

RUNNING HEAD: WORKING MEMORY TRAINING

No evidence of intelligence improvement after working memory training: A randomized,

placebo-controlled study

Thomas S. Redick¹

Zach Shipstead¹

Tyler L. Harrison¹

Kenny L. Hicks¹

David E. Fried²

David Z. Hambrick²

Michael J. Kane³

Randall W. Engle¹

1: Georgia Institute of Technology, School of Psychology

2: Michigan State University, Department of Psychology

3: University of North Carolina Greensboro, Department of Psychology

Corresponding Author:

Thomas S. Redick

4601 Central Avenue

Columbus, IN 47203

Email: tsredick@iupuc.edu

Phone: 812-348-7236

Abstract

Numerous recent studies seem to provide evidence for the general intellectual benefits of working memory training. In reviews of the training literature, Shipstead, Redick, and Engle (2010, in press) argued that the field should treat recent results with a critical eye. Many published working memory training studies suffer from design limitations (no-contact control groups, single measures of cognitive constructs), mixed results (transfer of training gains to some tasks but not others, inconsistent transfer to the same tasks across studies), and lack of theoretical grounding (identifying the mechanisms responsible for observed transfer). The current study compared young adults who received 20 sessions of practice on an adaptive dual n -back program (working memory training group) or an adaptive visual search program (active placebo-control group) with a no-contact control group that received no practice. In addition, all subjects completed pre-test, mid-test, and post-test sessions, comprising multiple measures of fluid intelligence, multitasking, working memory capacity, crystallized intelligence, and perceptual speed. Despite improvements on both the dual n -back and visual search tasks with practice, and despite a high level of statistical power, there was no positive transfer to any of the cognitive ability tests. We discuss these results in the context of previous working memory training research, and address issues for future working memory training studies.

Keywords: training; working memory; attention; intelligence; multitasking

No evidence of intelligence improvement after working memory training: A randomized, placebo-controlled study

The idea that a brief, inexpensive intervention can improve one's cognitive abilities is appealing and supported by some research investigations. Prominent examples of successful, scientifically validated interventions include reducing stereotype threat in African-American students (Cohen, Garcia, Apfel, & Master, 2006) and treating neuropsychological impairments in psychiatric patients (Neuropsychological Educational Approach to Remediation; Medalia & Freilich, 2008). Although research on cognitive interventions is not new (e.g., Thorndike & Woodworth, 1901), the advent of inexpensive and portable computerized devices has made such programs easily accessible, as witnessed by a recent proliferation of commercial cognitive training programs (e.g., Brain Age, BrainTwister, Cogmed, JungleMemory, Lumosity, Mindsparke Brain Fitness Pro, Posit Science Brain Fitness, Posit Science InSight, WMPPro). As a representative commercial example, Lumosity's website claims: "Based on extensive research, Lumosity improves memory, attention, processing speed, and problem-solving skills so you can feel more confident in your abilities" (<http://www.lumosity.com/how-we-help>; retrieved April 26, 2012).

What evidence is available that brief cognitive training programs actually lead to transfer, or positive gains, on non-trained fluid intelligence tests? In his landmark *Human Cognitive Abilities: A Survey of Factor-Analytic Studies*, Carroll (1993) reviewed previous educational and research interventions to improve intelligence (pp. 669-674). Carroll's summary of the literature indicated very limited success in fundamentally and permanently changing one's general intellectual abilities. More recently, based on the

limited (but rapidly growing) research to date, some researchers (Klingberg, 2010; Perrig, Hollenstein, & Oelhafen, 2009; Sternberg, 2008) are optimistic about the efficacy of computerized working memory (WM) training in increasing fluid aspects of intelligence (i.e., those related to reasoning and problem-solving). In contrast, other reviews (Conway & Getz, 2010; Morrison & Chein, 2011; Shipstead et al., 2010; in press) of the recent WM training literature have concluded that many of the training programs listed above have limited efficacy in improving intelligence and reasoning abilities. Shipstead et al. (2010) did note that the adaptive dual n -back training program (used in Brain Fitness Pro, BrainTwister, and Lumosity) held promise relative to other WM training programs. The commercial uses of the adaptive dual n -back task followed a report by Jaeggi, Buschkuhl, Jonides, and Perrig (2008; hereafter *JBJP*, 2008) of an improvement in intelligence test scores in healthy, young adults, after dual n -back practice. This study has been widely cited both in the psychological literature (cited 142 times as of April 26, 2012, according to ISI Web of Science), and in the mainstream media (Highfield, 2008; Shellenbarger, 2011; Wang & Aamodt, 2009), due, in part, to the authors' conclusion that "the finding that cognitive training can improve [fluid intelligence] is a landmark result because this form of intelligence has been claimed to be largely immutable" (p. 6832, *JBJP*, 2008). Indeed, the results led Robert Sternberg (2008) to proclaim that "fluid intelligence is trainable to a significant and meaningful degree" (p. 6791). Because of the potential importance of *JBJP* (2008) and subsequent research by Jaeggi, Studer-Leuthi et al. (2010; hereafter *JSBSJP*, 2010), we will begin by critically evaluating the evidence suggesting that adaptive dual n -back practice improves intelligence.

After reviewing this work, we will present the results of a new study that sought to address limitations of previous WM training studies. Shipstead et al. (2010) noted two particular design problems prevalent in the WM training literature: (a) the use of no-contact control groups; and (b) inadequate measurement of cognitive constructs by using single tasks. First, if the only comparison is between experimental (WM training) and control (no-contact) groups, there are a number of alternative explanations that can account for any observed differences on intelligence assessments after training (Campbell & Stanley, 1963). In addition, because no task is “process-pure”, in the sense that it captures only the construct of interest without measurement error, the use of a single task to represent an ability such as intelligence leaves open the possibility that non-intelligence components of test performance have been improved via training. The purpose of the current study was to address these and other issues in a comprehensive and systematic fashion, in order to answer the following question: Does repeated practice on an adaptive dual n-back task transfer to, and actually cause, improvements in intelligence, multitasking, and WM capacity?

Adaptive Dual N-back Training

Based on numerous studies indicating a strong relationship between WM capacity and higher-order cognition (for review, see Kane, Conway, Hambrick, & Engle, 2007), the logic of many training programs is that increasing WM capacity should lead to improvements in tests measuring related constructs, including selective attention (Klingberg et al., 2005), inhibition (Thorell, Lindqvist, Bergman, Bohlin, & Klingberg, 2009), updating (Dahlin, Nyberg, Bäckman, & Neely, 2008), reading comprehension (Chein & Morrison, 2010), and fluid intelligence (JBJP, 2008). Note, however, that this

logic is applicable only if the processes that are improved via WM training are the same processes shared between WM capacity and the targeted construct. For example, latent-variable studies indicate that WM capacity and fluid intelligence share approximately 50% of their variance (Kane, Hambrick, & Conway, 2005; Oberauer, Schulze, Wilhelm, & Süß, 2005). Of course, this also means that 50% of the variance in each construct is not shared, which leaves substantial room for WM improvements that would not have any effect on fluid intelligence (and for fluid intelligence increases that do not have WM improvement as their cause).

However, there is evidence that the processes involved in successful dual n -back task performance do overlap with processes needed to solve reasoning and fluid intelligence test items. In the dual n -back task (Jaeggi et al., 2007), subjects respond to the identity of aurally presented letters and the location of visually presented squares, with letters and squares presented simultaneously. Subjects decide whether the current stimuli (letter and/or square) match the ones presented n -back, with n varying between 1 and 4 across experimental blocks for all subjects. Dual n -back accuracy correlated positively with measures of fluid intelligence test performance in three different studies (Jaeggi, Buschkuhl, Perrig, & Meier, 2010, Study 3; JSBSJP, 2010, Study 1; Redick et al., 2012). Interestingly, dual n -back correlations with fluid intelligence tests were greater than with other WM tasks. In Jaeggi, Buschkuhl et al. (2010, Studies 1 and 2), dual n -back and Reading Span correlations were not statistically different from zero. In JSBSJP (2010, Study 1) and Redick et al. (2012), dual n -back correlations with Operation Span and other complex span measures of WM were significant but smaller than the dual n -back correlations with fluid intelligence tests within the same subjects. Although it is

unclear why dual n -back accuracy is weakly related to performance on other WM measures, of most importance here is that dual n -back accuracy is positively related to performance on fluid intelligence tests.

JBJP (2008) published the first report concluding that adaptive dual n -back training improved intelligence. In the training version of the task, n changed as a function of performance across the experimental session (starting with $n = 1$). Subjects performed 20 blocks of $n + 20$ trials in each session, for approximately 30 min of daily practice. Dual n -back performance increased as a function of “dosage”, or the number of sessions completed. Critically, JBJP (2008) also reported that trained subjects exhibited significantly larger gains on an intelligence test compared to no-contact control subjects who did not perform the dual n -back between the pre- and post-test sessions (Figure 1a).

While the results presented in Figure 1a seem compelling, the figure represents data collapsed across four separate studies (Figure 1b), in which different groups receiving either 8, 12, 17, or 19 sessions of dual n -back performance were compared to four separate groups of control subjects. This is not necessarily a problem, especially if the only difference among the studies (other than the subjects) was the number of dual n -back sessions completed. However, the four studies in JBJP (2008) differed in other important ways, and these lead to numerous interpretative challenges of the combined Figure 1a:

1. *Data were collapsed across different transfer tests administered under different time limits.* To assess transfer to intelligence, the 8-session groups completed the Raven Advanced Progressive Matrices (RAPM), and the 12-, 17- and 19-session groups performed the Bochumer Matrizen-Test (BOMAT). Although both tests assess matrix

reasoning (presenting 3x3 vs. 3x5 matrices, respectively), they are not comparable in length (18 vs. 29 items, respectively). In addition, the 19-session groups were given 20 min to complete BOMAT, whereas the 12- and 17-session groups received only 10 min (S. M. Jaeggi, personal communication, May 25, 2011). As shown in Figure 2, the use of the short time limit in the 12- and 17-session studies produced substantially lower scores than the 19-session study. We argue that it is inappropriate to simply average across the number of problems solved from the four tests to create Figure 1a.

2. *Procedural differences across the four studies.* Although the dual n -back groups differed systematically in the number of practice sessions performed between pre- and post-test, other procedural changes justify keeping the four studies separate. First, the 8-session study also included an active-control group that completed simple- and choice-RT tasks (Jaeggi, 2005).¹ Second, the 17-session study also included EEG recordings during performance of the dual n -back tasks in the pre- and post-test sessions and an extra non-practice session between pre- and post-test for both groups. Finally, in addition to either RAPM or BOMAT, JBJP (2008) reported that Reading Span (no transfer) and Digit Span (positive transfer), were administered during the pre- and post-test sessions. Subjects in the 19-session study also performed additional transfer tasks during multi-day pre- and post-test sessions, exhibiting *positive* transfer to Stroop and delayed free recall tasks; *negative* transfer to digit-symbol substitution test; and *no* transfer to visuospatial span, task-switching, immediate free recall, and semantic priming tasks (Jaeggi, 2005). Note also that in the 19-session study, positive transfer was observed for Reading Span. The procedural variations mentioned here

serve as additional reasons not to collapse the intelligence transfer results across the individual studies.

3. *Patterns of transfer differed across the four studies.* As noted in JBJP (2008), the ANCOVA results, using post-test score as the dependent variable and pre-test score as the covariate, were not significant for the 8- and 12-session dual n -back training groups. In fact, the 17-session study is the only one that is visually similar to Figure 1 collapsed across the four individual studies, with matched intelligence scores between control and training groups before training and substantial differences in intelligence scores between groups after training (see Figure 2). Although JBJP (2008) interpreted the results across the four studies as consistent with a dose-dependent relationship (Figure 1b), it is also correct to state that whereas two of the studies found evidence for dual n -back transfer to matrix reasoning, two of the studies did not.

The individual studies of JBJP (2008) are also based on very small samples (e.g., $n = 7, 8,$ or 11 in each group of the four studies). In a follow-up study, JSBSJP (2010, Study 2) assigned subjects to a dual n -back ($n = 25$), visuospatial single n -back ($n = 21$), or no-contact control ($n = 43$) group. Both n -back groups performed adaptive versions of the tasks for 20 sessions. All subjects completed both RAPM (11-minute time limit) and BOMAT (16-minute time limit) in counterbalanced order, among other measures, at pre- and post-test. Summarizing the intelligence transfer results, both the single n -back and dual n -back groups showed more improvement on RAPM (Figure 3a) and BOMAT (Figure 3b) than did the control group, although the effect of dual n -back training appeared stronger in RAPM (dual n -back: $d = .98$; control: $d = .10$) than in BOMAT

(dual n -back: $d = .49$; control: $d = .29$); indeed, the dual n -back BOMAT gain did not statistically differ from the control gain (S. M. Jaeggi, September 15, 2010).

Focusing on the dual n -back versus no-contact control group comparisons across the JBJP (2008) and JSBSJP (2010) studies, only the 17-session BOMAT results in JBJP (2008), and the RAPM results in JSBSJP (2010), show clear evidence for transfer to fluid intelligence after adaptive dual n -back training. When one considers also that the comparison to a no-contact control group maximizes the likelihood of observing an effect of training that is influenced by placebo, Hawthorne, and related motivational and expectancy-based effects (French, 1953; Shipstead et al., 2010), the evidence for intelligence transfer after WM training is less compelling.

