

# A signaling theory of self-handicapping

Yang Xiang,<sup>1,\*</sup> Samuel J. Gershman,<sup>1,2,3,†</sup> Tobias Gerstenberg<sup>4,†</sup>

<sup>1</sup>Department of Psychology, Harvard University

<sup>2</sup>Center for Brain Science, Harvard University

<sup>3</sup>Center for Brains, Minds, and Machines, MIT

<sup>4</sup>Department of Psychology, Stanford University

\*Corresponding author: [yyx@g.harvard.edu](mailto:yyx@g.harvard.edu)

†Equal senior authors

## Abstract

People use various strategies to bolster the perception of their competence. One strategy is *self-handicapping*, by which people deliberately impede their performance in order to protect or enhance perceived competence. Despite much prior research, it is unclear why, when, and how self-handicapping occurs. We develop a formal theory that chooses the optimal degree of self-handicapping based on its anticipated performance and signaling effects. We test the theory's predictions in two experiments ( $N = 400$ ), showing that self-handicapping occurs more often when it is unlikely to affect the outcome and when it increases the perceived competence in the eyes of a naive observer. With sophisticated observers (who consider whether a person chooses to self-handicap), self-handicapping is less effective when followed by failure. We show that the theory also explains the findings of several past studies. By offering a systematic explanation of self-handicapping, the theory lays the groundwork for developing effective interventions.

*Keywords:* Self-handicapping; Signaling; Competence; Bayesian inference; Attribution; Theory of mind.

One of the most important attributes of people is their competence—the ability to perform well in various aspects of life [1, 2, 3]. People use a variety of strategies to bolster others' perception of their competence [4, 5, 6, 7]. One strategy is *self-handicapping*, where a person deliberately impedes their performance to protect perceived competence in case of failure, or enhance it in case of success [8]. For example, a student might procrastinate before an exam and then use tiredness as an excuse for poor performance, rather than lack of ability.

Much work has documented self-handicapping [8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18]. In the context of academic learning, self-handicapping can be harmful. It decreases performance over time [19, 20, 21, 22, 23, 24], lowers well-being, self-esteem, academic and competence satisfaction, and intrinsic motivation [25, 26]. Therefore, it is important to better understand self-handicapping to design effective interventions.

Despite much prior research, self-handicapping remains poorly understood. Why do people self-handicap in some situations but not others [27]? Why does self-handicapping sometimes increase perceptions of competence [28] and sometimes not [29]? Past work explained the phenomenon as the self-handicapper signaling their competence by affecting the observer's causal attributions [30]. For example, failing the exam is attributed to tiredness instead of lacking competence (discounting principle), and succeeding despite having been tired leads to increased perceptions of competence (augmentation principle; e.g., 8, 31, 32, 33). However, these attributional principles do not explain when people choose to self-handicap, nor why observers form mixed impressions—some increasing their perceptions of competence and others don't. Existing theories have identified situational factors [27, 34, 35, 22] and personality traits [36, 37, 38, 11] that predict self-handicapping behaviors. However, these theories are largely informal and don't elucidate the computations that drive self-handicapping behaviors.

Here, we develop a signaling theory of self-handicapping that predicts when these attributional principles apply and how they influence the behavior. The theory explains why, when, and how self-handicapping occurs. The theory involves a naive observer, an actor, and a sophisticated observer. The naive observer evaluates the actor's competence based on their outcome and handicap. The actor seeks to impress the naive observer through strategic self-handicapping. The sophisticated observer considers the actor's decision whether to self-handicap and evaluates the actor's competence accordingly. This distinction between naive and sophisticated observers follows past work showing that observers who were previously self-handicapping actors think differently about the tactics of other actors [39].

We tested the theory in two experiments ( $N = 400$ ) that use a game show setting where actors' competence is judged by observers. Actors can choose to self-handicap. Participants played both actors and observers in different phases. We manipulated the level of observer sophistication by having participants play the role of an observer twice: once before they played the actor role, when they were "naive", and again afterward, when they were "sophisticated" and could think through how an actor might behave. Consistent with the theoretical predictions, we found that: (a) Participants were more likely to self-handicap when they were either very incompetent or very competent, but not when they were just good enough for the task; and (b) self-handicapping when the actor failed increased naive observers' evaluations, but less so for sophisticated observers. We additionally show that the theory captures several results from earlier self-handicapping studies.

## Results

### Theory

As illustrated in Figure 1, the theory involves: (a) a *naive observer* who evaluates an actor's competence; (b) an *actor* who seeks to impress the naive observer through strategic self-handicapping; and (c) a *sophisticated observer* who sees through the actor's intent and evaluates the actor's compe-

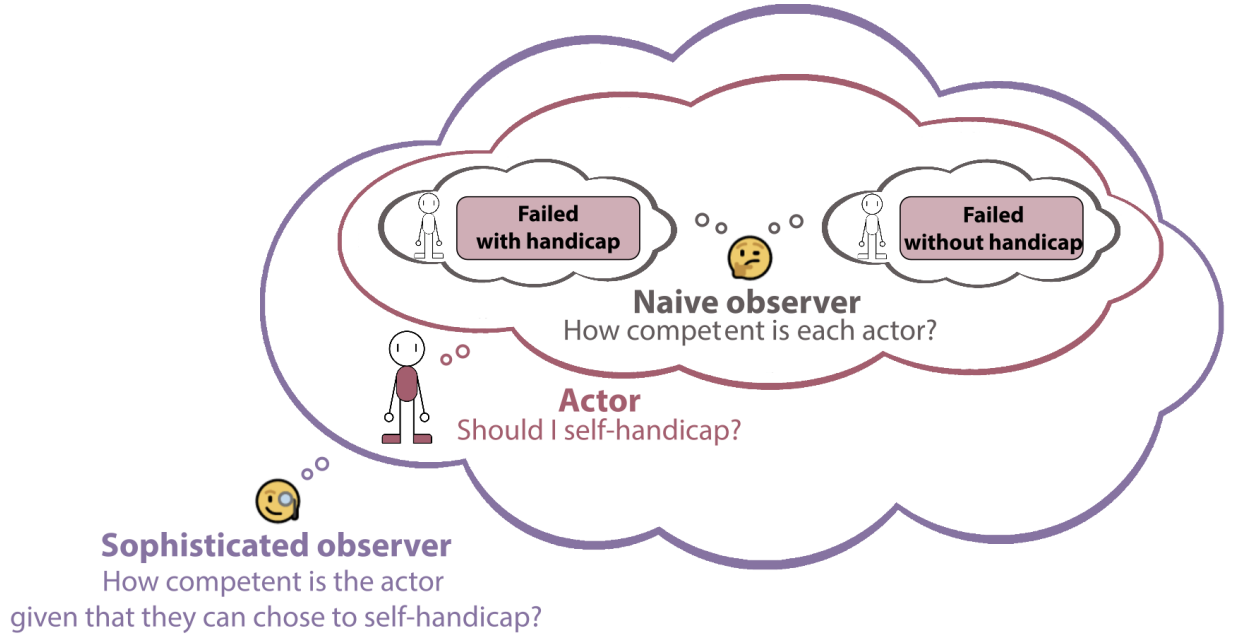


Figure 1: An illustration of the signaling theory of self-handicapping. The naive observer evaluates the actor’s competence based on whether they were handicapped ( $\gamma$ ) and whether they succeed or fail ( $s$ ). The actor decides whether to self-handicap by considering what the naive observer would infer about their competence. The sophisticated observer evaluates the actor’s competence knowing that they decided whether to self-handicap.

tence knowing that they chose whether to handicap. Suppose an actor with competence  $c$  performs a task of difficulty  $d$ , achieving outcome  $s \in \{0, 1\}$ , where  $s = 0$  indicates failure and  $s = 1$  success. The actor can choose to self-handicap, which creates an impediment such that only a proportion  $\gamma \in (0, 1]$  of their competence is applied to the task (e.g., taking a performance-inhibiting drug as in [8]).  $\gamma = 1$  means that the actor did not self-handicap at all—thus preserving full ability, whereas  $\gamma = 0$  means that the actor is completely unable to carry out a task (e.g., paralyzed).

**Naive observer**

The naive observer evaluates the actor’s competence based on the outcome  $s$  and the handicap factor  $\gamma$ , using Bayesian inference:

$$P(c|s, \gamma) \propto P(s|c, \gamma)P(c), \tag{1}$$

where  $P(c)$  is the naive observer's prior over competence, which we assume to be uniform, and  $P(s|c, \gamma)$  is the likelihood of success or failure given  $c$  and  $\gamma$ .

