistence characteristics in neural systems (Musslick et al., 2018), and the resulting need to trade off the amount of control that can be allocated to a single task against the time required to switch from one task to another. In this work, we introduced a formal analysis of the latter — that is, the tradeoff between cognitive stability and cognitive flexibility.

Applying perturbation theory to the network model described by Musslick et al. (2018), we formally defined cognitive stability in terms of the distance between attractors for competing control states, and defined cognitive flexibility in terms of the time to converge from one control attractor to the other. We showed that the two measures trade off against each other, and that the balance of this tradeoff is determined by the gain of the network’s activation function. We then examined whether human participants balance this tradeoff in a similar manner as a function of the demand for flexibility, by fitting the model to participants who were required to switch tasks at either a low or high frequency. We observed that participants who switched more frequently showed lower switch costs, suggesting that they became more cognitively flexible. Moreover, model fits showed that this could be explained by lower gain, and with it, higher constraints on control. Interestingly, fitted gains for all participants were lower compared to the gains that optimized accuracy in the face of the stability-flexibility tradeoff. This suggests that there may be other factors that limit control allocation.

Altogether, our analytic and empirical results provide a rationale for how participants should adapt to different demands for flexibility given a mechanistic model for how control is represented and allocated in a recurrent neural network model. A formal relationship between cognitive stability and cognitive flexibility may not only help interpret human behavior in terms of model fits but may also help identify neural correlates for both measures. For instance, the dynamics of steady-state visually evoked potentials (SSVEP) — used to index feature-specific attention (Müller et al., 2006) — may be characterized in terms of the evolving distance between

attractors of competing attentional states. Finally, the behavioral results replicate earlier work, showing that participants’ switch costs decrease as task switches become more frequent (Mayr, 2006; Monsell & Mizon, 2006). Furthermore, prior work suggests that participants trade off cognitive flexibility against higher incongruency costs in voluntary task switching scenarios when task switches are associated with a higher reward than task repetitions (Braem, 2017).

One interesting puzzle concerns the learning mechanisms that underlie rational adaptations to changing demands in cognitive flexibility. A computationally cheap, but inflexible approach is to learn the amount of control that should be exerted through model-free reinforcement (Lieder et al., 2018). Alternatively, humans may approximate the optimal tradeoff, by attaching a cost to the amount of control that can be allocated. From this perspective, the stability-flexibility tradeoff may provide a normative rationale for parameterizing the cost of cognitive control that is integral to recent theories of control allocation (Shenhav et al., 2013, 2017).

## References

- Allport, D. A. (1980). Attention and performance. *Cognitive psychology: New Directions, 1*, 12–153.
- Anderson, J. R. (2013). *The architecture of cognition*. Psychology Press.
- Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence, 35*(8), 1798–1828.
- Braem, S. (2017). Conditioning task switching behavior. *Cognition, 166*, 272–276.
- Caruana, R. (1997). Multitask learning. *Machine Learning, 28*(1), 41–75.
- Cools, R. (2015). The cost of dopamine for dynamic cognitive control. *Current Opinion in Behavioral Sciences, 4*, 152–159.
- Feng, S. F., Schwemmer, M., Gershman, S. J., & Cohen, J. D. (2014). Multitasking vs. multiplexing: toward a normative account of limitations in the simultaneous execution of control-demanding behaviors. *Cognitive, Affective, & Behavioral Neuroscience, 14*(1), 129–146.
- Fischer, R., & Plessow, F. (2015). Efficient multitasking: parallel versus serial processing of multiple tasks. *Frontiers in psychology, 6*, 1366.
- Franci, A., Srivastava, V., & Leonard, N. (2015). A realization theory for bio-inspired collective decision-making. *arXiv:1503.08526*.
- Golubitsky, & Schaeffer. (1985). *Singularities and groups in bifurcation theory*. Springer-Verlag New York.
- Goschke, T. (2000). Intentional reconfiguration and involuntary persistence in task set switching. *Control of Cognitive Processes: Attention and Performance XVIII, 18*, 331.
- Gray, R., Franci, A., Srivastava, V., & Leonard, N. (2018). Multiagent decision-making dynamics inspired by honey-