JSBSJP (2010) addressed a number of limitations of JBJP (2008), including: (a) using larger samples; (b) matching all training subjects on number of practice sessions; (c) administering both BOMAT and RAPM to all subjects as transfer measures; and (d) counterbalancing the order of BOMAT and RAPM versions across pre- and post-test sessions. However, the promising results of JSBSJP (2010) are limited by the use of a no-contact control group and only matrix-reasoning tests to measure intelligence. The present study sought to address these and other limitations in recent WM training studies.

Current Study

We followed several recommendations (Buschkuhl & Jaeggi, 2010; Shipstead et al., 2010; Sternberg, 2008) to critically examine the effectiveness of adaptive dual n -back training. As Sternberg (2008) argued, the JBJP (2008) results are so potentially important that replicating them across different laboratories and samples is necessary. The idea that relatively brief dual n -back training can increase an individual's intelligence has

important implications for applied contexts outside of the laboratory, such as educational practices and remediation for low-IQ individuals. So, one basic goal of the current study was to replicate the dual *n*-back training results by showing transfer to fluid intelligence.

We also followed Sternberg's (2008) recommendation to evaluate the efficacy of dual *n*-back training on "behaviors that extend beyond the realm of psychometric testing" (p. 6791). In previous work (Redick et al., 2012), we found that dual *n*-back accuracy was positively correlated with performance on different measures of multitasking ability. In addition, other studies have shown that WM tasks are strong predictors of multitasking performance (Bühner, König, Pick, & Krumm, 2006; Colom, Martínez-Molina, Shih, & Santacreu, 2010; Hambrick, Oswald, Darowski, Rench, & Brou, 2010). We predicted that successful dual *n*-back training could also increase multitasking performance, especially because the dual *n*-back task, itself, requires a form of multitasking.

In addition, we included a diverse sample of young adults to ensure that the results of JBJP (2008) and JSBSJP (2010) are not specific to those of above-average intelligence. As noted by Sternberg (2008), the University of Bern students in JBJP (2008) represent a selective sample. Although the majority of the current subjects were college students, we sampled across three different universities with varying academic profiles in an attempt to include subjects with a reasonably wide range of intelligence and WM abilities.

In order to draw conclusions about latent abilities, like intelligence, instead of particular tasks, like RAPM, we administered seventeen transfer measures assessing fluid intelligence, multitasking, WM capacity, crystallized intelligence, and perceptual speed, with multiple verbal and nonverbal measures of each construct. Because multiple causes

contribute to variance in performance on any single test, it's important to use numerous measures to rule out explanations based on task-specific abilities or processes. By creating factors for each of the constructs, we examined the efficacy of dual *n*-back training at the construct level. We also wanted to measure fluid intelligence with tasks other than matrix reasoning to ensure that any observed transfer was not due to the use of visuospatial materials in both the training and transfer tasks (Sternberg, 2008).

Our rationale for including multiple cognitive-ability constructs, beyond fluid intelligence, in the transfer sessions was not that we expected transfer to occur for all tasks or constructs. Rather, if fluid intelligence is actually improved via WM training, then tests showing the highest loadings on a general factor of intelligence (*g*) also should show the most transfer, and tests with the lowest *g* loadings would show the least; this is because fluid intelligence tests are more strongly associated with *g* than are other aspects of intelligence (Jensen, 1998; see Colom, Ángeles Quiroga et al., 2010, for similar logic). Therefore, after dual *n*-back training, the largest cognitive-ability improvements should be observed on the fluid intelligence tests (which have very high *g* loadings), and the smallest improvements should be observed on the perceptual speed tests (which tend to have lower *g* loadings; Marshalek, Lohman, & Snow, 1983).²

The use of an active control group is critical to elucidate the mechanisms responsible for transfer after any cognitive training (Buschkuhl & Jaeggi, 2010; Shipstead et al., 2010, in press; Sternberg, 2008). Although there are no firm guidelines about what makes for a good active-control condition, we agree with Sternberg (2008) that the task should be as adaptive and challenging as the dual *n*-back, but not thought or previously shown to depend heavily upon WM. In this way, subjects' motivations,

beliefs, expectations, and efforts would match between dual n -back and active-control groups, but their WM capacities (after training) would not. Therefore, we included an adaptive visual search training group in the current study, which we aimed to be as difficult and engaging as the dual n -back task and thus to serve as a placebo control. We chose visual search, in particular, because extensive studies with over 500 subjects have shown that individual differences in WM capacity are not related to performance in a variety of visual search tasks (Kane, Poole, Tuholski, & Engle, 2006). Because visual search performance is not likely to be determined by WM capacity, visual search training is unlikely to improve WM capacity. By comparing a visual search group that did not train WM to a dual n -back group that arguably did train WM, we can separate the potential transfer effects due to improving WM from those associated with placebo-type effects (Shipstead et al., 2010; Sternberg, 2008).

We also administered transfer sessions at three occasions (pre-, mid-, and post-test) to investigate the assertion by JBJP (2008) that the amount of dual n -back training dosage determines intelligence transfer after training (Figure 1b). Assessing this dose-response relationship within subjects, instead of between subjects as in JBJP (2008), allowed us to potentially trace the growth of improvements in cognitive abilities as a function of amount of training. Based on JBJP (2008), fluid intelligence gains for the dual n -back group during the mid-test session (after 10 training sessions) might be expected to be small or non-existent, but large by the post-test session (after 20 training sessions).

Predictions

We evaluated four possible transfer outcomes for the current study (Figure 4). The first (Figure 4a) is that the processes trained via dual n -back practice cause specific

improvements in fluid intelligence, whereas visual search practice does not. The second possibility (Figure 4b) is that visual search training produces improvement relative to the no-contact control group, but dual n -back training yields even greater improvement; this would indicate that the visual search practice produced placebo effects and/or that visual search training increased fluid intelligence somewhat. In any case, the Figure 4b results would still indicate that dual n -back training improves intelligence benefits over and above placebo effects. A third possibility (Figure 4c) is that the dual n -back and visual search training would increase fluid intelligence equivalently relative to the no-contact control groups. Equivalent transfer gains for the dual n -back and visual search training groups would indicate that the processes trained in the dual n -back are not specifically responsible for improving intelligence (leaving room for the possibility that cognitive training gains are entirely driven by placebo responses). Finally, the fourth possibility (Figure 4d) is the null hypothesis – none of the groups show differential improvement on the intelligence tests; this result would be consistent with the 8-session study of JBJP (2008), but here with 20 practice sessions.

Although the predictions above have used RAPM as the example dependent measure, they could extend to other fluid intelligence tests, and the multitasking and WM capacity measures, given that these measures have strong relationships to g (e.g., Hambrick et al., 2010). If we observed the same amount of transfer to crystallized intelligence and perceptual speed, which are predicted to have weaker relationships with g , then we could conclude that the cognitive processes trained in the adaptive dual n -back or visual search tasks are not specifically related to increases in fluid intelligence.

Method

Subjects

Subjects between 18 and 30 years old completed practice and transfer sessions at Georgia Institute of Technology or Michigan State University; the sample included students from those universities, students from Georgia State University, and a small number of non-students. In addition to the 75 subjects that completed all sessions, an additional 55 subjects completed at least the pre-test session. Thirty-six of these additional subjects began participation near the end of a semester. Most of these subjects wanted to continue participating after the semester break, but we determined that they could not continue in the study with a three-week absence during the training period. The other 19 subjects (10 dual n -back, 5 visual search, and 4 control) began the study but did not complete all sessions. Out of the 75 subjects that completed all sessions, two control subjects were excluded from data analysis because they received the same transfer test items at pre-test and post-test. Demographic information for the final sample is provided in Table 1.

Training Tasks

Adaptive dual n-back. Dual n -back training subjects performed 20 sessions of the adaptive task described previously. Subjects made button responses to visual (location of squares) and auditory (identity of letters) stimuli using their left and right index fingers, respectively. There were eight possible visual and auditory stimuli used. Simultaneous visual/auditory stimuli were presented for 500 ms, followed by a fixation screen for 2500 ms in which subjects could respond. On each trial, subjects had four response choices: (a) visual match/left key; (b) auditory match/right key; (c) visual and auditory matches/both

keys; and (d) no match/no response. Each block presented $n + 20$ trials, and subjects completed 20 blocks in each of the 20 sessions. Each block presented 4 visual targets, 4 auditory targets, 2 visual and auditory targets, and $14 + n$ nontargets.

Subjects' performance determined the level of n in the subsequent block. If the subject's visual and auditory accuracy was $\geq 90\%$ for the block, then n increased by one; if accuracy $\leq 70\%$, then n decreased by one. Any other combination of visual and auditory performance meant no change in n . The first dual n -back session provided subjects with detailed instructions and examples before initiating one block each of 2-back, 3-back, and 4-back trials as practice before proceeding to the adaptive task. For sessions 1-3, the adaptive task started at $n = 1$. For sessions 4-20, the adaptive task started at $n = 2$. The dependent variable was the mean level of n reached during each session, Consistent with JBJP (2008) and JSBSJP (2010), we used the mean n achieved during the session excluding blocks one through three, under the assumption that blocks one through three largely represented practice for most subjects.

Adaptive visual search. Visual-search training subjects performed 20 sessions of an adaptive task. On each trial, subjects reported whether a target F presented somewhere in the array was facing left or right with a left or right key-press, respectively; a leftward or rightward F always appeared amongst distractor stimuli (left- and right-facing E s, left- and right-tilted T s). After a 500 ms fixation, the visual search array appeared for 500 ms, followed by a 2500 ms mask during which subjects could respond. Each block presented 24 trials, with equal numbers of left-facing and right-facing F s. Subjects completed 20 blocks in each of the 20 sessions.

Subjects' performance determined both the number and type of distractors in the subsequent trial block. Figure 5 provides examples of the search arrays from different levels of task adaptation. Target Fs could appear anywhere in the array. On levels with homogeneous-distractor trials (e.g., level 1, level 3, etc.), the distractor identity changed across trials but was fixed within a trial (e.g., all right-facing *Es* on one trial; all left-tilted *Ts* on another trial). The levels were ordered such that an array of a given size always started with homogeneous distractors, the subsequent level was the same size array but with heterogeneous distractors (e.g., mix of right-facing *Es*, left-facing *Es*, left-tilted *Ts*, and right-tilted *Ts* as distractors on a given trial), and the next level was an increase in array by adding two more columns and rows of homogeneous distractors. If subjects' accuracy was $\geq 87.5\%$ for the block, then the level increased by one; if accuracy was $\leq 75\%$, then the level decreased by one. Any other accuracy led to no change in the level on the next block. The first visual search session provided subjects with detailed instructions and examples before presenting one block each of 2x2 homogeneous (level 1), 2x2 heterogeneous (level 2), and 4x4 homogeneous (level 3) trials as practice before proceeding to the adaptive task. For all sessions, the adaptive visual search task started at level 1. The dependent variable was the mean level reached during each session.

Transfer tasks

Computerized, alternate forms of the fluid and crystallized intelligence tests were created by taking the original items from each test and creating three test versions. The logic is similar to that used for previous studies that have divided tests into even and odd items to use for pre- and post-test (Chein & Morrison, 2010; Colom, Ángeles Quiroga et al., 2010; JBJP 2008; JSBSJP, 2010). However, because we had pre-, mid-, and post-test

sessions, we divided the fluid and crystallized intelligence tests into thirds to create three alternate versions of each measure (for similar strategy with RAPM, see Basak, Boot, Voss, & Kramer, 2008). Because the fluid and crystallized intelligence measures were all power tests, with item difficulty approximately increasing with item number, we chose a “snake” procedure to distribute the items across versions so that the tests would have equal difficulty (see Supplemental Materials, Table S1). In addition, time limits were only used for RAPM and Raven Standard Progressive Matrices – all other tests were given with no time limit. For the multitasks and perceptual speed tasks, parallel versions were created by generating unique versions from a pool of items. For the WM capacity tasks, we did not explicitly create alternate versions – because the items and order of items are generated randomly each time these programs are started, each test administration is always an alternate version.

Tests were presented in one of three different orders (A, B, C) across sessions (see Supplemental Materials, Table S2). Test orders were generated so that tests of the same construct did not occur consecutively within a session, and that each test appeared toward the beginning, middle, and end of a session across the three sessions. Test version order was counterbalanced across subjects using a Latin Square procedure (ABC, BCA, CAB).

RAPM (Fluid intelligence, spatial; Raven, Raven, & Court, 1998). Items present abstract shapes and patterns in a 3 x 3 matrix. The shape in the bottom-right location is missing, and subjects must select from the eight possible choices the item that best completes the overall pattern both vertically and horizontally. Subjects had 10 min to complete the test. The number of correct responses out of 12 is used as the dependent variable.

Raven Standard Progressive Matrices (Fluid intelligence, spatial; Raven et al., 1998). The test is similar to the advanced version, but the individual must select from either six or eight possible choices the item that best completes the matrix. Subjects had 15 min to complete the test. The number of correct responses out of 20 is used as the dependent variable.

Cattell Culture-Fair Test (Fluid intelligence, spatial; Cattell, 1973). The Cattell test is composed of four sub-tests (*Series Completion, Odd Elements, Matrix Completion, Dot Task*) of spatial reasoning tasks. The number of correct responses out of 19 is used as the dependent variable.