The likelihood of success given  $c$  and  $\gamma$  follows a logistic function:

$$P(s = 1|c, \gamma) = \frac{1}{1 + e^{-k((\gamma c - d) - b)}}, \quad (2)$$

where  $k$  controls the steepness of the curve and  $b$  adjusts the position of the sigmoid midpoint.  $\gamma$  controls what proportion of an actor's competence is used to perform the task. The probability of success thus depends on the difference between this fractional competence and the task difficulty  $d$ .

### **Actor**

We assume that the actor has two goals: (1) maximizing their perceived competence, and (2) succeeding at the task. We formalize the first goal as follows:

$$\mathbb{E}[\hat{c}|c, \gamma] = \sum_s P(s|c, \gamma) \mathbb{E}[c|s, \gamma], \quad (3)$$

where  $\hat{c} = \mathbb{E}[c|s, \gamma] = \int_c P(c|s, \gamma) c \, dc$  is the naive observer's perception of the actor's competence after observing  $s$  and  $\gamma$ . The actor maximizes the observer's perception of their competence by taking into account each possible perception  $\mathbb{E}[c|s, \gamma]$  given each potential outcome  $s$  and the selected handicap factor  $\gamma$ , weighted by the probability of that outcome occurring.

The second goal is maximizing the likelihood of success:

$$\mathbb{E}[s|c, \gamma] = \sum_s P(s|c, \gamma) s = P(s = 1|c, \gamma) \quad (4)$$

The two goals are then combined by a weight parameter  $w \in [0, 1]$  that controls the relative weight the actor places on maximizing perceived competence (the first goal) versus performance (the

second goal). We use  $Q_c(\gamma)$  to denote the value of choosing  $\gamma$  when competence is  $c$ :

$$Q_c(\gamma) = w\mathbb{E}[\hat{c}|c, \gamma] + (1 - w)\mathbb{E}[s|c, \gamma]. \quad (5)$$

We assume that the actor chooses  $\gamma$  to optimize the choice value  $Q_c(\gamma)$ . To allow for some stochasticity in choice behavior, we assume a softmax choice probability:

$$P(\gamma|c) \propto \exp[\tau Q_c(\gamma)], \quad (6)$$

where  $\tau \geq 0$  is an inverse temperature parameter that controls choice stochasticity by scaling the choice values. Smaller  $\tau$  produces more stochasticity.

### ***Sophisticated observer***

A sophisticated, “mentalizing” observer considers the actor’s decision process and recognizes that  $\gamma$  provides information about  $c$  even before  $s$  is observed. Therefore, for a sophisticated observer,  $P(c)$  in Equation 1 is replaced with  $P(c|\gamma)$ :

$$P(c|s, \gamma) \propto P(s|c, \gamma)P(c|\gamma), \quad (7)$$

where  $P(c|\gamma)$  is computed by incorporating information gained from the actor’s choice of  $\gamma$ :

$$P(c|\gamma) \propto P(\gamma|c)P(c). \quad (8)$$

Note that in order to compute  $P(\gamma|c)$ , it is necessary for the sophisticated observer to infer the actor’s choice value, which in turn requires the actor to infer the observer’s beliefs, leading to an infinite recursion. In practice, we cut off this recursion after 1 step.

## Experiments

In each experiment, 200 participants read vignettes about “Hidden Genius”, a game show where actors answered general knowledge questions and judges evaluated their competence. Actors were assigned 20 questions. They could choose to self-handicap and only be evaluated on a random subset of 10 questions. Passing required giving at least 8 correct answers, regardless of the number of questions evaluated. However, the judges did not know the exact scores; they only knew whether each actor was evaluated on 10 or 20 questions, and whether they passed (i.e., whether the 10 or 20 answers contained at least 8 correct responses).

Each experiment consisted of three blocks, illustrated in Figure 2. In the first block (Figure 2A), participants played naive judges who thought that actors *could not* choose how many of their answers were evaluated. Participants evaluated four combinations of outcomes and handicaps: An actor could pass with 20 answers evaluated (“Pass20”) or 10 (“Pass10”), or fail with 20 answers evaluated (“Fail20”) or 10 (“Fail10”). We showed participants four actors with different results on the same screen and participants answered the question “How competent is each contestant?” on a sliding scale that ranged from “Not competent at all” (coded as 0) to “Extremely competent” (coded as 100).

In the second block (Figure 2B), participants played the role of 11 actors whose average accuracy in practice tests ranged from 0% to 100% (in steps of 10%). The actors were presented in random order, and participants answered the question “How many answers should this contestant choose to be evaluated on?” on a sliding scale ranging from “Definitely 10 answers” (coded as 100% probability of self-handicapping) to “Definitely 20 answers” (coded as 0% probability of self-handicapping), with the middle being “Unsure” (coded as 50% probability of self-handicapping). The actors’ goal differed across two experiments. In Experiment 1, the actors’ goal was to maximize their competence evaluations. In Experiment 2, the actors’ goal was to maximize their chances of succeeding.

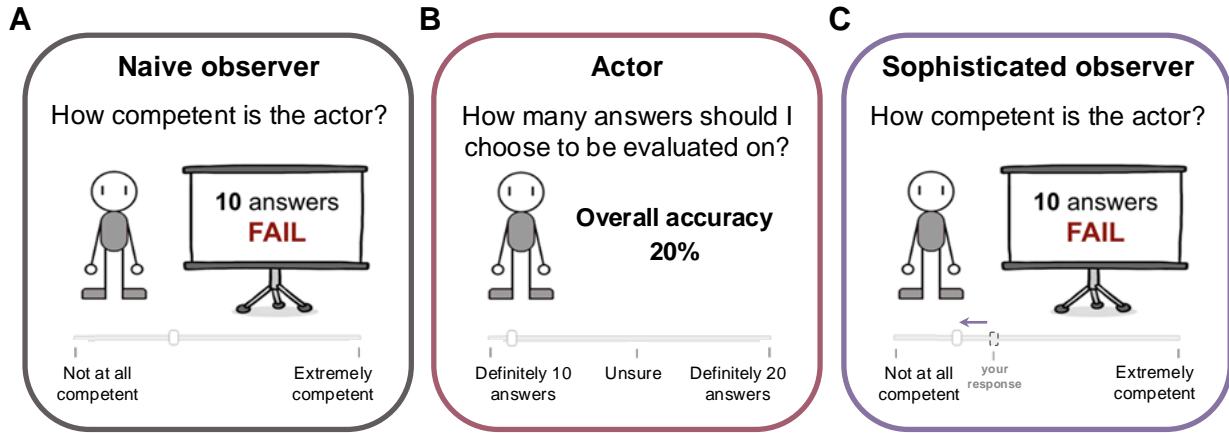


Figure 2: An illustration of the experiments. (A) In Block 1, participants played the role of naive observers and evaluated the actors’ competence based on the results. (B) In Block 2, participants played the role of actors with different competencies and decided whether to be evaluated on 10 or 20 answers. In Experiment 1, their goal was to maximize perceived competence. In Experiment 2, their goal was to maximize chances of succeeding. (C) In Block 3, participants played the role of sophisticated observers and adjusted their previous evaluations. The “your response” indicates the participant’s response in Block 1.

In the third block (Figure 2C), participants played sophisticated judges who knew that actors could choose—while the actors thought the judges didn’t know—and re-evaluated the four actors’ competence from the first block. We showed participants four actors with different results on the same screen and participants answered the question “How competent is each contestant?” on a sliding scale that ranged from “Not competent at all” (coded as 0) to “Extremely competent” (coded as 100). The sliders were initialized at the responses from the first block and participants were able to update their evaluations if they wanted to. The experimental design and analyses were preregistered (see <https://aspredicted.org/f4h3-f4xv.pdf>).

### Self-handicapping increases naive observers’ evaluations

We first analyzed data from the naive observer block (Block 1). Figure 3 shows that participants rated actors who passed the test as more competent, and the “Pass10” actor as the most competent ( $M = 83.64$ ,  $SD = 14.50$  in Experiment 1,  $M = 82.87$ ,  $SD = 15.67$  in Experiment 2), followed by the “Pass20” actor ( $M = 74.15$ ,  $SD = 15.70$  in Experiment 1,  $M = 80.80$ ,  $SD = 17.42$  in Experiment 2).