- bees. *IEEE Transactions on Control of Network Systems*, 5(2), 793–806.
- Kurzban, R., Duckworth, A., Kable, J. W., & Myers, J. (2013). An opportunity cost model of subjective effort and task performance. *Behav. Brain Sci*, 36(6), 661–679.
- Lieder, F., Shenhav, A., Musslick, S., & Griffiths, T. L. (2018). Rational metareasoning and the plasticity of cognitive control. *PLoS Comput. Biol.*, 14(4), e1006043.
- Liljenström, H. (2003). Neural stability and flexibility: a computational approach. *Neuropsychopharmacology*, 28(S1), S64.
- Mayr, U. (2006). What matters in the cued task-switching paradigm: Tasks or cues? *Psychonomic Bulletin & Review*, 13(5), 794–799.
- Meyer, D. E., & Kieras, D. E. (1997). A computational theory of executive cognitive processes and multiple-task performance: Part I. Basic mechanisms. *Psychol. Rev.*, 104(1), 3.
- Monsell, S., & Mizon, G. A. (2006). Can the task-cuing paradigm measure an endogenous task-set reconfiguration process? *Journal of Experimental Psychology: Human Perception and Performance*, 32(3), 493.
- Müller, M., Andersen, S., Trujillo, N., Valdes-Sosa, P., Malinowski, P., & Hillyard, S. (2006). Feature-selective attention enhances color signals in early visual areas of the human brain. *Proceedings of the National Academy of Sciences*, 103(38), 14250–14254.
- Musslick, S., Dey, B., Özcimder, K., Patwary, M. M. A., Willke, T. L., & Cohen, J. D. (2016). Controlled vs. automatic processing: A graph-theoretic approach to the analysis of serial vs. parallel processing in neural network architectures. In *Proceedings of the 38th annual conference of the Cognitive Science Society* (pp. 1547–1552). Philadelphia, PA.
- Musslick, S., Jang Jun, S., Shvartsman, M., Shenhav, A., & Cohen, J. D. (2018). Constraints associated with cognitive control and the stability-flexibility dilemma. In *Proceedings of the 40th Annual Meeting of the Cognitive Science Society* (pp. 806–811). Madison, WI.
- Musslick, S., Saxe, A., Özcimder, K., Dey, B., Henselman, G., & Cohen, J. D. (2017). Multitasking capability versus learning efficiency in neural network architectures. In *Proceedings of the 39th Annual Meeting of the Cognitive Science Society* (pp. 829–834). London, UK.
- Musslick, S., Shenhav, A., Botvinick, M. M., & Cohen, J. D. (2015). A computational model of control allocation based on the expected value of control. In *The 2nd Multidisciplinary Conference on Reinforcement Learning and Decision Making*. Edmonton, Can.
- Navarro, D. J., & Fuss, I. G. (2009). Fast and accurate calculations for first-passage times in wiener diffusion models. *Journal of Mathematical Psychology*, 53(4), 222–230.
- Navon, D., & Gopher, D. (1979). On the economy of the human-processing system. *Psychol. Rev.*, 86(3), 214.
- Posner, M., & Snyder, C. (1975). attention and cognitive control. In *Information processing and cognition: The Loyola symposium* (pp. 55–85).
- Rajananda, S., Lau, H., & Odegaard, B. (2018). A random-dot kinematogram for web-based vision research. *Journal of Open Research Software*, 6(1).
- Ratcliff, R. (1978). A theory of memory retrieval. *Psychol. Rev.*, 85(2), 59.
- Rogers, R. D., & Monsell, S. (1995). Costs of a predictable switch between simple cognitive tasks. *Journal of Experimental Psychology: General*, 124(2), 207.
- Sagiv, Y., Musslick, S., Niv, Y., & Cohen, J. D. (2018). Efficiency of learning vs. processing: Towards a normative theory of multitasking. In *Proceedings of the 40th Annual Conference of the Cognitive Science Society* (pp. 1004–1009). Madison, WI.
- Salvucci, D. D., & Taatgen, N. A. (2008). Threaded cognition: an integrated theory of concurrent multitasking. *Psychol. Rev.*, 115(1), 101.
- Salvucci, D. D., Taatgen, N. A., & Borst, J. P. (2009). Toward a unified theory of the multitasking continuum: From concurrent performance to task switching, interruption, and resumption. In *Proceedings of the sigchi conference on human factors in computing systems* (pp. 1819–1828).
- Servan-Schreiber, D., Printz, H., & Cohen, J. D. (1990). A network model of catecholamine effects: gain, signal-to-noise ratio, and behavior. *Science*, 249(4971), 892–895.
- Shenhav, A., Botvinick, M. M., & Cohen, J. D. (2013). The expected value of control: an integrative theory of anterior cingulate cortex function. *Neuron*, 79(2), 217–240.
- Shenhav, A., Musslick, S., Lieder, F., Kool, W., Griffiths, T. L., Cohen, J. D., & Botvinick, M. M. (2017). Toward a rational and mechanistic account of mental effort. *Annual Review of Neuroscience*, 99–124.
- Shiffrin, R. M., & Schneider, W. (1977). Controlled and automatic human information processing: II. Perceptual learning, automatic attending and a general theory. *Psychol. Rev.*, 84(2), 127.
- Ueltzhöffer, K., Armbruster-Genç, D. J., & Fiebach, C. J. (2015). Stochastic dynamics underlying cognitive stability and flexibility. *PLoS Comput. Biol.*, 11(6), e1004331.
- Wendt, M., & Kiesel, A. (2008). The impact of stimulus-specific practice and task instructions on response congruency effects between tasks. *Psychological Research*, 72(4), 425–432.