Paper Folding (Fluid intelligence, spatial; Ekstrom, French, Harman, & Dermen, 1976). Items present a square piece of paper on the left of the problem. The markings indicate that the paper has been folded a certain number of times, with a hole or holes then punched through the paper. Subjects decide which one of the five response choices depicts what the piece of paper would look like when completely unfolded. The number of correct responses out of 6 is used as the dependent variable.

Letter Sets (Fluid intelligence, verbal; Ekstrom et al., 1976). Each item presents five sets of four letters, and subjects induce a rule that applies to the composition and ordering of four of the five letter sets, and then indicate the set that violates the rule. The number of correct responses out of 10 is used as the dependent variable.

Number Series (Fluid intelligence, numeric; Thurstone, 1938). Each item presents a series of numbers, and subjects identify the response option that continues the sequence. The number of correct responses out of 5 is used as the dependent variable.

Inferences (Fluid intelligence, verbal; Ekstrom et al., 1976). Each item presents a one-to-three sentence passage and subjects choose the response option that is a logical necessity following only from the information provided. The number of correct responses out of 6 is used as the dependent variable.

Analogies (Fluid intelligence, verbal; Berger, Gupta, Berger, & Skinner, 1990). Each item presented an analogy in the format of *A is to B as C is to D*, with either the *C* or the *C is to D* missing. Subjects chose which of the five response options best completed the analogy. The number of correct responses out of 8 is used as the dependent variable.

SynWin (Multitasking; Elsmore, 1994). A visual display with four simultaneous sub-tasks is presented, with each sub-task in its own quadrant (Supplemental Materials, Figure S1): (a) Probe-recognition: Five letters are presented briefly at the beginning of the task. Throughout the rest of the task, a probe letter is presented every 10 s, and the subject makes a yes/no decision whether the probe was in the memory set. Points were earned for correct answers, and points were subtracted for incorrect answers. (b) Arithmetic: Subjects must add 2 three-digit numbers, and a new problem is shown after an answer is submitted. Points were earned for correct answers, and points were subtracted for incorrect answers. (c) Visual monitoring: A fuel gauge continuously drops gradually from 100 to 0 and must be reset to 0 by clicking on the gauge. Points were earned for responding before the gauge reaches 0, and points were subtracted if the gauge reached 0. (d) Auditory monitoring: High- and low-frequency tones occur every 10 s and subjects click on the quadrant when the rarely-occurring high-frequency tone is presented. Points were earned for hits, and points were subtracted for misses or false

alarms. Subjects were supposed to complete two, 5-minute blocks of the task, but because of experimenter errors, several subjects were not administered the second block of the task during each session. Therefore, SynWin performance was based on the first block in each session, which all subjects completed. The subject's score is determined by a formula that combines the points earned across all four sub-tasks, and this composite score is used as the dependent variable.

ControlTower (Multitasking; Redick et al., 2012). This multitask contains a primary comparison task with distractor tasks that interrupt primary-task performance (Supplemental Materials, Figure S2). The primary task is to search through side-by-side arrays of numbers, letters, and symbols. Certain elements of the left array are highlighted, and the appropriate items in the corresponding row of the right array must be clicked by the subject. For numbers, subjects click on the matching numbers in the right array. For letters, subjects click on the letter that precedes it alphabetically in the right array. For symbols, subjects click on the relevant symbols in the right array by referring to a consistently-mapped symbol codebook. During this primary task, several distractor tasks occur that interrupt performance. For the radar task, subjects click on the radar when a blip occurs inside of a specific area of the radar. For the airplane task, requests for landing are presented via headphones and the subject decides if one of three runways is clear for landing. For the color task, subjects press one of three buttons depending on the color that flashes. For the problem-solving task, subjects solve auditory questions by clicking the correct answer among the three response options provided. Subjects completed one, 10-minute block of the task. The subject's score on the primary task is

determined by the number of correct comparisons (numbers, letters, symbols) minus incorrect comparisons, which is used as the dependent variable.

ATClab (Multitasking; modified from Fothergill, Loft, & Neal, 2009). Each trial presents a display of four to ten planes that move dynamically along various flight paths, traveling at varying rates of speed (Supplemental Materials, Figure S3). Subjects are given a maximum of 1 min to make two to four yes/no decisions about whether two or three specific planes clustered together are in conflict given their current positions, flight path, and speed. For example in Figure S3, planes 3, 4, and 5 will be in conflict because plane 5's flight path intersects too closely with planes 3 and 4 given the speed of each plane. The proportion of correct conflict decisions (45 decisions across 15 trials) is used as the dependent variable.

Symmetry Span (WM capacity, spatial; modified from Redick et al., in press). Subjects made a vertical symmetry judgment about a black-and-white figure via mouse-click, and then were presented with a red square location within a 4x4 matrix that is to be remembered. The task is the same as the one described by Redick et al. (in press), with the exception that longer list lengths were used to try to avoid performance at ceiling at post-test. After three to six symmetry-square elements, subjects recalled the red squares in the order in which they were presented, by clicking on a blank 4x4 matrix. The total number of squares out of 54 recalled in the correct order across 12 trials (3 trials each of 3, 4, 5, and 6 elements) is used as the dependent variable.

Running Letter Span (WM capacity, verbal; Broadway & Engle, 2010). Subjects recalled the final n letters that were presented sequentially every 500 ms, where n equaled 3-7 across trials. Subjects saw $n+0$, $n+1$, and $n+2$ letters in each trial, and clicked on the

letters in a fixed response grid. The total number of letters out of 75 recalled in the correct serial order across 15 trials is used as the dependent variable.

Vocabulary (Crystallized intelligence, verbal; Zachary, 1986). Each item presented a word and one of the four response options was a synonym. The total number correct out of 13 is used as the dependent variable.

General Knowledge (Crystallized intelligence, verbal; Ekstrom et al., 1976). Items presented trivia questions about literature, world history, geography, and other topics, each with four response options. The total number correct out of 10 is used as the dependent variable.

Letter Comparison (Perceptual speed, verbal; computerized version of Salthouse & Babcock, 1991). Subjects decided whether sets of three, six, or nine consonants on either side of a line were the same or different. If the sets were the same, subjects clicked *SAME*; if sets differed, subjects clicked *DIFFERENT*. Subjects had 30 s to complete as many comparisons as possible, with the total correct across two 30-s administrations used as the dependent variable.

Number Comparison (Perceptual speed, numeric; computerized version of Salthouse & Babcock, 1991). This task was identical to Letter Comparison, but numbers were used instead of letters.

Procedure

Experimenters did not inform subjects that they were participating in a training study, nor did they give an indication that subjects should expect any aspect of performance to improve (in contrast to coaching methods used in commercial programs such as Cogmed). If subjects inquired about the study's purpose, they were told that the

researchers were investigating the effects of practice on memory and attention tasks (a generic description applicable to both training and control subjects because there were a minimum of three test sessions for all subjects). During recruitment, we informed potential subjects that they should be available multiple times over the course of four to five weeks to complete the study.

For the pre-test, mid-test, and post-test sessions, subjects were compensated \$40 per completed session; subjects completing all three transfer sessions received a 10% bonus (\$12). On average, subjects took 2 hours and 20 minutes to complete the pre-test session, and about 1 hour and 40 minutes to complete the mid- and post-test sessions.³ Pre-test sessions included collecting demographic information and longer instructions for the multitasks (including demonstration video).

Subjects in the two training groups completed an additional 20 practice sessions, each of which took between 30-40 min. Subjects could not complete more than one experimental session per day, with a maximum of seven sessions per week. Dual *n*-back and visual search subjects completed all 20 practice sessions (and the mid-test session) in an average of 46 ($SD = 13.7$) and 41 ($SD = 10.7$) days, respectively; the time to complete the training did not differ for the two groups, $t(51) = 1.59, p = .12$. Dual *n*-back and visual search group subjects performed mid-test sessions after 10 practice sessions.⁴ Control subjects performed mid- and post-test sessions at approximately the same interval of days as the training groups, and the intervals did not differ for the three groups: Pre- to Mid-test, $F(2, 70) = 1.23, p = .30$; Pre- to Post-test, $F(2, 70) = 1.58, p = .21$. Compensation for each practice session was \$10, with a 10% bonus at the end of the study for completing all practice sessions (\$20).

We assigned subjects to the dual n -back, visual search, and no-contact control groups as follows. We first assigned subjects to groups such that the pre-test performance on the three multitasks was not significantly different between groups. Our attempt to match groups suffered, however, when subjects dropped out and were replaced by another on our standby list. Moreover, to ensure that we had adequate statistical power, we tested an extra wave of subjects where we randomly assigned them to one of the two training groups without consideration of their pre-test data. Despite these complications, one-way ANOVAs on all 17 pre-test measures indicated no significant differences among the no-contact, visual search, and dual n -back groups at pre-test (all p 's > .14). Of particular interest, the three groups performed similarly at pre-test on the two WM measures: Symmetry Span, $F(2, 70) = 0.13, p = .88$, and Running Letter Span, $F(2, 70) = 0.11, p = .90$, despite no individual performing the exact same order of items and memoranda as any other subject in any group. To further assess whether the groups were different at pre-test, we additionally conducted Bonferroni-corrected post-hoc comparisons, and all between-group comparisons were non-significant (all p 's > .17). Note the inherent limitation in this pre-test comparison is that we are relying upon the failure to reject the null hypothesis, an important consideration not just for the current research but for any training study arguing for a lack of pre-test group differences, especially if the sample size is small.

All subjects completed a survey after the last task in the post-test session. Questions focused on the amount of perceived improvement in several categories, strategies used during the practice sessions, and self-reported engagement during the

experiment. Most questions used a 1-4 point rating scale, although open-ended answers were allowed for two questions, too.

Design and Analyses

We evaluated dual n -back and visual search performance via repeated-measures ANOVAs with Practice (20) as the within-subjects factor. We evaluated transfer performance on the ability tasks via factorial ANOVAs with Group (3) as the between-subjects factor and Session (3) as the within-subjects factor. Significant Group x Session interactions were decomposed using simple effects analyses focusing on the effects of Group and Session independently. Partial eta-squared (η_p^2) is reported as index of effect size. Because of the number of analyses that were conducted, an alpha of .01 was used for all transfer analyses (two-tailed); however, the transfer results are identical with an alpha of .05.

Results

Practice Effects

Looking at Figure 6, both training groups improved with practice. For the dual n -back group (Figure 6a), the practice effect was significant, $F(19, 437) = 18.77, p < .01, \eta_p^2 = .45$. We also observed substantial individual differences in dual n -back performance and dual n -back improvement across the 24 subjects: One subject reached $n = 10$, whereas another subject maxed out at only $n = 4$. As a crude index, we also examined individual differences in improvement by comparing the subjects' maximum n of session 20 to their maximum n of session 1. Twenty-two of the twenty-four subjects achieved a higher n in session 20 compared to session 1, but whereas one subject's maximum n improved by 6, the next highest improvement by any other subject was 3.

Practice also significantly improved performance for the visual search group, $F(19, 475) = 17.30, p < .01, \eta_p^2 = .41$ (Figure 6b), from an approximate mean level of 5 (6x6 homogeneous) to 7 (8x8 homogeneous). Again, we observed substantial individual differences in visual search performance and visual search improvement across the 29 subjects. For example, whereas one subject reached level 12, indicating accurate discrimination of a left- or right-facing F amongst 123 distractors from a masked array, another subject reached asymptote at level 6 (35 distractors). Twenty-one subjects achieved a higher level in session 20 compared to session 1, with seven subjects improving by 4 levels.

Transfer Data

Descriptive statistics for each of the transfer tasks are presented in Table 2. Because the results clearly indicate no transfer effects for either the dual n -back or visual search group relative to the control group, for brevity, the significance testing results for the transfer data are provided in Table 3. Out of 17 ANOVAs, there were no significant Group x Session interactions.⁵

Given our emphasis on using multiple indicators of a construct instead of single tasks, we also examined transfer performance as a function of ability z -composites for the tasks representing fluid intelligence (separate spatial and nonspatial factors), multitasking, WM capacity, crystallized intelligence, and perceptual speed. We calculated composites in a manner analogous to that reported in Jaeggi, Buschkuhl, Shah, and Jonides (2011). Specifically, for each task, two different standardized gain scores were calculated for each group for the intervals of pre- to mid-test, and pre- to post-test. The standardized gain scores were created by taking the gain for each subject

and dividing by the pre-test SD for the entire sample, collapsing across the groups. To create the composite gain scores, we averaged the standardized gain scores across the relevant constructs. The statistical results are presented in Table 4. Overall, the z-composite analyses confirm the lack of transfer observed in the individual task analyses. Note that the marginal results for the gain from pre- to mid-test for the multitasking composite represent greater improvement for the no-contact control group relative to the other two groups (see Table 2).

To facilitate comparisons with previous research, we re-analyzed all transfer data excluding the visual search group. The results matched the full analyses – no transfer for the dual *n*-back group relative to the no-contact control group. In addition, we re-analyzed all transfer data excluding the mid-test session, comparing only the pre- and post-test sessions, and, again, we found no transfer for the dual *n*-back group relative to the no-contact or active control groups. Finally, we conducted ANCOVAs for each transfer task, using the pre-test score as the covariate. The results were qualitatively the same as the ANOVA results presented in Table 3.⁶

A post-hoc power analysis (G*Power; Faul, Erdfelder, Lang, & Buchner, 2007) indicated that with our sample size, we had sufficient power to detect a significant Group (between-subjects) x Session (within-subjects) interaction, if it was present in our transfer data. Our power to detect a large ($f = .40$) or medium ($f = .25$) effect was $> .99$, based on our sample size and the use of the within-subjects correlation of $r = .53$, which was the average correlation among the repeated measures across all 17 transfer tasks (the default value in G*Power is $r = .50$). We also re-ran the power analyses using a correlation among repeated measures of $r = .30$, which was the lowest observed correlation among

pre-, mid-, and post-test performance, ignoring test version (Paper Folding). Using $r = .30$, the power to detect a large or medium effect was $> .95$. Note that we based our power analysis on a medium or large effect of dual n -back training, given the previous literature (JBJP, 2008: $d = .68$; JSBSJP, 2010: $d = .98$ for RAPM, $d = .49$ for BOMAT).