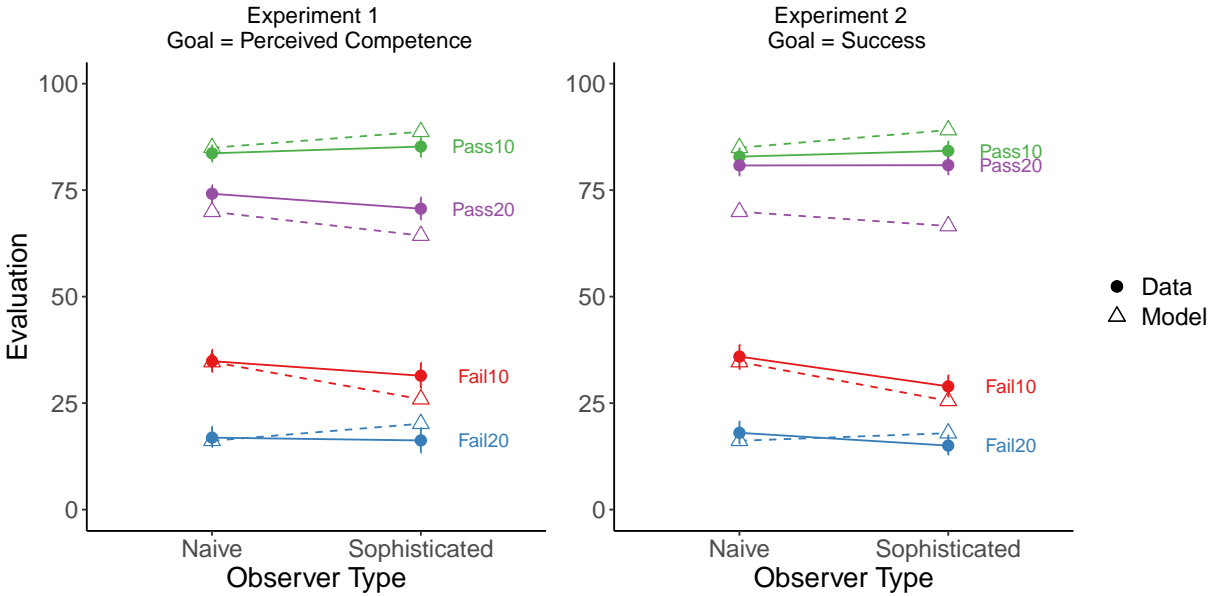


Figure 3: Observer evaluations of actors’ competence when their goal was to maximize perceived competence (Experiment 1) and when their goal was to maximize chances of succeeding (Experiment 2). In both experiments, naive observers rated actors who passed with 10 answers evaluated (“Pass10”) more competent than actors who passed with 20 answers evaluated (“Pass20”), followed by actors who failed with 10 answers evaluated (“Fail10”) and 20 (“Fail20”). Sophisticated observers rated “Fail10” actors less competent than naive observers. These patterns were captured by the model. Error bars indicate bootstrapped 95% confidence intervals.

Among actors who failed, the “Fail10” actor ( $M = 34.85$ ,  $SD = 19.20$  in Experiment 1,  $M = 35.93$ ,  $SD = 20.12$  in Experiment 2) was rated as more competent than the “Fail20” actor ( $M = 16.9$ ,  $SD = 17.29$  in Experiment 1,  $M = 18.03$ ,  $SD = 19.00$  in Experiment 2). In other words, actors are considered more competent if they pass compared to if they fail, and when the outcome is the same, actors with handicaps are perceived as more competent. The model captures this pattern (see Figure 3).

We further confirmed this finding with a Bayesian mixed-effects regression predicting participants’ competence evaluations in the naive observer block with the outcome (pass or fail), the number of answers the actor was evaluated on (10 or 20), and intercept, along with random intercept and slopes for each regressor grouped by participants. The results are summarized in Table 1. For both

experiments, we found a credible positive effect of outcome, indicating that participants rated actors who passed as more competent than actors who failed. We also found a credible negative effect of the number of answers evaluated in both experiments, meaning that actors who self-handicapped received higher evaluations than actors who did not. In summary, when the outcome is held the same, self-handicapping increases competence evaluations in the eyes of a naive observer who does not think the handicap is strategic.

Table 1: Estimates of a Bayesian mixed effects regression that was fitted for the following model:  $\text{Naive observer evaluation} \sim \text{Outcome} + \text{Answers} + (1 + \text{Outcome} + \text{Answers} \mid \text{Participant})$ . Outcome was coded as 0 = ‘fail’, and 1 = ‘pass’. Answers was coded as 0 = ‘10 answers’, and 1 = ‘20 answers’. ‘Estimate’ shows the mean of the posterior distribution, ‘Est. Error’ the standard deviation of the posterior distribution, and ‘95% CrI’ the 95% credible interval.

	Estimate	Est. Error	95% CrI
<b>Experiment 1</b>			
(Intercept)	59.25	0.88	[57.52, 60.93]
Outcome	37.51	1.03	[35.50, 39.52]
Answers	-13.72	1.22	[-16.07, -11.39]
<b>Experiment 2</b>			
(Intercept)	59.38	0.88	[57.67, 61.07]
Outcome	38.81	1.19	[36.49, 41.11]
Answers	-9.95	1.14	[-12.18, -7.69]

Note that the model predictions looked similar between the two experiments. This is because the naive observers infer the actors’ competence with the same information (outcome and handicap), without thinking about the actors’ goals. Participants, on the other hand, evaluated the “Pass20” actor as more competent when their goal was to maximize chances of succeeding (Experiment 2), compared to when their goal was to maximize the naive observers’ perceptions of their competence (Experiment 1). This pattern was not captured by the model.

### Actors self-handicap when they are very incompetent or very competent

In the actor block (Block 2), participants decided whether to self-handicap based on their com-

petence and their goal. In Experiment 1, the actors' goal was to achieve the highest competence evaluations. Based on the naive observer's inferences, actors should choose not to self-handicap when the handicap would affect the outcome—that is, if the actors would pass without self-handicapping but fail after self-handicapping. On the other hand, due to the effect of the handicap, the actors should choose to self-handicap when the handicap wouldn't affect the outcome—that is, if the actors would have at least 8 correct responses regardless of whether they were evaluated on 10 or 20 of their answers. This was indeed what we saw in the data from the actor block. The left panel of Figure 4 shows that participants were more likely to self-handicap when they were very incompetent and couldn't pass even with their full ability (accuracy between 0% and 30%) or when they were very competent and could still pass even with the handicap (accuracy between 80% and 100%), but not when they were moderately competent and self-handicapping could affect the outcome (accuracy between 40% and 70%). To test for this non-monotonic effect, we fit a Bayesian mixed-effects model predicting participants' probability of self-handicapping with the actors' quadratic accuracy, accuracy, and intercept, with random intercept and slopes for each regressor grouped by participants. As predicted, we found a credible positive effect for the quadratic term (see Table 2).

In Experiment 2, the actors' goal was to maximize their chances of succeeding. Therefore, we predicted that participants should avoid self-handicapping because the handicap would hinder their chances of passing. As predicted, participants were less likely to self-handicap for almost all competence levels except for when accuracy was 100% and the probability of self-handicapping was around chance ( $M = 50.62\%$ ,  $SD = 41.15\%$ ; Figure 4, right panel). This provided a stark contrast to Experiment 1, where participants were more likely to self-handicap, particularly when an actor's accuracy was low. Notably, the curve was still U-shaped (credible quadratic effect; see Table 2) presumably because, when the actors were very competent or very incompetent, both actions (self-handicapping or not) produced similar expected values and similar choice probability according to Equation 6, therefore it mattered less whether the actors self-handicapped. As shown