Survey Data

The dual n -back and visual search groups did not differ in their ratings for either effort, (dual n -back: $M = 3.13$, $SD = 0.63$; visual search: $M = 3.32$, $SD = 0.61$), $F(1, 49) = 1.21$, $p = .28$, $\eta_p^2 = .024$, nor enjoyableness, (dual n -back: $M = 2.09$, $SD = 0.79$; visual search: $M = 2.39$, $SD = 0.79$), $F(1, 49) = 1.90$, $p = .18$, $\eta_p^2 = .037$, for the repeated practice sessions. All subjects were asked about their perceived improvement as a function of study participation. When asked if they thought that their performance had improved by the third session, the three groups did not differ, $F(2, 69) = 1.46$, $p = .24$, $\eta_p^2 = .040$, with all but two subjects perceiving that their performance improved either “moderately so” or “very much so”. Proportions of subjects who endorsed specific improvements in each of several abilities are presented in Table 5. Chi-square tests indicated a difference across the three groups in the proportion of “Yes” responses for memory ($p = .02$) and intelligence ($p = .06$), with the dual n -back group more likely to report memory and intelligence changes than the visual search and control groups. Numerically, the visual search group had the highest rates of endorsement to changes in attention and perception, although the chi-square tests were not significant for these categories. Dual n -back subjects also endorsed changes to their everyday functioning. Ten of 23 dual n -back subjects said the study changed the way they carried out their daily activities, compared to only five of 49 combined visual search and control subjects.

When given the opportunity to elaborate on these changes, dual n -back subjects offered, for example, “My ability to multitask has improved”, “Better short term memory when doing tasks”, “I think it helps me focus better in class and while studying”, and “How to memorize orders at work”.

Finally, we asked the dual n -back and visual search groups about their strategies for the practice sessions, first as an open-ended question about what advice they would give to a friend who was just beginning the study. The most consistent response across the visual search and dual n -back groups involved getting sufficient rest before beginning the sessions. Other suggested dual n -back strategies included: (a) grouping items into sets of three; (b) visualizing the letter inside the blue square; (c) forgetting old items; and (d) giving more attention to the auditory task instead of the visual task. The dual n -back subject that reached $n = 10$ listed both the chunking and visualization strategies. Suggested visual search strategies included: (a) fixating on the central location; (b) unfocusing attention to passively encode the array; (c) moving their head further away from the computer screen; and (d) resting during the intertrial interval and then preparing for the array when the subsequent fixation point appeared. The visual search subject that reached level 12 listed the passive-encoding strategy. After the open-ended question, subjects were given five task-specific strategies and had to rate from 1 (“Almost never”) to 4 (“Almost always”) the degree to which they used the strategies listed during the training task. Analyses of the open-ended and forced-response strategy questions showed that there were no systematic relationships between the self-reported strategy used and training-task performance – subjects who reached high levels of performance reported using similar strategies as those who reached low levels of performance.

Discussion

Our study yielded three main findings. First, subjects improved with practice on both the dual n -back and visual search tasks. Second, training-group subjects showed no transfer to any of the ability measures, consistent with the prediction outlined in Figure 4d. Third, dual n -back trained subjects reported subjective improvements in various aspects of cognition in the absence of any objective evidence for change.

Of importance, we observed dual n -back improvement with practice that was consistent with the training gains shown in previous dual n -back training studies with young adults. If we had not obtained such training gains, then the null transfer effects obviously would have been uninformative. We were also successful in designing an adaptive, active-control treatment (visual search training) that yielded the same amount of experimental contact as the dual n -back task, as well as similar self-reported effort and enjoyment. Despite the performance improvements on the dual n -back and visual search tasks, no positive transfer to any of the intelligence, multitasking, WM capacity, and perceptual speed tasks was observed (although transfer was not expected for the crystallized intelligence and perceptual speed tasks). In addition, we did not find any evidence of a dose-dependent relationship between the amount of dual n -back training practice and fluid intelligence gains. That is, based on JBJP (2008), fluid intelligence gains for the dual n -back group during the mid-test session might have been expected to be small or non-existent, but large by the post-test session. However, we did not observe this pattern for any of the fluid intelligence measures.

WM training and transfer to fluid intelligence

Clearly, then, the question is: *Why didn't we observe fluid intelligence transfer for the dual n-back group?* One possible answer is that we didn't observe transfer simply because WM transfer effects to intelligence are actually not commonly observed. After completing the current study, we became aware of another adaptive dual *n*-back training study that is relevant for understanding our results. First, Seidler et al. (2010) assigned subjects to a dual *n*-back ($n = 29$) or knowledge-trainer active-control ($n = 27$) group.⁷ Subjects in the knowledge-trainer group answered multiple-choice and vocabulary questions. Although the knowledge-trainer control task was not adaptive based on the subject's performance, it provided a similar amount of contact with experimenters as did the dual *n*-back task. All subjects completed a minimum of 17 practice sessions, and multiple transfer measures during pre- and post-test sessions. Despite significant dual *n*-back practice improvements, there was no significant transfer for the dual *n*-back group versus the active-control group on BOMAT, RAPM, or verbal analogies, although transfer was observed on Operation Span.

A recent review (Morrison & Chein, 2011) of the broader WM training literature with young adult subjects detailed: (a) four studies reporting significant transfer to reasoning and intelligence measures (Klingberg, Forssberg, & Westerberg, 2002; JBJP 2008; Olesen, Westerberg, & Klingberg, 2004; Westerberg & Klingberg, 2007); (b) three published studies reporting *no* significant transfer to reasoning and intelligence measures (Chein & Morrison, 2010; Dahlin et al., 2008; Owen et al., 2010); and (c) one study reporting significant transfer to some intelligence measures but not others (Schmiedek, Lovden, & Lindenberger, 2010). Two of the significant transfer studies in the review (Klingberg et al., 2002; Olesen et al., 2004) had training group samples of $n = 4$ and 3,

respectively. The subjects in Olesen et al. (2004) were the same as those in Westerberg and Klingberg (2007; T. Klingberg, personal communication, February 14, 2010). Note that the positive intelligence transfer observed in JSBSJP (2010), and the lack of transfer observed in Seidler et al. (2010), were not included in Morrison and Chein's review. In addition, Morrison and Chein's (2011) assessment of training benefits may also have been unwittingly biased because of the file-drawer problem (Rosenthal, 1979), in which non-significant transfer results such as those described in the current research are less likely to be published.

However, a recent meta-analysis by Melby-Lervåg and Hulme (in press) indicates that even when considering published studies, few appropriately-powered empirical studies have found evidence for transfer from various WM training programs to fluid intelligence. Melby-Lervåg and Hulme reported that WM training showed evidence of transfer to verbal and spatial WM tasks ($d = .79$ and $.52$, respectively). When examining the effect of WM training on transfer to nonverbal abilities tests in 22 comparisons across 20 studies, they found an effect of $d = .19$. Critically, a moderator analysis showed that there was no effect ($d = .00$) in the 10 comparisons that used a treated control group, and there was a medium effect ($d = .38$) in the 12 comparisons that used an untreated control group.

More specifically examining the efficacy of adaptive dual n -back training in young adults, there are two results reporting transfer to RAPM and/or BOMAT when compared to a no-contact control group (JBJP, 2008; JSBSJP, 2010), and one result of no transfer to RAPM and/or BOMAT when compared to an active-control group (Seidler et al., 2010). Our data show no transfer to RAPM or other measures of fluid intelligence

when compared to either a no-contact control or active-control group. On the whole, then, our lack of significant fluid intelligence transfer results may not be that surprising.

In addition, other WM training studies have used children, older adults, or special populations (e.g., stroke patients, children with ADHD) as subjects. As with the young adults in JSBSJP (2010), certain studies of older adults have shown preliminary evidence that WM training can transfer to untrained measures of fluid intelligence. For example, Borella, Carretti, Riboldi, and De Bini (2010) trained older adults on a WM-span-like task, and compared them to a control group of older adults that completed questionnaires instead. They reported positive transfer to several tasks, including the Cattell Culture-Fair Test. However, other studies have been less optimistic. For example, three other recent studies reported no transfer to different versions of the Raven Progressive Matrices after WM-span training in older adults (Brehmer, Westerberg, & Bäckman, 2012; Richmond, Morrison, Chein, & Olson, 2011; Zinke, Zeintl, Eschen, Herzog, & Kliegel, 2012). Likewise, a recent review (Shipstead et al., in press) found that the majority of published studies using developmental and patient samples have not observed transfer to fluid intelligence after WM training. The Melby-Lervåg and Hulme (in press) meta-analysis confirmed that age was not a significant moderator of transfer to nonverbal abilities after WM training. In fact, young children ($d = .03$) and adolescent children ($d = -.05$) showed no evidence of transfer to nonverbal intelligence, inconsistent with the idea that younger children may be more susceptible to WM training and improvements in intelligence because of their increased brain plasticity relative to adults. Again, our findings do not appear to be an aberration – there is little evidence for transfer from WM training to fluid intelligence.

Limitations

We acknowledge a limitation in our data related to three of the fluid intelligence tasks. Specifically, for Number Series, Paper Folding, and Inferences, the mean pre-test scores were close to the maximum possible score, and this was likely due to our use of shortened versions of these tests which left five or six items per test version. Given this limitation the reader may put less emphasis on the non-significant results from these three transfer tasks. However, this ceiling-effect problem does not affect the interpretation of any of the other 14 transfer measures. Moreover, we re-analyzed the spatial and verbal fluid intelligence composite standardized gain scores, after removing Number Series, Paper Folding, and Inferences. None of the ANOVAs on the composites were significant (Spatial fluid intelligence: mid-test, $p = .54$; post-test, $p = .55$; Verbal fluid intelligence: mid-test, $p = .15$; post-test, $p = .78$).

In addition, it was difficult to assess the reliability of the shortened intelligence measures we used. Ideally, we would use the pre-test data in order to calculate Cronbach's alpha, before either the effects of practice or the specific interventions could influence performance (as could occur on the mid- and post-test sessions). However, because subjects performed a different version (A, B, or C) at pre-test, the sample sizes for each version ($N = 20, 24, \text{ and } 29$) were too small to calculate useful Cronbach's alpha information. Likewise, we could instead calculate test-retest reliability, but this would need to be done using only the no-contact control group, which again had a limited sample size ($N = 20$) for a meaningful test-retest correlation.

As alternative measures of reliability, the g factor loadings (Footnote 2) indicate there was substantial systematic variance in our shortened tests. For example, Raven

Standard and RAPM had the highest loadings (.71, .67) and Letter and Number Comparison had the lowest loadings (.21, .16). Note that the two Raven tasks might be suspected of having low reliabilities due to shortening the test, but the *g*-loadings indicate that such suspicions are not warranted. In addition, we calculated the squared multiple correlations as an estimate of reliability reflecting the communalities of the 17 transfer measures. These values ranged from .23 (Vocabulary) to .69 (Number Comparison). None of the values obtained were particularly low – for similar analysis and comparison, see Engle, Tuholski, Laughlin, and Conway (1999).

Note also that shortened measures of tasks such as RAPM have been used in the literature, without adverse effect upon reliability. For example, Arthur and Day (1994) developed a 12-item version of RAPM. Arthur and Day reported in a sample of $N = 461$ young adults that their 12-item version of RAPM had a Cronbach's alpha of .69 and a test-retest reliability of .75. In addition, Basak et al. (2008) divided the RAPM into three sections of 12 items each in order to have a pre-, mid-, and post-test administration. Basak et al. found evidence of transfer to RAPM performance (a Group x Session interaction), indicating again that our use of a shortened RAPM administration did not preclude the possibility of observing transfer.

One final note on the point of reliability: In the current context – which was an *experiment* (and not an individual-differences study) – what is most important is whether or not there are changes produced in the transfer measures as a function of the intervention. That is, in this experimental design, our interest is in between-subjects variability, and more specifically between-subjects variability in the pre-to-post difference score. We know that several experimental effects in cognitive psychology have

low reliability in terms of internal consistency (e.g., Stroop effect difference scores), yet we still appropriately use these measures in experimental research because we want to see whether a particular between-subjects manipulation (e.g., vocal vs. manual responses; proportions of congruent vs. incongruent trials) produces an observable change in the dependent variable (for more on this issue, see Salthouse, Siedlecki, & Krueger, 2006).