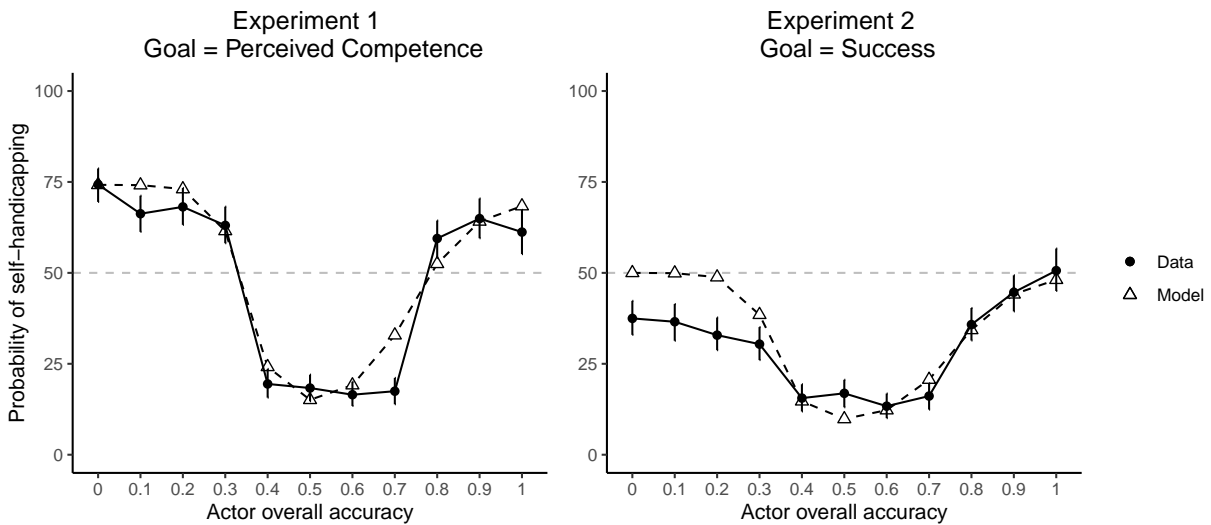


Figure 4: Actors' probability of handicapping as a function of accuracy. Passing required at least 8 correct responses, regardless of whether 10 or 20 answers were evaluated. In Experiment 1, actors were more likely to self-handicap when they were very incompetent or very competent, but less likely to do so when they were just competent enough for the task. In Experiment 2, actors overall preferred not to self-handicap, although the probability of self-handicapping was slightly higher when the actors were very incompetent or very competent. The model captures these patterns. Error bars indicate bootstrapped 95% confidence intervals.

Table 2: Estimates of a Bayesian mixed effects regression that was fitted for the following model:  $\text{Probability of self-handicapping} \sim \text{Accuracy}^2 + \text{Accuracy} + (1 + \text{Accuracy}^2 + \text{Accuracy} \mid \text{Participant})$ . ‘Estimate’ shows the mean of the posterior distribution, ‘Est. Error’ the standard deviation of the posterior distribution, and ‘95% CrI’ the 95% credible interval.

	Estimate	Est. Error	95% CrI
<b>Experiment 1</b>			
(Intercept)	86.50	2.82	[80.99, 91.97]
Accuracy <sup>2</sup>	198.15	11.07	[176.45, 219.91]
Accuracy	-215.54	10.89	[-236.79, -194.34]
<b>Experiment 2</b>			
(Intercept)	44.63	2.42	[39.96, 49.41]
Accuracy <sup>2</sup>	120.33	9.54	[102.16, 139.55]
Accuracy	-113.54	8.77	[-131.10, -96.94]

in Figure 4, the model captured the patterns in both experiments.

### **Self-handicapping is less effective with sophisticated observers when the actor fails**

How is self-handicapping perceived by observers who know that actors were able to choose whether to self-handicap? Sophisticated observers recognize that the actors’ choices provide information about their competence even before observing the outcome. By thinking about the actors’ decision, they realize that, in both Experiments 1 and 2, an actor who chooses to self-handicap is more likely to be either very incompetent or very competent, whereas an actor who chooses not to self-handicap is more likely somewhere in between the two extremes (see Figure 4). This realization should in turn affect sophisticated observers’ perception of the actors’ competence; in particular, for actors who failed, self-handicapping is more likely a strategy used to discount the failure and they should be perceived as less competent than previously thought by a naive observer.

As predicted, Figure 5 shows that sophisticated observers in both experiments provided lower evaluations of the “Fail10” actor’s competence than the naive observers ( $M = -3.42$ ,  $SD = 18.21$  in Experiment 1,  $M = -7.00$ ,  $SD = 18.22$  in Experiment 2). This effect was confirmed by a Bayesian

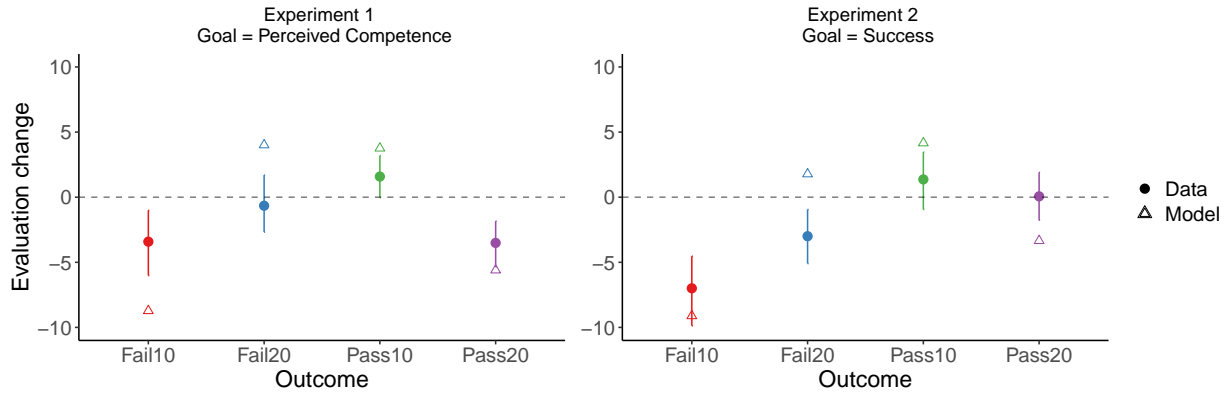


Figure 5: Changes in observer evaluations of actors’ competence (sophisticated minus naive). Notably, in both experiments, the “Fail10” actor was rated less competent by sophisticated observers. The model captures this pattern. Error bars indicate bootstrapped 95% confidence intervals.

mixed-effects model regressing evaluations of the “Fail10” actor’s competence on observer type (naive or sophisticated) and intercept, with random intercept grouped by participants (due to convergence issues, here we deviated slightly from our preregistration and did not include the random slope for observer type grouped by participants). We found a credible negative effect of observer type in both experiments (see Table 3), meaning that sophisticated observers’ evaluations of actors who self-handicapped and failed were lower than naive observers’. Since evaluation change towards the “Fail10” actor was our main interest and the only analysis we preregistered (see <https://aspredicted.org/f4h3-f4xv.pdf>), we include the regression outputs for the other three actors in the supplement.

### **The signaling theory of self-handicapping captures key patterns in past work**

Here, we demonstrate that the theory generalizes to past work, which used different setups and methods. We reviewed the 168 papers cited in [40], the latest comprehensive survey of the self-handicapping literature. We selected studies that: (a) investigated self-handicapping empirically; (b) involved actual instances or descriptions of self-handicapping rather than questionnaires that measured the tendency to self-handicap; and (c) provided the necessary data for the model to

Table 3: Estimates of a Bayesian mixed effects regression that was fitted for the following model:  $\text{Fail}_{10} \text{ actor} \sim \text{Observer type} + (1 \mid \text{Participant})$ .  $\text{Observer type}$  was coded as 0 = ‘naive’, and 1 = ‘sophisticated’. ‘Estimate’ shows the mean of the posterior distribution, ‘Est. Error’ the standard deviation of the posterior distribution, and ‘95% CrI’ the 95% credible interval.

	Estimate	Est. Error	95% CrI
<b>Experiment 1</b>			
(Intercept)	34.85	1.48	[31.99, 37.72]
Observer type	-3.41	1.30	[-6.01, -0.79]
<b>Experiment 2</b>			
(Intercept)	35.96	1.40	[33.17, 38.72]
Observer type	-7.00	1.36	[-9.67, -4.31]

generate predictions. This meant that studies where participants were actors needed to include a direct measure of competence or a related proxy, individually or by competence group, and a categorical outcome that reflected actual competence rather than fabricated feedback (e.g., non-contingent success on unsolvable tasks). Studies where participants were observers needed to report competence judgments for each combination of outcome and handicap condition, instead of separately by each factor (e.g., an average competence judgment of all actors who failed, whether or not they self-handicapped). This procedure yielded four studies: Luginbuhl and Palmer [28], Rhodewalt et al. [29], Tice [9], and Tice and Baumeister [11]. We modeled the first three, but excluded Tice and Baumeister [11] because it was an earlier version of Tice [9] that manipulated fewer variables and had less information to constrain the model parameters. The results of the remaining three studies fit with our experimental results: Luginbuhl and Palmer [28] captures how naive observers perceive actors’ competence; Tice [9] captures how actors make strategic choices; and Rhodewalt et al. [29] captures how sophisticated observers view the actors.

We emphasize the theory’s ability to qualitatively capture the patterns found in past work, predicting changes in patterns across experimental conditions, rather than evaluating its quantitative fit. Because these papers had very different setups and measurements, we had to make assumptions

in order to generate theoretical predictions, which we spell out below. Another issue we had to deal with was lack of data—the studies only presented participants’ mean responses. We therefore hand-tuned the parameters to bring the predictions quantitatively closer to the data. Note that information from the papers is too sparse to strongly constrain the parameter values, but we know that the parameter values have to be in a certain range to produce certain outcomes. For example, to achieve 75 points out of 100 as in Luginbuhl and Palmer [28],  $\gamma$  has to be relatively large; by way of illustration, it would be unlikely for even the most competent person ( $c = 100$ ) to get at least 75 points if their effective competence is only 20 ( $\gamma = 0.2$ ). However, the same qualitative patterns arise for a broad range of the possible parameter values.

### **Naive observers’ evaluations**

In two experiments, Luginbuhl and Palmer [28] showed participants videotapes of an actor the night before an exam; a friend repeatedly asked the actor to go to a movie that night, with the actor repeatedly saying no and reiterating that he had to study. In the self-handicapping condition, participants additionally watched a second clip showing that, after the friend asked once more, the actor agreed to go. The actor’s reluctance to go to the movie makes it plausible that participants did not consider the actor’s eventual decision to go to the movie as a strategic choice aimed at maximizing competence evaluations; thus this study captures how a naive observer would perceive an actor who did not choose whether to self-handicap.

As indicated in Luginbuhl and Palmer [28], we translated “receiving Grade A/C/F” to a success rule of getting at least 95/75/55 points out of 100. We used the “predicted future test score” as a proxy for competence. The aggregated data from both experiments and model predictions are shown in Figure 6A. Both participants and the model evaluated the actor who received higher grades as more competent, and evaluated the self-handicapper as slightly more competent than the non-self-handicapper when they received the same grade.

### **Actors’ strategy**



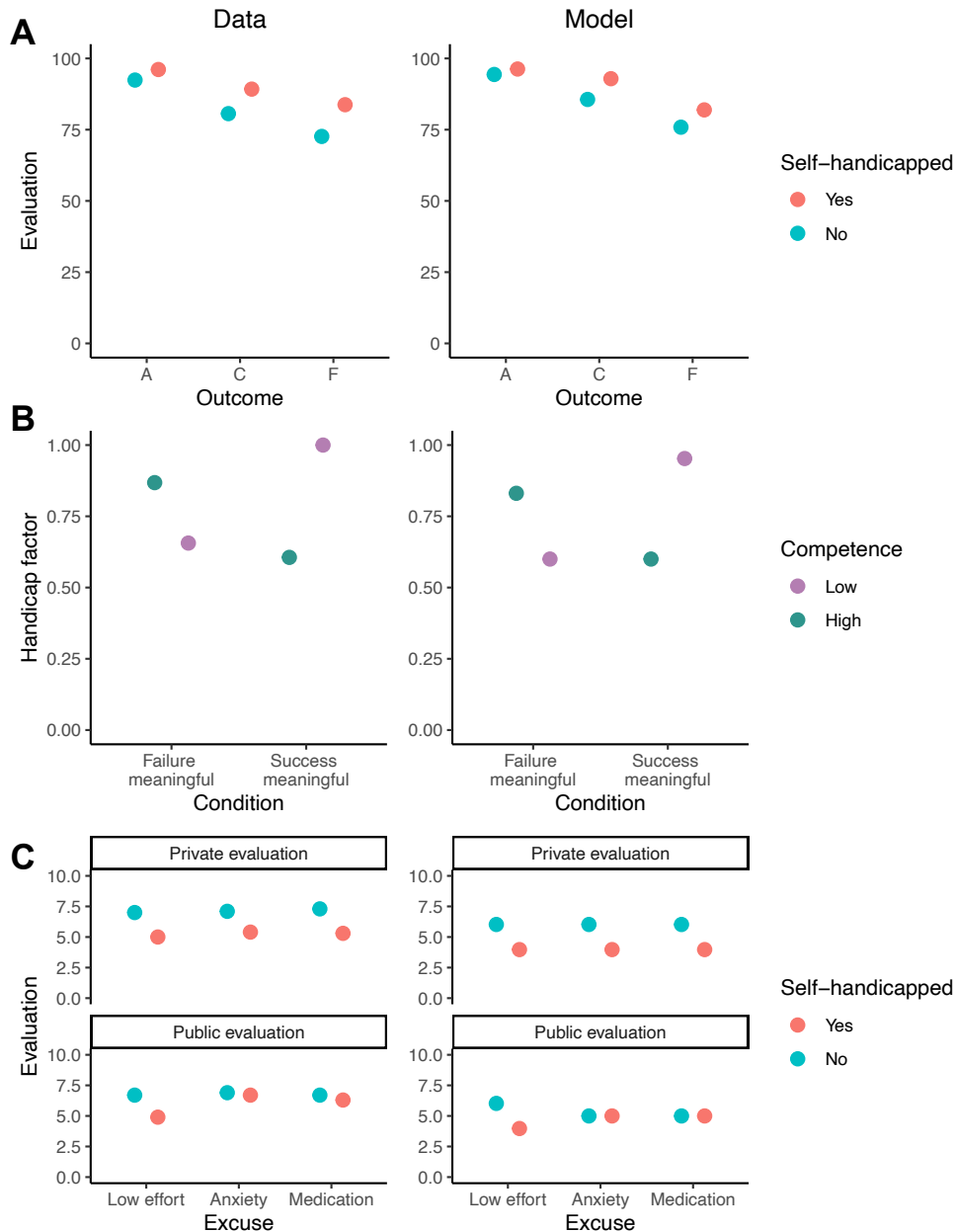


Figure 6: Modeling results of previous studies. Only participants-averaged responses were available. For illustration purposes, data were aggregated across experiments. (A) Experiment 1 and Experiment 2 in Luginbuhl and Palmer [28], where participants evaluated the competence of actors who self-handicapped or not and received a grade of either A, C, or F, corresponding to at least 95, 75, or 55 points. (B) Study 1 and Study 2 in Tice [9], where participants decided how much to self-handicap (amount of practice before important evaluation and level of distraction from a tape during a test) when failure or success was meaningful. (C) Rhodewalt et al. [29], where participants evaluated the competence of actors who self-handicapped because of low effort, anxiety, or medication.

In two experiments, Tice [9] studied how trait self-esteem (which we treat as a proxy of competence, an implicit assumption made by 9) affected how much participants self-handicapped (note that only Study 1 and Study 2 in Tice [9] were actual instances of self-handicapping; the remaining studies measured self-handicapping tendencies with questionnaires). The response variables were number of seconds practiced before an important evaluation (Study 1) and level of distraction from a tape during a test (Study 2), which we divided by the maximum response in each study to convert them to a percentage. Each study manipulated whether failure or success was meaningful. “Failure meaningful” meant that only failure would reveal information about participants’ competence, and “success meaningful” meant that only success would reveal information about participants’ competence. This was modeled by equating  $\mathbb{E}[c|s = 1, \gamma]$  with the prior  $\mathbb{E}[c]$  in the “failure meaningful” condition and  $\mathbb{E}[c|s = 0, \gamma] = \mathbb{E}[c]$  in the “success meaningful” condition. In other words, no information about competence is gained if the outcome was not meaningful.

We computed the predicted handicap factors and compared them to how much participants actually handicapped themselves. Aggregated data from the two studies in Figure 6B show that how much participants handicapped depended on which outcome was meaningful; low competence led to more handicapping to avoid failure when failure was meaningful, whereas high competence led to more handicapping to enhance success when success was meaningful. The theory captures this pattern.