WM training and transfer to WM

Although general or broad transfer after repeated practice on a specific task may not be typical, the idea that transfer may occur to tasks that share ‘identical elements’ (Thorndike & Woodworth, 1901) is reasonable. Thus, while the lack of evidence for “far-transfer” from WM training to fluid intelligence may not be surprising, especially given the broader history of intelligence training research (Carroll, 1993), we also did not find evidence for “near-transfer” to two WM span tasks. Across studies, WM training typically leads to changes in untrained WM tasks (Morrison & Chein, 2010; Shipstead et al., in press). Subjects that train on adaptive simple or complex span measures of WM have exhibited transfer to other untrained versions of simple and complex span measures of WM (Bergman Nutley et al., 2011; Klingberg et al., 2005; Chein & Morrison, 2010). Likewise, within the adaptive *n*-back training literature, subjects that train on single or dual *n*-back typically show transfer to untrained versions of the task (JSBSLP, 2010; Seidler et al., 2010).

However, as mentioned in the Introduction, *n*-back tasks and span measures of WM are weakly related to each other (Jaeggi, Buschkuhl et al., 2010; JSBSLP, 2010; Kane et al., 2007; Oberauer, 2005), despite both types of measures correlating with fluid intelligence tasks such as RAPM. If there are no (or few) overlapping processes between

the two types of WM measures, then improvement on one task (n -back) would not likely improve performance on the other (span tasks). Table 6 summarizes the relevant studies that have examined n -back training and transfer across WM task types. There are many differences among the studies (training procedures, sample sizes, transfer tasks etc.), but the results show that transfer across types of WM tasks is inconsistent. Clearly, further work is necessary to understand what different WM processes n -back and span tasks tap, and how these processes overlap with other constructs such as fluid intelligence (Burgess, Gray, Conway, & Braver, 2011; Schmiedek, Hildebrandt, Lövdén, Wilhelm, & Lindenberger, 2009).

Variables that affect transfer

If WM training, and more specifically dual or single n -back training, can actually cause real improvements in fluid intelligence, then the diversity of transfer results across studies indicates that there are important boundary variables that can mediate or moderate training effectiveness. A recent study with children that trained on an adaptive single n -back task identified the amount of n -back improvement as a critical variable determining whether or not transfer to intelligence was observed (Jaeggi, Buschkuhl, Jonides, & Shah, 2011). When Jaeggi et al. (2011) halved the n -back training group based on the amount of improvement observed, the children with the biggest gain showed transfer relative to an active-control group, whereas the children with smaller gains did not. We therefore attempted a similar analysis, by dividing our dual n -back subjects into high and low improvement groups, using a median split on the difference score of mean dual n -back level in sessions 19-20 versus sessions 1-2. This *post-hoc* analysis is limited by sample size (only 12 and 12 subjects in the high- and low-improvement groups,

respectively), but with that caveat in mind, no significant Group (high dual *n*-back improvement, low dual *n*-back improvement, no-contact control) effects were obtained for the fluid intelligence, multitasking, and WM composite standardized gains (Table 7). A similar median-split analysis for the visual search group (15 and 14 subjects in the high- and low-improvement groups, respectively) also produced no significant Group effects on the composite standardized gains (Table 7).

We also correlated the amount of training gain and transfer gain for the same four standardized gain composites (Table 7). Dual *n*-back improvement was not associated with fluid intelligence gains; it was marginally correlated with WM capacity improvement but, surprisingly, visual search improvement was also correlated with improvement on the verbal fluid intelligence tasks (Supplemental Materials, Figure S4). Other WM training studies (Chein & Morrison, 2010; Jaeggi et al., 2011) reporting significant correlations between training change and transfer change suffer from the same limitations as our data for such correlational analyses – small sample sizes and the influence of subjects who performed worse at post-test than pre-test on the transfer tasks (i.e., negative value for transfer gain) and performed worse at the end of training than the beginning of training on the training task (i.e., negative value for training gain). Indeed, in our data, the correlation between visual search change and verbal fluid intelligence change was no longer significant, $r(28) = .25, p = .20$, after removing the lone subject who had negative values on both variables.

Other studies reporting a relationship between WM training gain and fluid intelligence test improvement have been equivocal. In a training study of young adults using adaptive versions of a neutral and an emotional dual *n*-back task, Schweizer,

Hampshire, and Dalgleish (2011) reported a nearly significant correlation between training gain and Raven Standard Progressive Matrices gain. However, Loosli, Buschkuehl, Perrig, and Jaeggi (2012) trained children using an adaptive WM span task and tested fluid intelligence using the Test of Nonverbal Intelligence. There was no transfer, and also no relationship between training gain and matrix reasoning test improvement from pre- to post-test in this study. A potential relationship between the amount of training improvement and the amount of transfer is intuitively appealing. In fact, the data from previous DNB training studies could be re-analyzed to see whether the amount of training improvement affected the amount of intelligence transfer. Of course, if such associations hold up to replication, it would then be important to understand why only some individuals benefit from the training intervention (it would also be important to provide clear evidence that this correlation reflected a causal relation between gain scores).

Clearly, the amount of *n*-back improvement observed is only one possible variable that might affect the presence or amount of transfer. Others include the pre-training ability level of the sample, the size of the samples, the number and duration of training sessions, the transfer test(s) used, the administration method of the transfer tests, the spacing of the training sessions, the motivation of the subjects, the subjects' knowledge about the goals of the study, and the experimenters' influences on subjects' behavior. Because training studies are difficult to conduct in terms of time (here, 23 sessions per training-group subject) and can be financially costly (here, \$352 per training-group subject), it is critical that as many factors as possible are ruled out.

Although a detailed discussion of all of the aforementioned variables is outside the scope of the current article (for a review, see Shipstead et al., in press), three variables warrant further comment here. First, perhaps our lack of transfer represented a lack of motivation by our participants. Motivation is not easily measured, and we agree that subjects who are not motivated to perform either the training or transfer sessions would severely impact the ability to detect improvements as a function of training. Indeed, research has shown that motivation can account for non-ability variance in performance on intelligence tests (Duckworth, Quinn, Lynam, Loeber, & Stouthamer-Loeber, 2011). Note, however, that in addition to self-reported effort not differing between our two training groups, we observed significant increases in performance on the training tasks for both the dual n -back and visual search groups, and significant session effects on six of the transfer measures. Moreover, the amount of dual n -back practice gain in our sample (mean n session 1 = 2.3; mean n session 20 = 4.1) was slightly greater than the dual n -back group in Anguera et al. (2012), but lower than JBJP (2008; 19-session group) and JSBSLP (2010). Although it is potentially important to eventually understand why different patterns of improvement are observed across WM training studies, we argue that no matter the cause, there is nothing fundamentally different about the performance observed in our current experiment compared to the existing studies using adaptive dual n -back training.

Second, given the previous literature on aptitude-treatment interactions (Cronbach & Snow, 1977), and broad cognitive differences between individuals high and low in WM capacity, it would not be unreasonable to think that WM training (treatment) may also vary in its effectiveness depending on a number of factors, including initial WM

level (aptitude). Although the current study was not set up to address the presence of aptitude-treatment interactions directly, we examined initial ability level (fluid intelligence and WM capacity) to see whether pre-test transfer-task scores were related to training gain (a similar analysis to one reported in Jaeggi, Buschkuhl, Jonides, & Shah, 2011). We found that there were no significant correlations between the amount of dual n -back gain and pre-test scores on the WM capacity factor, $r(24) = .14, p = .51$, Fluid intelligence-Spatial factor, $r(24) = .14, p = .51$, or Fluid intelligence-Verbal factor, $r(24) = .29, p = .17$. Note, however, that the sample sizes are rather small for such analyses, and the dual n -back gain variable was a difference score. Other research on intelligence interventions examined the role of pre-existing individual differences in specific intellectual abilities (for review, see Carroll, 1993). Studies by Salomon (1974) and Kyllonen, Lohman, and Snow (1984) demonstrated that pre-training individual differences in verbal and/or spatial abilities interacted with the type of training program, indicating that certain training methods were more or less effective for certain individuals. This is an important consideration for future WM training research as well.

Third, given the diverse nature of WM training programs and procedures being used in research and in commercial applications, empirical analyses of the best practices would be helpful to maximize training efficacy. How many sessions of training are necessary? Schmiedek et al. (2010) administered 100 training sessions, whereas other studies have exhibited transfer after only 3 sessions (Borella et al., 2010). Many WM training programs (Cogmed, n -back) use an adaptive procedure to continually challenge subjects. However, other studies have shown transfer with a static training regimen (Schmiedek et al., 2010). In addition, some WM training procedures involve practicing

on many different tasks during and across practice sessions (Klingberg et al., 2005; Schmiedek et al., 2010), with the idea that diverse practice will consequently cause broader and more general cognitive transfer. In contrast, the single and dual *n*-back training studies use the same task throughout training. Finally, physical-exercise training regimens have been related to cognitive improvements, especially in older adults (Colcombe & Kramer, 2003). More research examining the combined efficacy of exercise and WM training (e.g., Fabre, Chamari, Mucci, Masse-Biron, & Prefaut, 2002) might lead to more effective training procedures, and provide some information about the underlying physiological mechanisms of WM training.

What does WM training actually train?

As indicated by the subjective survey responses, subjects believed that certain aspects of their cognition had been affected by the experiment, even though objectively none of the transfer measures reflected differences over and above practice effects. Our questionnaire findings thus appear to indicate so-called “illusory placebo effects”, whereby trained subjects report subjective improvement in the absence of any objective improvement (see Pratkanis, Eskenazi, & Greenwald, 1994). The potential for such illusions should raise interpretative concerns whenever WM-trained subjects are compared to no-contact controls, as attendant motivational and self-efficacy changes might improve transfer-task performance of trained subjects, even in the absence of any underlying “ability” improvements (Shipstead et al., 2010, in press). This might include persistence in solving the difficult items of intelligence tests – instead of giving up when problems become harder later in the test, the individuals that believe the training has improved their abilities may be more likely to continue to attempt to solve the problem.

Of course, we cannot rule out the possibility that dual n -back training actually did change individual's daily-life abilities, because we did not attempt to verify these behaviors in this study. Other cognitive training studies have shown little to no evidence for positive transfer to the performance of everyday functioning (e.g., Willis et al., 2006). However, in a recent Cogmed training study with young and older adult subjects, transfer was observed to simple span tasks similar to those included in the training program, but not to the Stroop task or Raven Standard Progressive Matrices (Brehmer et al., 2012). Interestingly, subjects in the adaptive training groups self-reported fewer cognitive problems after completing the study, as evidenced by a significant Group by Session interaction on the Cognitive Failures Questionnaire. For both our study and Brehmer et al. (2012), the self-report results could be interpreted as reflecting subjects' implicit ideas about the processes involved in the training tasks, because practice improvements were observed on the training tasks.

More generally, we think that the self-report strategy results highlight the lack of knowledge about what is being trained in dual n -back and other WM training programs. As outlined elsewhere (Shipstead et al., 2010), understanding the mechanisms responsible for transfer in WM training studies is an important goal. Such understanding may require further task-analytic studies of the training procedures in order to isolate the cognitive processes that are being trained or improved. As indicated in our survey data, self-reported strategy use differed not only between subjects but also within subjects; more generally, it seems reasonable to ask whether a subject who has attained an n of 10 via dual n -back training is engaging many (any?) of the same mental processes and strategies in the task as is a subject who attains an n of only 4. A fundamental

understanding of the processes involved in the performance of the dual n -back and other WM training programs is thus important to understanding *why* transfer occurs (if and when it does). Research on the physiological (Jaeggi et al., 2007) and psychometric (Jaeggi, Buschkuhl et al., 2010; JSBSLP, 2010, Study 1; Redick et al., 2012) properties of the dual n -back are informative, but these studies did not use the adaptive versions of the dual n -back tasks, nor did they assess performance across multiple sessions.

Finally, we strongly advocate the full report of all transfer tasks and experimental conditions assessed in WM training studies so that any significant transfer results can be interpreted within the context of measures that did not show significant transfer. The transfer results for any one particular measure must be statistically analyzed and interpreted within the full pattern of results, avoiding selective reporting and uncorrected comparisons. Having full knowledge of previous WM training procedures and results (whether statistically significant or not) will help future researchers identify the mechanisms responsible for transfer and understand the boundary conditions of WM training. On a broader level, in order to minimize the file-drawer problem, publication of studies that do not find transfer adds to the overall context of interpreting the efficacy of WM training.

Conclusion

A critical re-examination of dual n -back training studies indicated the need for replication and extension of previous research, in line with the recommendations of Buschkuhl and Jaeggi (2010), Shipstead et al. (2010), and Sternberg (2008). Subjects in an adaptive dual n -back training group were compared to both an adaptive visual search training (placebo control) group and a no-contact control group. Despite significant

improvements on the training tasks, subjects showed no positive transfer to fluid intelligence, multitasking, WM capacity, crystallized intelligence, or perceptual speed tasks. The current study presents a pessimistic view of the effects of dual n -back practice, but future research might identify variables that maximize the potential intellectual benefits of WM training.

In press - JEP.G

References

- Anguera, J. A., Bernard, J. A., Jaeggi, S. M., Buschkuhl, M., Benson, B. L., Jennett, S., et al. (2012). The effects of working memory resource depletion and training on sensorimotor adaptation. *Behavioural Brain Research*, *228*, 107-115.
- Arthur, Jr., W., & Day, D. V. (1994). Development of a short form for the Raven Advanced Progressive Matrices test. *Educational and Psychological Measurement*, *54*, 394-403.
- Basak, C., Boot, W. R., Voss, M. W., & Kramer, A. F. (2008). Can training in a real-time strategy video game attenuate cognitive decline in older adults? *Psychology and Aging*, *23*, 765-777.
- Berger, F.R., Gupta, W.B., Berger, R.M., & Skinner, J. (1990). *Air Force Officer Qualifying Test (AFOQT) form P: Test Manual (AFHRL-TR-89-56)*. Brooks Air Force Base, TX: Manpower and Personnel Division, Air Force Human Resources Laboratory.
- Bergman Nutley, S., Söderqvist, S., Bryde, S., Thorell, L. B., Humphreys, K., & Klingberg, T. (2011). Gains in fluid intelligence after training non-verbal reasoning in 4-year old children: A controlled, randomized study. *Developmental Science*, *14*, 591-601.
- Borella, E., Carretti, B., Riboldi, F., & De Bini, R. (2008). Working memory training in older adults: Evidence of transfer and maintenance effects. *Psychology and Aging*, *25*, 767-778.