### **Sophisticated observers’ evaluations**

Rhodewalt et al. [29] asked participants to evaluate the competence of actors based on their cartoon captions on a scale of  $-5$  to  $5$ . To apply the theory, we shifted the ratings to a 0-10 scale. Participants listened to a 3-min audiotape allegedly from a non-self-handicapping actor and a self-handicapping actor. The non-self-handicapping actor expressed 10 generic statements about the task (e.g., “These are interesting”), whereas the self-handicapping actor expressed 7 generic statements plus 3 condition-specific statements, either about not trying hard (low effort condition),

or performance being affected by anxiety (anxiety condition), or medication making them very sleepy (medication condition). Rhodewalt et al. [29] asked participants to listen to the tape and categorize the statements by content, indicating how many times each actor mentioned anxiety, medication, etc. Participants also knew that the actor knew that their statements were being recorded and would be listened to. These features of the setup—the participants carefully thinking about the actors’ excuses and knowing that the actor knows that the excuses would be considered—gave us confidence that participants were thinking about the possibility of actors using strategic excuses to maximize their competence perceptions; thus this study captures how a sophisticated observer would evaluate an actor’s competence.

The study additionally manipulated whether the evaluation was public. In the *private* evaluation condition, participants were told that the actors would never see their evaluations. In the *public* evaluation condition, participants were told that they would meet the actors in person and explain their evaluation to the actors. We assumed that participants more carefully considered the actors’ excuses in the public evaluation condition since the judgments mattered more.

The validity of the excuses depended on their perceived controllability. As Rhodewalt et al. [29] themselves pointed out, an actor may not be able to control their anxiety or medication, but they can control how much effort they exert. Following this assumption, the model predicted that participants were naive observers in the public evaluation condition when actors’ gave less controllable excuses (anxiety and medication) and maintained their prior belief about self-handicappers’ competence, which was the mean of a uniform distribution. By contrast, in all other situations, participants were sophisticated observers—updating their belief about the self-handicappers’ competence from their choice of whether to self-handicap as in Equation 8. Even though theoretically, knowing that someone intentionally self-handicapped by itself reveals that they could be either very incompetent or very competent, we inferred from the experiment—which labeled average competence ( $c = 5$ ) as “can’t tell (the competence)” —that participants were evaluating the actors relative to a task

difficulty of 5. The excuse statements such as “I can hardly keep my eyes open” suggested that the handicap heavily affected the actors. Since a severe handicap would make it unlikely for even the most competent actor to succeed at a task that is difficult for an average person, participants were more likely to infer that the self-handicapping actors were at the lower end of the spectrum. As shown in Figure 6C, both participants and the model gave lower evaluations to the self-handicapper in the private evaluation condition and in the public evaluation condition when the excuses were more controllable (low effort), whereas the two actors were judged similarly in the public evaluation condition when the excuses were less controllable (anxiety and medication).

## **Discussion**

We presented a signaling theory of self-handicapping that explains when people self-handicap and how that is perceived by naive and sophisticated observers. We showed that this theory generates predictions in line with behavioral data and is capable of explaining existing results in the literature. Specifically, we found that self-handicappers were perceived as more competent than non-self-handicappers for the same outcome and, relatedly, actors were more likely to self-handicap when the handicap wouldn't affect the outcome. However, sophisticated observers who reasoned about actors' decision whether to self-handicap perceived them as less competent when they failed the task, compared to naive observers. These patterns were predicted by the theory.

Previous accounts of self-handicapping have largely been informal, explaining the phenomenon as actors capitalizing on discounting and augmentation principles, without specifying exactly when these principles apply and how they determine the actors' behavior and the observers' inferences. As the first formal, normative account of self-handicapping, the signaling theory makes precise predictions about how an actor's competence affects their self-handicapping choices and how the self-handicapping behavior affects observers' perception of the actors' competence.

The theory draws a distinction between naive and sophisticated observers, which potentially rec-

onciles a debate in the literature on the effectiveness of self-handicapping. Past work has been inconclusive regarding how observers perceive self-handicappers' competence. While some studies showed that self-handicapping increases perception of competence [28], some showed the opposite [29]. Self [27] pointed out that, in order for self-handicapping to be effective, the handicapper should not appear to desire factors that hinder their performance, and that self-handicapping faces the disapproval of perceivers who detect the use of this strategy (see also 41). The signaling theory formalizes these ideas and provides an intuitive explanation for when and why self-handicapping is effective in managing perceived competence. The differences in naive and sophisticated observers' perceptions might be numerically small in our experiments, but the signaling theory is able to capture larger differences, as evidenced by the modeling results of previous studies.

The theory can be extended to deal with additional sources of uncertainty. For simplicity, we assumed that actors had perfect knowledge about their competence. While past work has shown that people can form representations of their own competence given limited amounts of data [42, 43], they might still be somewhat uncertain, and this may affect when they use these strategies (the inflection points in Figure 4, left panel). We also told observers how exactly the handicap affected the actor (i.e., the handicap factor  $\gamma$ ), which might not always be observable, and could be confounded with effort, which has been modeled in a similar way in previous work [44, 45, 46, 47]. An open question is how self-handicapping relates to the ability-effort trade-off (i.e., for a given level of success, greater attribution to ability would occur when effort is lower [48]).

We assumed that the actor only had two possible goals (or some weighted combination of the two)—maximizing competence perception (which we tested in Experiment 1) and maximizing chances of success (tested in Experiment 2). People may have other goals, such as deceiving oneself in order to maintain a positive self-image [49], or obtaining reliable diagnostic feedback about one's competence [3]. Past work has also shown that, while self-handicapping helps with competence evaluations, it can make actors seem less reliable or favorable [50, 28, 41]. How people

make these trade-offs should be studied in future work.

Finally, our findings have several implications for designing better interventions against academic self-handicapping. We showed that situational factors alone can explain why students might choose to self-handicap, without having to postulate individual differences in the tendency to self-handicap. We also showed that, even when the goal is to maximize chances of success (Experiment 2)—instead of observers' impressions of competence (Experiment 1)—self-handicapping is more likely to occur in certain situations than others. When the goal was to maximize chances of success, participants overall preferred not to self-handicap. However, it's possible that people would nonetheless choose to do so in certain situations. For example, when self-handicapping reduces effort (e.g., practicing less) or produces pleasurable feelings (e.g., using substances and drinking alcohol), students may be more willing to self-handicap. Inspired by our findings, one possible solution is for teachers to provide students with tasks that are just right for them. As shown in Figure 4, regardless of the actor's goal, self-handicapping is least likely when the task is close to the actor's competence and would thus affect the outcome. Figuring out the right task difficulty would require teachers to reason about the students' competence and assign tasks accordingly [51, 52]. Teachers may also shift the focus from vertical comparisons (i.e., comparing the performance of different students within a single period) to horizontal comparisons (i.e., comparing each student's performances over time). Past work has shown that evaluative threat induces self-handicapping [53, 54, 55, 56, 57]. This change would emphasize the goal of learning rather than appearing competent or succeeding at a task. The present work provides a theoretical and empirical basis for new interventions that target self-handicapping, helping people to realize their full potential.

## Methods

This research was approved by the Harvard Institutional Review Board. All data, code, materials, and links to experiments are publicly available at <https://github.com/yyyxiang/>

self-handicapping. The research questions, methods, and analyses were preregistered at <https://aspredicted.org/f4h3-f4xv.pdf>.

## **Participants**

We recruited 200 participants for Experiment 1 (78 Female, 120 Male, 1 Non-binary, 1 Other; mean age 44 years, range 23–77 years) and 200 participants for Experiment 2 (87 Female, 112 Male, 1 Non-binary; mean age 44 years, range 21–76 years) via Amazon’s Mechanical Turk platform (MTurk). This sample size was selected based on a power analysis on a pilot study of Experiment 1 with 29 participants, which revealed that at least 176 participants are required to detect an effect of evaluation change in the ‘Fail10’ condition with 90% power. We decided to be conservative and collect 200 participants for each experiment. This decision was preregistered. We did not exclude any participant or observation. Participants received \$4 for completing the experiment.

## **Procedure**

A complete list of task instructions and comprehension check questions is included in the supplement. Participants were first introduced to the game show “Hidden Genius”, then completed a few comprehension check questions about the game show. They subsequently completed three blocks as a naive judge, actors, and a sophisticated judge, respectively. They answered a few comprehension check questions right before each block, and were allowed to proceed only after correctly answering all of the questions. They were directed back to the relevant set of instructions each time they failed a comprehension check.

In the sophisticated observer block, we allowed participants to update their evaluations after they learned that actors could choose whether to self-handicap. We set the response slider selection button to start at their responses in the naive observer block (whereas in the naive observer block, the slider selection button was hidden until the slider was clicked). Participants updated their evaluations by dragging the selection button. They could also choose to click the selection button

if they wanted to keep the same evaluation. To remind participants of their responses in the actor block, we showed them a summary table of their responses next to the response sliders. Additionally, right before the sophisticated observer block, we asked participants three reflection questions about what number of answers a more competent, averaged-skilled, or less competent actor should choose based on their previous responses. The experimental materials are included in the Supplement.