- Brehmer, Y., Westerberg, H., & Bäckman, L. (2012). Working-memory training in younger and older adults: Training gains, transfer, and maintenance. *Frontiers in Human Neuroscience, 6*(63), 1-7.
- Broadway, J. M., & Engle, R. W. (2010). Validating running memory span: Measurement of working memory capacity and links with fluid intelligence. *Behavior Research Methods, 42*, 563-570.
- Bühner, M., König, C., Pick, M., & Krumm, S. (2006). Working memory dimensions as differential predictors of the speed and error aspect of multitasking performance. *Human Performance, 19*, 253-275.
- Burgess, G. C., Gray, J. R., Conway, A. A., & Braver, T. S. (2011). Neural mechanisms of interference control underlie the relationship between fluid intelligence and working memory span. *Journal of Experimental Psychology: General, 140*, 674-692.
- Buschkuehl, M., & Jaeggi, S. M. (2010). Improving intelligence: A literature review. *Swiss Medical Weekly, 140*, 266-272.
- Campbell, D., & Stanley, J. (1963). *Experimental and quasi-experimental designs for research*. Chicago, IL: Rand-McNally.
- Carroll, J. B. (1993). *Human cognitive abilities: A survey of factor-analytic studies*. New York: Cambridge University Press.
- Cattell, R. B. (1973). *Measuring intelligence with the Culture Fair tests*. Champaign, IL: Institute for Personality and Ability Testing.

- Chein, J. & Morrison, A. (2010). Expanding the mind's workspace: Training and transfer effects with a complex working memory span task. *Psychonomic Bulletin & Review*, *17*, 193-199.
- Cohen, G. L., Garcia, J., Apfel, N., & Master, A. (2006). Reducing the racial achievement gap: A social-psychological intervention. *Science*, *313*, 1307-1310.
- Colcombe, S. & Kramer, A.F. (2003). Fitness effects on the cognitive function of older adults: A meta-analytic study. *Psychological Science*, *14*, 125-130.
- Colom, R., Ángeles Quiroga, M., Shih, P.-C., Martínez, K., Burgaleta, M., Martínez-Molina, A., et al. (2010). Improvement in working memory is not related to increased intelligence scores. *Intelligence*, *38*, 497-505.
- Colom R., Martínez-Molina A., Shih P., & Santacreu J. (2010). Intelligence, working memory, and multitasking performance. *Intelligence*, *38*, 543-551,
- Conway, A. R. A., & Getz, S. J. (2010). Cognitive ability: Does working memory training enhance intelligence? *Current Biology*, *20*, R362–R364.
- Cronbach, L. J., & Snow, R. E. (1977). *Aptitudes and instructional methods: A handbook for research on interactions*. Oxford England: Irvington.
- Dahlin, E., Nyberg, L., Bäckman, L., & Neely, A. S. (2008). Plasticity of executive functioning in young and older adults: Immediate training gains, transfer, and long-term maintenance. *Psychology and Aging*, *23*, 720–730.
- Duckworth, A., Quinn, P. D., Lynam, D. R., Loeber, R., & Stouthamer-Loeber, M. (2011). Role of test motivation in intelligence testing. *Proceedings of the National Academy of Sciences of the United States of America*, *108*, 7716-7720.

- Elsmore, T. F. (1994). SYNWORK: A PC-based tool for assessment of performance in a simulated work environment. *Behavior Research Methods, Instrumentation, and Computers*, 26, 421-426.
- Ekstrom, R. B., French, J. W., Harman, H. H., & Dermen, D. (1976). *Manual for kit of factor-referenced cognitive tests*. Princeton, NJ: Educational Testing Service.
- Engle, R. W., Tuholski, S. W., Laughlin, J. E., & Conway, A. R. A. (1999). Working memory, short-term memory and general fluid intelligence: A latent-variable approach. *Journal of Experimental Psychology: General*, 128, 309-331.
- Fabre, C., Chamari, K., Mucci, P., Masse-Biron, J., & Prefaut, C. (2002). Improvement of cognitive function by mental and/or individualized aerobic training in healthy elderly subjects. *International Journal of Sports Medicine*, 23, 415-421.
- Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G*Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, 39, 175-191.
- Fothergill, S., Loft, S. D., & Neal, A. (2009). ATC-labAdvanced: An air traffic control simulator with realism and control. *Behavior Research Methods*, 41, 118-127.
- French, J. R. P. (1953). Experiments in field settings. In L. Festinger & D. Katz (Eds.), *Research Methods in the Behavioral Sciences* (pp. 98-135). New York: Holt, Rinehart, and Wilson.
- Hambrick, D. Z., Oswald, F. L., Darowski, E. S., Rench, T. A., & Brou, R. (2010). Predictors of multitasking performance in a synthetic work paradigm. *Applied Cognitive Psychology*, 24, 1149-1167.

- Highfield, R. (2008, April 28). 'Brain training' games do work, study finds. The Telegraph. <http://www.telegraph.co.uk/science/science-news/3340967/Brain-training-games-do-work-study-finds.html>
- Jaeggi, S., Buschkuhl, M., Jonides, J., & Shah, P. (2011). Short- and long-term benefits of cognitive training. *Proceedings of the National Academy of Sciences of the United States of America*, *108*, 10081-10086.
- Jaeggi, S. M. (2005). *Capacity limitations in human cognition: Behavioural and biological considerations*. Unpublished doctoral dissertation, University of Bern, Switzerland.
- Jaeggi, S. M., Buschkuhl, M., Etienne, A., Ozdoba, C., Perrig, W. J., & Nirrko, A. C. (2007). On how high performers keep cool brains in situations of cognitive control. *Cognitive, Affective, & Behavioral Neuroscience*, *7*, 75-89.
- Jaeggi, S. M., Buschkuhl, M., Jonides, J., & Perrig, W. J. (2008). Improving fluid intelligence with training on working memory. *Proceedings of the National Academy of Sciences of the United States of America*, *105*, 6829-6833.
- Jaeggi, S. M., Buschkuhl, M., Perrig, W. J., & Meier, B. (2010). The concurrent validity of the *N*-back task as a working memory measure. *Memory*, *18*, 394-412.
- Jaeggi, S. M., Studer-Luethi, B., Buschkuhl, M., Su, Y.-F., Jonides, J., & Perrig, W. J. (2010). The relationship between *n*-back performance and matrix reasoning - Implications for training and transfer. *Intelligence*, *38*, 625-635.
- Jensen, A. R. (1998). *The g factor: The science of mental ability*. Westport, CT: Praeger.

- Kane, M. J., Hambrick, D. Z., & Conway, A. R. A. (2005). Working memory capacity and fluid intelligence are strongly related constructs: Comment on Ackerman, Beier, and Boyle (2005). *Psychological Bulletin*, *131*, 66-71.
- Kane, M. J., Conway, A. R. A., Hambrick, D. Z., & Engle, R. W. (2007). Variation in working-memory capacity as variation in executive attention and control. In A. R. A. Conway, C. Jarrold, M. J. Kane, A. Miyake, & J. N. Towse (Eds.), *Variation in working memory* (pp. 21-48). New York: Oxford University Press.
- Kane, M. J., Conway, A. R. A., Miura, T. K., & Colflesh, G. J. H. (2007). Working memory, attention control, and the n-back task: A question of construct validity. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *33*, 615-622.
- Kane, M. J., Poole, B. J., Tuholski, S. W., & Engle, R. W. (2006). Working memory capacity and the top-down control of visual search: Exploring the boundaries of 'executive attention'. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *32*, 749-777.
- Klingberg, T. (2010). Training and plasticity of working memory. *Trends in Cognitive Science*, *14*, 317-324.
- Klingberg, T., Forssberg, H., & Westerberg, H. (2002). Training of working memory in children with ADHD. *Journal of Clinical and Experimental Neuropsychology*, *24*, 781-791.
- Klingberg, T., Fernell, E., Olesen, P., Johnson, M., Gustafsson, P., Dahlström, K., Gillberg, C. G., Forssberg, H., & Westerberg, H. (2005). Computerized training of working memory in children with ADHD – A randomized, controlled, trial.

- Journal of the American Academy of Child and Adolescent Psychiatry*, 44, 177-186.
- Kyllonen, P. C., Lohman, D. F., & Snow, R. E. (1984). Effects of aptitudes, strategy training, and task facets on spatial task performance. *Journal of Educational Psychology*, 76, 130-145.
- Li, S.-C., Schmiedek, F., Huxhold, O., Röcke, C., Smith, J., & Lindenberger, U. (2008). Working memory plasticity in old age: Practice gain, transfer, and maintenance. *Psychology and Aging*, 23, 731-742.
- Loosli, S. V., Buschkuehl, M., Perrig, W. J., & Jaeggi, S. M. (2012). Working memory training improves reading processes in typically developing children. *Child Neuropsychology*, 18, 62-78.
- Marshalek, B., Lohman, D. F., & Snow, R. E. (1983). The complexity continuum in the radex and hierarchical models of intelligence. *Intelligence*, 7, 107-127.
- Medalia, A., & Freilich, B. (2008). The Neuropsychological Educational Approach to Cognitive Remediation (NEAR) model: Practice principles and outcome studies. *American Journal of Psychiatric Rehabilitation*, 11, 123-43.
- Melby-Lervåg, M., & Hulme, C. (in press). Is working memory training effective? A meta-analytic review. *Developmental Psychology*.
- Morrison, A., & Chein, J. (2011). Does working memory training work? The promise and challenges of enhancing cognition by training working memory. *Psychonomic Bulletin & Review*, 18, 46-60.

- Oberauer, K. (2005). Binding and inhibition in working memory: Individual and age differences in short-term recognition. *Journal of Experimental Psychology: General*, 134, 368-387.
- Oberauer, K., Schulze, R., Wilhelm, O., & Süß, H. (2005). Working memory and intelligence – Their correlation and their relation: Comment on Ackerman, Beier, and Boyle (2005). *Psychological Bulletin*, 131, 61-65.
- Olesen, P. J., Westerberg, H., & Klingberg, T. (2004). Increased prefrontal and parietal activity after training of working memory. *Nature Neuroscience*, 7, 75-79.
- Owen, A. M., Hampshire, A., Grahn, J. A., Stenton, R., Dajani, S., Burns, A. S., et al. (2010). Putting brain training to the test. *Nature*, 465, 775-776.
- Perrig, W. J., Hollenstein, M., & Oelhafen, S. (2009). Can we improve fluid intelligence with training on working memory in persons with intellectual disabilities? *Journal of Cognitive Education and Psychology*, 8, 148-164.
- Pratkanis, A. R., Eskenazi, J., & Greenwald, A. G. (1994). What you expect is what you believe (but not necessarily what you get): A test of the effectiveness of subliminal self-help audiotapes. *Basic and Applied Social Psychology*, 15, 251-276.
- Raven, J., Raven, J. C., Court, J. H. (1998). *Manual for Raven's Progressive Matrices and Vocabulary Scales*. New York: Psychological Corporation.
- Redick, T. S., Broadway, J. M., Meier, M. E., Kuriakose, P. S., Unsworth, N., Kane, M. J., & Engle, R. W. (in press). Measuring working memory capacity with automated complex span tasks. *European Journal of Psychological Assessment*.

- Redick, T. S., Shipstead, Z., Meier, M. E., Evans, J., Hicks, K. L., Unsworth, N. et al. (2012). *Cognitive predictors of multitasking performance*. Manuscript in progress.
- Richmond, L. L., Morrison, A. B., Chein, J. M., & Olson, I. R. (2011). Working memory training and transfer in older adults. *Psychology and Aging, 26*, 813-822.
- Salomon, G. (1974). Internalization of filmic schematic operations in interaction with learners' aptitudes. *Journal of Educational Psychology, 66*, 499-511.
- Salthouse, T. A., & Babcock, R. L. (1991). Decomposing adult age differences in working memory. *Developmental Psychology, 27*, 763-776.
- Schmiedek, F., Hildebrandt, A., Lövdén, M., Wilhelm, O., & Lindenberger, U. (2009). Complex span versus updating tasks of working memory: The gap is not that deep. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 35*, 1089-1096.
- Schmiedek, F., Lövdén, M., & Lindenberger, U. (2010). Hundred days of cognitive training enhance broad abilities in adulthood: Findings from the COGITO study. *Frontiers in Aging Neuroscience, 2*, 1-10.
- Schweizer, S., Hampshire, A., & Dalgleish, T. (2011). Extending brain-training to the affective domain: Increasing cognitive and affective executive control through emotional working memory training. *PLoS ONE, 6*, e24372.
- Seidler, R. D., Bernard, J. A., Buschkuhl, M., Jaeggi, S., Jonides, J. & Humfleet, J. (2010). *Cognitive training as an intervention to improve driving ability in the older adult* (Technical Report No. M-CASTL 2010-01). Ann Arbor, MI.