## Model Fitting

Table 4 summarizes the parameter values for all the simulations. Below we provide a justification for each of them.

Table 4: Summary table of parameter values for all the simulations.

	$c$	$\gamma$	$d$	$k$	$b$	$\tau$	$w$
<b>New experiments</b>							
Experiment 1	[0,20]	{0.5,1}	8	1.2	-0.8	0.3	
Experiment 2	[0,20]	{0.5,1}	8	1.2	-0.8	2.5	0
<b>Past experiments</b>							
Luginbuhl and Palmer [28]	[0,100]	{0.8,1}	{95,75,55}	0.3	-3	-	-
Tice [9]	[0,10]	{0.6,1}	3	2	0	15	1
Rhodewalt et al. [29]	[0,10]	{0.5,1}	5	1	0	0.5	1

### New experiments

Task difficulty  $d = 8$ , handicap factor  $\gamma \in \{0.5, 1\}$ , and competence  $c \in [0, 20]$  (defined as the number of correct answers with the actor’s full capacity) were prescribed by the task instructions. Note that the discrete set of  $\gamma$  values was constrained by the experimental setup; this assumption is not required by the theory. We assumed a uniform prior over competence  $P(c)$  to avoid making assumptions about participants’ prior beliefs about how competent the actors were in general. The model had three free parameters, fitted to the actor block data. Two of them were the logistic growth rate  $k$  and the x-value of the sigmoid midpoint  $b$  in Equation 2, which we assumed to be



shared across both experiments. The last free parameter was the inverse temperature parameter in Equation 6 that controls the stochasticity of  $\gamma$  selection. These parameters were fit to participant-averaged data in the actor block. Data from the other two blocks were compared to the model but not used to fit the model. Because the choice values in the two experiments had different scales, we fit the inverse temperature parameter to each experiment separately. The shared logistic function parameters were fit on data pooled across two experiments and the inverse temperature on individual experiments using the Nelder-Mead optimization algorithm. Loss was measured as the total sum of squared error across two experiments.

### **Past experiments**

When possible, we used the ranges and values specified in the papers. For parameters we couldn't infer from the papers, we hand-tuned them in ranges that would make the experimental conditions possible. We assumed a uniform prior competence distribution for consistency. For Luginbuhl and Palmer [28], we set  $c \in [0, 100]$  as the points an actor could get with their full capacity, task difficulty  $d = 95, 75, \text{ and } 55$  as indicated in the paper, and  $\gamma \in \{0.8, 1\}$  so that the corresponding outcomes were possible based on the participants' evaluations. For the simulation of Tice [9], the empirical handicap factor was calculated by dividing the raw data (number of seconds practiced and level of distraction from a tape) by the maximum response in each study. Since the average was 0.8 and it is a weighted combination of the different values in the handicap factor space, we set  $\gamma \in \{0.6, 1\}$  so that it aligns with the empirical data. We also set  $c \in [0, 10]$  for simplicity. Because it was unclear how competent participants were, we had to make some assumptions about this and the task difficulty. We assumed that the ground truth was  $c = 2$  for incompetent participants and  $c = 9$  for highly competent participants. Task difficulty  $d$  was set to 3 so that it was too hard for incompetent participants and easy for highly competent participants even with the handicap.  $w$  was set to 1 according to the traditional definition of self-handicapping, meaning that the actors' only goal was to maximize evaluations. For Rhodewalt et al. [29], since the original competence range

was  $[-5, 5]$ , we increased all the ratings by 5 to ensure  $c \geq 0$ . Thus,  $c \in [0, 10]$ . The paper labeled the middle of the scale as “can’t tell (the competence)”, so we used that as the task difficulty (i.e.,  $d = 5$ ). The excuse statements such as “I can hardly keep my eyes open” suggested that the handicap heavily affected the actors, thus we set  $\gamma \in \{0.5, 1\}$  so that it was hard for self-handicappers to succeed. This study didn’t tell observers the outcome, so the dependent variable is the expectation of competence ( $\mathbb{E}[c]$  or  $\mathbb{E}[c|\gamma]$ , depending on the condition). Again,  $w$  was set to 1. For all three simulations, the two free parameters for the logistic function ( $k$  and  $b$ ) and the inverse temperature parameter were hand-tuned to bring model predictions quantitatively closer to the data.

## Data availability

All data are publicly available at <https://github.com/yyyxiang/self-handicapping>.

## Code availability

All code are publicly available at <https://github.com/yyyxiang/self-handicapping>.

## Acknowledgments

We thank Thomas Icard, Arnav Verma, Adani Abutto, and Yiqiao Wang for helpful discussions. This research was funded by National Science Foundation (DRL-2024462) to S.J.G, a McMaster Fund grant to Y.X. from the Harvard University Department of Psychology, and a Mind, Brain, Behavior Graduate Student Award to Y.X. from Harvard University.

## References

- [1] Norman H Anderson. Likableness ratings of 555 personality-trait words. *Journal of personality and social psychology*, 9(3):272, 1968. 2

- [2] Robert W White. Motivation reconsidered: the concept of competence. *Psychological review*, 66(5):297, 1959. 2
- [3] Leon Festinger. A theory of social comparison processes. *Human relations*, 7(2):117–140, 1954. 2, 21
- [4] Charles R Snyder and Raymond L Higgins. Excuses: Their effective role in the negotiation of reality. *Psychological bulletin*, 104(1):23, 1988. 2
- [5] Gifford W Bradley. Self-serving biases in the attribution process: A reexamination of the fact or fiction question. *Journal of personality and social psychology*, 36(1):56, 1978. 2
- [6] Thomas A Wills. Downward comparison principles in social psychology. *Psychological bulletin*, 90(2):245, 1981. 2
- [7] A Tesser. Toward a self-evaluation maintenance model of social behavior. *Advances in experimental social psychology/Academic Press*, 1988. 2
- [8] Steven Berglas and Edward E Jones. Drug choice as a self-handicapping strategy in response to noncontingent success. *Journal of personality and social psychology*, 36(4):405, 1978. 2, 4
- [9] Dianne M Tice. Esteem protection or enhancement? self-handicapping motives and attributions differ by trait self-esteem. *Journal of personality and social psychology*, 60(5):711, 1991. 2, 15, 17, 18, 24, 25
- [10] Frederick Rhodewalt and James Davison Jr. Self-handicapping and subsequent performance: Role of outcome valence and attributional certainty. *Basic and Applied Social Psychology*, 7(4):307–322, 1986. 2
- [11] Dianne M Tice and Roy F Baumeister. Self-esteem, self-handicapping, and self-presentation: The strategy of inadequate practice. *Journal of Personality*, 58(2):443–464, 1990. 2, 15

- [12] Ted Thompson and Anna Richardson. Self-handicapping status, claimed self-handicaps and reduced practice effort following success and failure feedback. *British Journal of Educational Psychology*, 71(1):151–170, 2001. 2
- [13] Brett L Beck, Susan R Koons, and Debra L Milgrim. Correlates and consequences of behavioral procrastination: The effects of academic procrastination, self-consciousness, self-esteem and self-handicapping. *Journal of social behavior and personality*, 15(5):3, 2000. 2
- [14] Joseph R Ferrari and Dianne M Tice. Procrastination as a self-handicap for men and women: A task-avoidance strategy in a laboratory setting. *Journal of Research in personality*, 34(1):73–83, 2000. 2
- [15] Jerald Greenberg. Unattainable goal choice as a self-handicapping strategy. *Journal of Applied Social Psychology*, 15(2):140–152, 1985. 2
- [16] Raymond L Higgins and Robert N Harris. Strategic “alcohol” use: Drinking to self-handicap. *Journal of Social and Clinical Psychology*, 6(2):191–202, 1988. 2
- [17] Jalie A Tucker, Rudy E Vuchinich, and Mark B Sobell. Alcohol consumption as a self-handicapping strategy. *Journal of Abnormal Psychology*, 90(3):220, 1981. 2
- [18] Jessi L Smith, Tiffany Hardy, and Robert Arkin. When practice doesn’t make perfect: Effort expenditure as an active behavioral self-handicapping strategy. *Journal of Research in Personality*, 43(1):95–98, 2009. 2
- [19] Shannon A Gadbois and Ryan D Sturgeon. Academic self-handicapping: Relationships with learning specific and general self-perceptions and academic performance over time. *British Journal of Educational Psychology*, 81(2):207–222, 2011. 2
- [20] Tim Urdan. Predictors of academic self-handicapping and achievement: Examining achieve-