- Shellenbarger, S. (2011, November 29). Ways to inflate your IQ. New York Times.
http://online.wsj.com/article_email/SB10001424052970203935604577066293669642830-1MyQjAxMTAxMDIwOTEyNDkyWj.html?mod=wsj_share_email
- Shipstead, Z., Redick, T. S., & Engle, R. W. (2010). Does working memory training generalize? *Psychologica Belgica*, 3-4, 245-276.
- Shipstead, Z., Redick, T. S., & Engle, R. W. (in press). Is working memory training effective? *Psychological Bulletin*.
- Sternberg, R. J. (2008). Increasing fluid intelligence is possible after all. *Proceedings of the National Academy of Sciences of the United States of America*, 105, 6791-6792.
- Thorell, L. B., Lindqvist, S., Bergman, S., Bohlin, G., & Klingberg, T. (2009). Training and transfer effects of executive functions in preschool children. *Developmental Science*, 11, 969-976.
- Thorndike, E. L., & Woodworth, R. S. (1901). The influence of improvement in one mental function upon the efficiency of other functions (I). *Psychological Review*, 8, 247-261.
- Thurstone, L. L. (1938). *Primary mental abilities*. Chicago: University of Chicago Press.
- Wang, S., & Aamodt, S. (2009, March 10). Guest column: Can we increase our intelligence? New York Times. <http://judson.blogs.nytimes.com/2009/03/10/guest-column-can-we-increase-our-intelligence/>
- Westerberg, H., & Klingberg, T. (2007). Changes in cortical activity after training of working memory – a single-subject analysis. *Physiology & Behavior*, 92, 186-192.

- Willis, S. L., Tennstedt, S. L., Marsiske, M., Ball, K., Elias, J., Mann Koepke, K., et al. (2006). Long-term effects of cognitive training on everyday functional outcomes in older adults. *JAMA*, *296*, 2805-2814.
- Zachary, R. A. (1986). *Shipley Institute of Living Scale: Revised Manual*. Los Angeles: Western Psychological Services.
- Zinke, K., Zeintl, M., Eschen, A., Herzog, C., & Kliegel, M. (2012). Potentials and limits of plasticity induced by working memory training in old-old age. *Gerontology*, *58*, 79-87.

In press - JEP.G

Acknowledgements

The research was conducted with support by a grant from the Office of Naval Research (Award No. N000140910129). We gratefully acknowledge Susanne Jaeggi for providing the materials for the adaptive dual n -back task, along with providing data and answers to several questions about previous studies. We thank Nash Unsworth and Gene Brewer for providing helpful comments on earlier versions of the manuscript. We thank Forough Azimi, Jameil Bailey, Kati Frampton, Rakshya Khatri, Christine Kim, Princy Kuriakose, Dakota Lindsay, Shannon Lukey, Robyn Marshall, Kevin Strika, Maureen van der Water, Tamara Ware, and Allison Vercellone for their assistance in data collection. We thank Marlyssa Fillmore for assistance in preparing Table 2.

Correspondence concerning this article should be sent to Thomas S. Redick, Indiana University Purdue University Columbus, 4601 Central Avenue, Columbus, IN 47203 (email to tsredick@iupuc.edu).

Footnotes

¹ Jaeggi (2005) refers to a dissertation that presents the full results of the 8- and 19-session studies included in JBJP (2008). This dissertation was downloaded from the following URL (http://www.zb.unibe.ch/download/eldiss/05jaeggi_s.pdf) and has been archived at <http://webcitation.org/662gjXMf3>.

² Our plan in designing the study was to use a correlated-vectors approach (Jensen, 1998) to interpret the pattern of expected transfer, similar to the logic used in Colom, Ángeles Quiroga et al. (2010). We conducted this analysis, but because we observed no transfer for any tasks, it seemed unnecessary and redundant. We examined the pre-test data for all subjects (N = 123 with data on all 17 pre-test measures), and extracted one factor (principal axis factoring) from all of the tests. The relative ordering of the tests indicated that, as predicted, the bulk of the fluid intelligence and multitasks had high loadings on the general factor, and the perceptual speed tasks had the lowest loadings.

1. Raven Standard (Gf)	.71
2. Raven Advanced (Gf)	.67
3. Cattell's (Gf)	.66
4. Letter Sets (Gf)	.62
5. Control Tower (Multi)	.62
6. Number Series (Gf)	.60
7. SynWin (Multi)	.56
8. Running Letter Span (WM)	.52
9. General Knowledge (Gc)	.52
10. ATClab (Multi)	.52
11. Paper Folding (Gf)	.48
12. Symmetry Span (WM)	.46
13. Inferences (Gf)	.46
14. Analogies (Gf)	.45
15. Vocabulary (Gc)	.29
16. Number Comparison (PS)	.21
17. Letter Comparison (PS)	.16

³ Although these are the averages of all of the subjects in the final sample, it is possible that individuals in a particular group may have spent more or less time working on the 17 transfer tasks during the pre-, mid-, and post-test sessions. However, the duration of the sessions did not differ among the groups: Pre-test, $F(2, 70) = 0.39, p = .68$; Mid-test, $F(2, 70) = 0.26, p = .77$; Post-test, $F(2, 70) = 1.46, p = .24$.

⁴ Due to experimenter error, one dual *n*-back subject completed the mid-test session after eight practice sessions.

⁵ Using an alpha of .05, the Group x Session interaction for Raven Standard Progressive Matrices was significant (Table 2). However, as shown in Figure A.2, the interaction did not reflect improvements for the dual *n*-back or visual search group relative to controls; simple main effects analyses (one-way ANOVAs) indicated that the three groups did not significantly differ from each other at pre-, mid-, or post-test (all p 's $> .26$). Repeated-measures ANOVAs calculated separately for each group revealed that the interaction was due to a significant effect of Session only for the no-contact control group, $F(2, 38) = 5.02, p = .01, \eta_p^2 = .21$; for the dual *n*-back and visual search groups, F 's < 1 .

⁶ The only slight difference between the ANCOVA and ANOVA results came from the Control Tower task, where the Group x Session interaction ($p = .037$) indicated a potential difference between the groups in the nature of their session effect from mid- to post-test after covarying pre-test performance. Similar to the RSPM results (Footnote 5), simple main effects analyses (one-way ANOVAs) indicated that the three groups did not significantly differ from each other at pre-, mid-, or post-test (all p 's $> .59$). Also similar

to the RSPM results, inspection of Table 2 indicates that the direction of the difference between the mid- and post-test sessions likely drove the ANCOVA interaction result. Whereas the control and dual *n*-back groups showed slight numerical increases from mid- to post-test, the visual search group showed a numerically larger improvement. However, note that numerically the visual search group showed a slightly lower mid-test score than the other two groups, indicating more room for improvement from mid- to post-test. Most importantly, there is no evidence here for greater Control Tower transfer gain from dual *n*-back training than from visual search training or no contact.

⁷ Although Seidler et al. (2010) is a technical report, the data subsume the transfer results reported in Experiment 2 of the recently published Anguera et al. (2012) journal article (personal communication, R. D. Seidler, December 27, 2011). Anguera et al. (2012) reported that: “Participants completed two days of pre-testing as part of a larger study. Of relevance here, the test battery included a working memory assessment using an *n*-back task (*n* = 3 and 4) with abstract shapes, an automated operation span task, as well as the card rotation task and the digit symbol substitution task from Experiment 1, and finally, a visuomotor adaptation task.” (p. 110). Anguera et al. reported positive transfer to the *n*-back and Operation Span tasks.

Seidler et al. (2010) included these transfer tests, but also listed seven other tests that did not exhibit significant dual *n*-back transfer: BOMAT, RAPM, verbal analogies, visual arrays comparison, Attention Network Test, motor sequence learning task, and various conditions of a driving simulator task. Perhaps most importantly, the technical report provides an attempt at replication as the reason for including the fluid intelligence measures: “Type 2 tests included Raven’s matrices (Raven et al., 1990), which is a

standardized test of fluid intelligence, and the BOMAT and verbal analogies tests of intelligence (Hossiep et al., 1995). We have previously shown that working memory training transfers to performance on this task (Jaeggi et al., 2008), and we included it here for the sake of replication” (p. 7).

The Seidler et al. (2010) technical report was downloaded from the following URL (<http://m-castl.org/files/2010-01SeidlerReport.pdf>) and has been archived at <http://www.webcitation.org/662goIQRW>.

In press - JEP.G

Table 1

Demographic information

Group	N	M/F	Age	GT	GSU	MSU	Other
Dual <i>N</i> -back	24	10/14	21.1 (2.7)	9	7	7	1
Visual Search	29	12/17	20.7 (2.5)	9	11	8	1
Control	20	10/10	21.2 (2.5)	7	7	5	1

Note. GT: Georgia Tech student; GSU: Georgia State student; MSU: Michigan State student; Other: Not currently attending one of these three colleges.

Table 2

Mean Performance for the Transfer Tasks as a Function of Training Group and Session

Task	<u>Control</u>			<u>Visual Search</u>			<u>Dual <i>n</i>-back</u>		
	Pre	Mid	Post	Pre	Mid	Post	Pre	Mid	Post
RAPM	6.65 (2.18)	6.45 (2.50)	6.00 (3.00)	6.52 (3.04)	6.07 (2.87)	6.24 (3.34)	7.04 (2.48)	6.17 (2.28)	6.25 (3.08)
RSPM ^a	17.15 (2.39)	15.85 (2.50)	16.85 (2.35)	16.66 (2.53)	16.34 (2.30)	16.45 (2.47)	16.30 (2.67)	16.74 (2.54)	16.09 (2.61)
Cattell	11.95 (2.63)	11.75 (2.05)	11.45 (2.65)	10.72 (2.79)	11.07 (2.07)	11.24 (2.25)	12.00 (2.38)	11.71 (2.29)	11.38 (2.45)
Paper Folding	4.05 (1.70)	4.50 (1.43)	4.00 (1.26)	4.41 (1.38)	4.00 (1.60)	4.52 (1.34)	3.79 (1.47)	4.46 (1.69)	4.33 (1.34)
Letter Sets	6.85 (1.90)	6.75 (2.29)	6.80 (2.22)	7.79 (1.84)	6.90 (2.16)	6.83 (2.19)	7.08 (2.45)	7.17 (1.52)	7.04 (2.14)

WM TRAINING 65

Number Series	4.20	3.75	3.70	3.59	3.76	3.52	3.96	3.92	3.75
	(0.83)	(0.85)	(1.22)	(1.32)	(1.35)	(1.24)	(0.96)	(1.18)	(1.19)
Inferences	4.35	4.30	4.45	4.41	4.03	4.24	3.67	4.04	4.04
	(1.31)	(1.78)	(1.54)	(1.40)	(1.84)	(1.60)	(1.97)	(1.60)	(1.65)
Analogies	4.90	4.65	3.90	4.83	4.45	4.38	4.46	4.46	3.79
	(1.59)	(1.42)	(1.71)	(1.65)	(1.64)	(1.66)	(1.62)	(1.53)	(1.50)
SynWin	352.40	682.50	701.50	461.14	625.76	729.14	480.28	581.88	655.08
	(626.95)	(190.05)	(214.60)	(252.22)	(205.12)	(193.33)	(218.08)	(231.66)	(201.34)
ControlTower ^b	29.62	32.47	34.05	29.90	29.63	37.43	29.20	31.41	34.26
	(10.45)	(11.22)	(11.02)	(11.38)	(13.87)	(15.63)	(8.10)	(11.54)	(11.29)
ATClab ^c	0.72	0.73	0.72	0.73	0.71	0.75	0.74	0.75	0.75
	(.09)	(.12)	(.12)	(.14)	(.12)	(.12)	(.13)	(.09)	(.12)
Symmetry Span	25.60	30.25	28.90	24.55	27.28	26.76	25.88	32.29	31.54
	(9.30)	(9.34)	(12.14)	(11.15)	(12.53)	(10.99)	(8.75)	(9.92)	(11.80)
Running Span	38.50	40.90	43.00	39.52	39.52	42.34	37.96	40.13	42.21
	(8.65)	(8.81)	(9.06)	(13.03)	(12.42)	(12.49)	(12.69)	(10.86)	(11.94)

WM TRAINING 66

Vocabulary	10.10	10.70	10.35	10.00	10.38	9.79	10.33	10.04	10.50
	(1.21)	(1.08)	(1.79)	(1.79)	(1.43)	(1.74)	(1.13)	(1.40)	(1.41)
Knowledge	6.75	6.30	6.20	5.90	6.17	6.10	6.25	6.04	6.29
	(1.83)	(2.18)	(1.61)	(1.99)	(1.91)	(1.78)	(2.03)	(1.23)	(1.81)
Letter Comparison	18.75	20.65	20.85	19.93	20.45	20.45	19.04	19.92	21.38
	(3.77)	(3.50)	(3.72)	(4.08)	(5.68)	(5.38)	(4.84)	(4.03)	(3.61)
Number Comparison	28.90	31.15	31.00	29.14	29.52	29.93	28.83	28.58	29.00
	(5.21)	(4.61)	(4.24)	(6.11)	(5.92)	(7.02)	(5.54)	(4.95)	(5.43)

Note. ^a N = 23 for dual *n*-back group due to experimenter error during mid-test session. ^b N = 19 for dual *n*-back group due to computer problem during post-test session. ^c N = 28 for visual search group and N = 19 for control group due to computer problem during post-test session.