- ment goals, classroom goal structures, and culture. *Journal of educational psychology*, 96(2): 251, 2004. 2
- [21] Miron Zuckerman, Suzanne C Kieffer, and C Raymond Knee. Consequences of self-handicapping: Effects on coping, academic performance, and adjustment. *Journal of personality and social psychology*, 74(6):1619, 1998. 2
- [22] Malte Schwinger, Linda Wirthwein, Gunnar Lemmer, and Ricarda Steinmayr. Academic self-handicapping and achievement: A meta-analysis. *Journal of educational psychology*, 106(3):744, 2014. 2
- [23] Tim Urdan, Carol Midgley, and Eric M Anderman. The role of classroom goal structure in students' use of self-handicapping strategies. *American Educational Research Journal*, 35(1):101–122, 1998. 2
- [24] Jari-Erik Nurmi, Tiina Onatsu, and Tarja Haavisto. Underachievers' cognitive and behavioural strategies-self-handicapping at school. *Contemporary Educational Psychology*, 20(2):188–200, 1995. 2
- [25] Miron Zuckerman and Fen-Fang Tsai. Costs of self-handicapping. *Journal of personality*, 73(2):411–442, 2005. 2
- [26] Sanna Eronen, Jari-Erik Nurmi, and Katariina Salmela-Aro. Optimistic, defensive-pessimistic, impulsive and self-handicapping strategies in university environments. *Learning and Instruction*, 8(2):159–177, 1998. 2
- [27] Elizabeth A Self. Situational influences on self-handicapping. In *Self-Handicapping: The Paradox That Isn't*. Plenum Press, 1990. 2, 21
- [28] James Luginbuhl and Randall Palmer. Impression management aspects of self-handicapping:

- Positive and negative effects. *Personality and Social Psychology Bulletin*, 17(6):655–662, 1991. 2, 15, 16, 17, 21, 24, 25
- [29] Frederick Rhodewalt, David M Sanbonmatsu, Brian Tschanz, David L Feick, and Ann Waller. Self-handicapping and interpersonal trade-offs: The effects of claimed self-handicaps on observers' performance evaluations and feedback. *Personality and Social Psychology Bulletin*, 21(10):1042–1050, 1995. 2, 15, 17, 18, 19, 21, 24, 25
- [30] Harold H Kelley. The processes of causal attribution. *American psychologist*, 28(2):107, 1973. 2
- [31] R. M. Arkin and A. H. Baumgardner. Attribution. basic issues and applications. In *Self-handicapping*. Orlando, FL: Academic Press, 1985. 2
- [32] R. M. Arkin and K. C. Oleson. Attribution and social interaction: The legacy of edward e. jones. In *Self-handicapping*. Washington, DC: American Psychological Association, 1998. 2
- [33] Frederick Rhodewalt and Michael W Tragakis. Self-handicapping and school: Academic self-concept and self-protective behavior. In *Improving academic achievement*, pages 109–134. Elsevier, 2002. 2
- [34] James A Shepperd and Robert M Arkin. Determinants of self-handicapping: Task importance and the effects of preexisting handicaps on self-generated handicaps. *Personality and Social Psychology Bulletin*, 15(1):101–112, 1989. 2
- [35] Andrew J Elliot, Francois Cury, James W Fryer, and Pascal Huguet. Achievement goals, self-handicapping, and performance attainment: A mediational analysis. *Journal of Sport and Exercise Psychology*, 28(3):344–361, 2006. 2
- [36] Scott R Ross, Kelli E Canada, and Marcus K Rausch. Self-handicapping and the five factor

- model of personality: Mediation between neuroticism and conscientiousness. *Personality and individual differences*, 32(7):1173–1184, 2002. 2
- [37] Jennifer L Bobo, Keila C Whitaker, and Kamden K Strunk. Personality and student self-handicapping: A cross-validated regression approach. *Personality and Individual Differences*, 55(5):619–621, 2013. 2
- [38] Harry Prapavessis and J Robert Grove. Self-handicapping and self-esteem. *Journal of applied sport Psychology*, 10(2):175–184, 1998. 2
- [39] Deborah S Smith and Michael J Strube. Self-protective tendencies as moderators of self-handicapping impressions. *Basic and Applied Social Psychology*, 12(1):63–80, 1991. 3
- [40] Lilla Török, Zsolt Péter Szabó, and László Tóth. A critical review of the literature on academic self-handicapping: Theory, manifestations, prevention and measurement. *Social Psychology of Education*, 21:1175–1202, 2018. 14
- [41] Edward R Hirt, Sean M McCrea, and Hillary I Boris. “i know you self-handicapped last exam”: Gender differences in reactions to self-handicapping. *Journal of Personality and Social Psychology*, 84(1):177, 2003. 21
- [42] John G Nicholls and Arden T Miller. Reasoning about the ability of self and others: A developmental study. *Child development*, pages 1990–1999, 1984. 21
- [43] Julia A Leonard, Julia Sandler, Amanda Nerenberg, Aidan Rubio, Laura Schulz, and Allyson Mackey. Preschoolers are sensitive to their performance over time. In *CogSci*, 2020. 21
- [44] Yang Xiang, Natalia Vélez, and Samuel J Gershman. Collaborative decision making is grounded in representations of other people’s competence and effort. *Journal of Experimental Psychology: General*, 152(6):1565, 2023. 21

- [45] Yang Xiang, Natalia Vélez, and Samuel J Gershman. Optimizing competence in the service of collaboration. *Cognitive Psychology*, 150:101653, 2024. 21
- [46] Yang Xiang, Jenna Landy, Fiery A Cushman, Natalia Vélez, and Samuel J Gershman. Actual and counterfactual effort contribute to responsibility attributions in collaborative tasks. *Cognition*, 241:105609, 2023. 21
- [47] Yang Xiang, Jenna Landy, Fiery Cushman, Natalia Vélez, and Samuel J Gershman. People reward others based on their willingness to exert effort. *Available at SSRN 4766719*, 2024. 21
- [48] Fritz Heider. *The psychology of interpersonal relations*. New York: Wiley, 1958. 21
- [49] George A Quattrone and Amos Tversky. Causal versus diagnostic contingencies: On self-deception and on the voter's illusion. *Journal of personality and social psychology*, 46(2):237, 1984. 21
- [50] Maurice J Levesque, Charles A Lowe, and Catherine Mendenhall. Self-handicapping as a method of self-presentation: An analysis of costs and benefits. *Current Research in Social Psychology*, 6(15):1–13, 2001. 21
- [51] Patrick Shafto, Noah D Goodman, and Thomas L Griffiths. A rational account of pedagogical reasoning: Teaching by, and learning from, examples. *Cognitive psychology*, 71:55–89, 2014. 22
- [52] Alicia M Chen, Andrew Palacci, Natalia Vélez, Robert D Hawkins, and Samuel J Gershman. A hierarchical Bayesian model of adaptive teaching. *Cognitive Science*, 48(7):e13477, 2024. 22
- [53] Julie Suhr and Christina Wei. Symptoms as an excuse: Attention deficit/hyperactivity disorder symptom reporting as an excuse for cognitive test performance in the context of evaluative threat. *Journal of Social and Clinical Psychology*, 32(7):753–769, 2013. 22



- [54] Jeff Stone. Battling doubt by avoiding practice: The effects of stereotype threat on self-handicapping in white athletes. *Personality and Social Psychology Bulletin*, 28(12):1667–1678, 2002. 22
- [55] CR Snyder et al. Adler’s psychology (of use) today: Personal history of traumatic life events as a self-handicapping strategy. *Journal of personality and social psychology*, 48(6):1512, 1985. 22
- [56] Edward R Hirt, Sean M McCrea, and Charles E Kimble. Public self-focus and sex differences in behavioral self-handicapping: Does increasing self-threat still make it “just a man’s game?”. *Personality and Social Psychology Bulletin*, 26(9):1131–1141, 2000. 22
- [57] Sarah Tandler, Malte Schwinger, Kristina Kaminski, and Joachim Stiensmeier-Pelster. Self-affirmation buffers claimed self-handicapping? a test of contextual and individual moderators. *Psychology*, 2014, 2014. 22