Table 3

Significance testing results for the transfer measures

Task	Group			Session			Group x Session		
	<i>F</i>	<i>p</i>	η_p^2	<i>F</i>	<i>p</i>	η_p^2	<i>F</i>	<i>p</i>	η_p^2
<u>Fluid Intelligence (Spatial)</u>									
Raven Advanced	0.05	.95	.001	2.00	.14	.028	0.33	.86	.009
Raven Standard	0.06	.94	.002	1.53	.22	.022	2.85	.03	.076
Cattell	1.08	.35	.030	0.28	.75	.004	0.92	.45	.026
Paper Folding	0.11	.90	.003	0.75	.47	.011	1.92	.11	.052
<u>Fluid Intelligence (Verbal/Numeric)</u>									
Letter Sets	0.31	.74	.009	1.14	.32	.016	1.08	.37	.030
Number Series	0.70	.50	.020	1.56	.21	.022	0.78	.54	.022
Inferences	0.70	.50	.020	0.20	.82	.003	0.67	.61	.019
Analogies	0.46	.64	.013	6.08	.00	.080	0.69	.60	.019
<u>Multitasking</u>									
SynWin ^a	0.16	.85	.005	29.95	.00	.300	1.81	.13	.049
ControlTower	0.26	.98	.001	17.28	.00	.200	1.96	.10	.054
ATClab	0.23	.78	.007	0.32	.73	.005	0.80	.53	.023
<u>Working Memory Capacity</u>									
Symmetry Span	1.02	.37	.028	10.30	.00	.128	0.70	.59	.020
Running Letter Span	0.24	.98	.001	8.64	.00	.110	0.38	.82	.011

Crystallized Intelligence

Vocabulary	0.57	.57	.016	0.68	.51	.010	1.62	.17	.044
------------	------	-----	------	------	-----	------	------	-----	------

General Knowledge	0.37	.70	.010	0.17	.85	.002	0.68	.61	.019
-------------------	------	-----	------	------	-----	------	------	-----	------

Perceptual Speed

Letter Comparison	0.02	.98	.001	5.55	.01	.073	0.98	.42	.027
-------------------	------	-----	------	-------------	------------	-------------	------	-----	------

Number Comparison	0.58	.57	.016	1.54	.22	.022	0.75	.56	.021
-------------------	------	-----	------	------	-----	------	------	-----	------

Note. ^a At the MSU testing location, N = 21 subjects were administered the same test version of SynWin at pre-, mid-, and post-test. Data were re-analyzed using only subjects that performed unique versions of SynWin at all three transfer sessions (N = 14/21/17 for C/VIS/DNB, respectively). The interpretation of the significance tests were the same as listed above with the full data. **Bold** entries indicate values significant at alpha = .01.

Table 4

Inferential results of the transfer composite standardized gain scores

Construct	Mid	Post
Fluid Intelligence (Spatial)	$F(2, 69) = 0.83, p = .44$	$F(2, 70) = 0.39, p = .68$
Fluid Intelligence (Verbal)	$F(2, 70) = 1.51, p = .23$	$F(2, 70) = 0.62, p = .54$
Multitasking	$F(2, 70) = 3.09, p = .05$	$F(2, 67) = 1.44, p = .24$
Working Memory	$F(2, 70) = 1.89, p = .16$	$F(2, 70) = 0.89, p = .41$
Crystallized Intelligence	$F(2, 70) = 1.44, p = .24$	$F(2, 70) = 0.16, p = .86$
Perceptual Speed	$F(2, 70) = 1.15, p = .32$	$F(2, 70) = 0.86, p = .43$

Table 5

Post-test survey data

Topic	Dual <i>N</i> -Back	Visual Search	Control	$\chi^2(2)$
Attention	52%	72%	50%	3.27 ($p = .20$)
Intelligence	65%	41%	30%	5.73 ($p = .06$)
Language	4%	3%	10%	1.06 ($p = .59$)
Memory	78%	45%	40%	8.01 ($p = .02$)
Perception	35%	59%	45%	2.98 ($p = .23$)
Daily activities	43%	10%	10%	10.51 ($p < .01$)

Note. Due to experimenter error, survey data were not available for one dual *n*-back subject and two control subjects. However, survey data were included for the two control subjects who received the same transfer test items at pre-test and post-test. The format of the question for each topic was “Do you feel that your participation in this study has changed your _____?”.

Table 6
WM training and transfer effects across task types

Study	Training	<i>N</i> -back transfer	Span transfer
<i>N</i> -back training			
Jaeggi et al. (2008) ^a	Dual	Dual	Reading
JSBSLP (2010)	Single, Dual	Single	Operation
Seidler et al. (2010)	Dual	Single, Dual	Operation
Li et al. (2008)			
Young/Older adults	Single	Single/Single	Operation/Operation Rotation/Rotation
Schmiedek et al. (2010)			
Young/Older adults	Single ^b	Single/Single	Reading/Reading Counting/Counting Rotation/ Rotation

Note. Tasks in **bold** produced significant ($p < .05$) Group (Training vs. Control) by Session (Pre- vs. Post-test) interactions. ^a Significant transfer for 19-session group reported in Jaeggi (2005). ^b Task was part of a battery of training tasks administered to subjects.

Table 7

Standardized gain composite transfer results based on amount of training gain

Composite	<u>Pre-to-Post</u>			<u>Correlation</u>	
	<i>F</i>	<i>p</i>	η_p^2	<i>r</i>	<i>p</i>
	Dual <i>n</i> -back (<i>n</i> = 24)				
Fluid Intelligence (Spatial)	0.83	.44	.039	.24	.26
Fluid Intelligence (Verbal)	0.68	.51	.032	-.19	.37
Multitasking	0.58	.56	.029	.30 ^a	.17
Working Memory Capacity	1.68	.20	.076	.39	.06
	Visual search (<i>n</i> = 29)				
Fluid Intelligence (Spatial)	0.66	.52	.028	-.10	.60
Fluid Intelligence (Verbal)	0.32	.73	.014	.36	.06
Multitasking	0.84	.44	.037	-.07 ^b	.71
Working Memory Capacity	2.21	.12	.088	-.19	.33

Note. ^a N = 23. ^b N = 28.

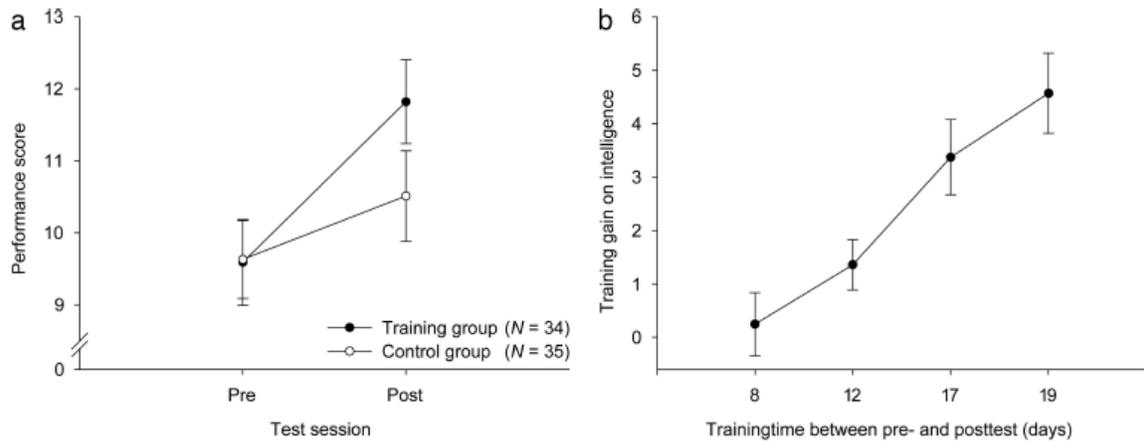


Figure 1. Matrix reasoning test performance as a function of group and session (a) and RAPM or BOMAT gain as a function of the number of dual n -back sessions completed (b). Reprinted from “Improving fluid intelligence with training on working memory,” by S. M. Jaeggi, M. Buschkuhl, J. Jonides, & W. J. Perrig, 2008, *Proceedings of the National Academy of Sciences*, 105, p. 6831. Copyright 2008 by the National Academy of Sciences.

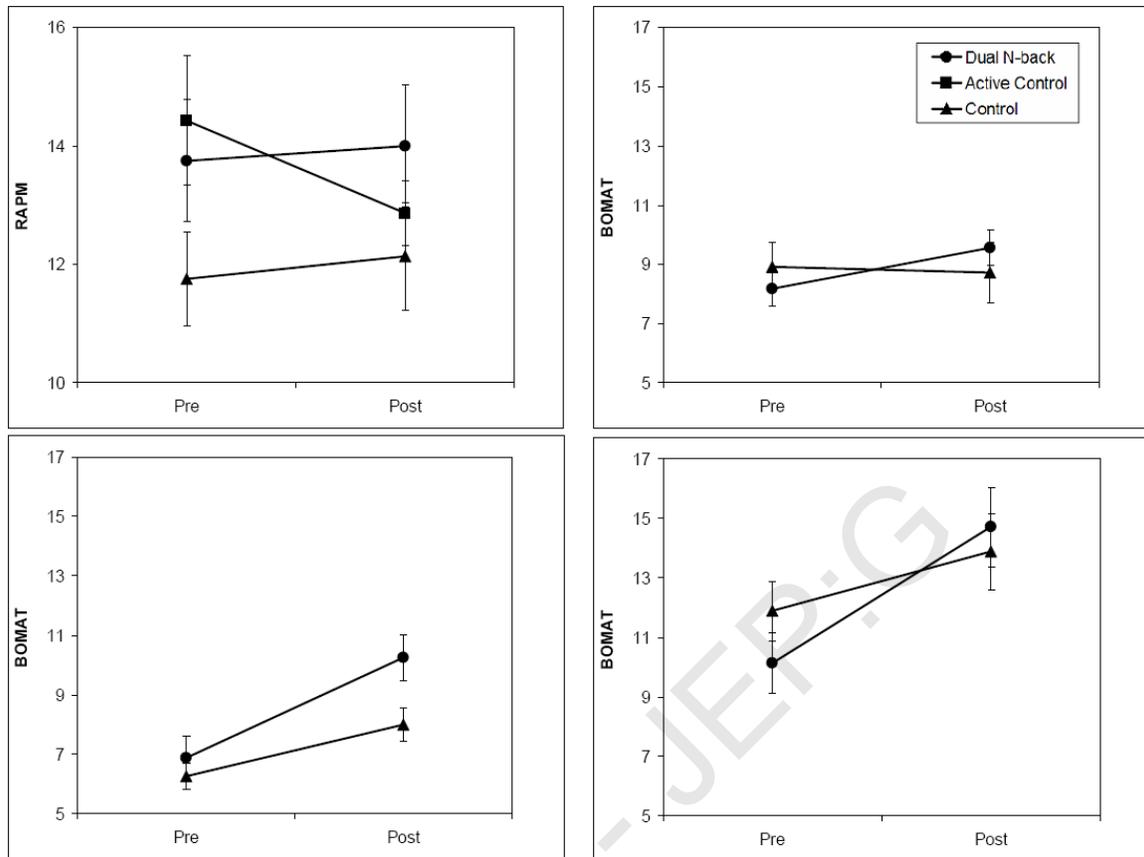


Figure 2. RAPM performance for the 8-session study (top left) and BOMAT performance for the 12- (top right), 17- (bottom left), and 19-session (bottom right) studies. Data provided by S. M. Jaeggi (personal communication, March 31, 2012). Error bars represent ± 1 standard error of the mean.

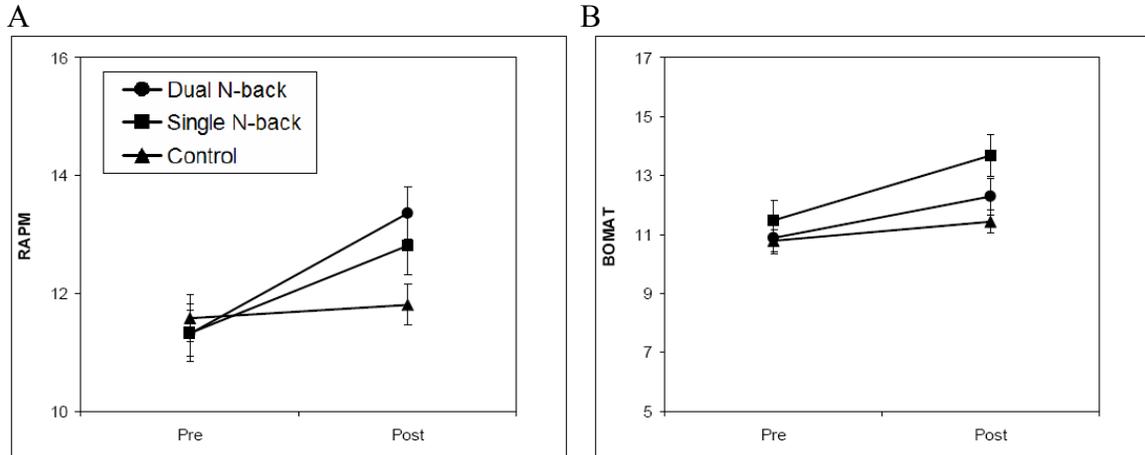


Figure 3. RAPM (a) and BOMAT (b) performance for the three groups in JSBSJP (2010). Error bars represent ± 1 standard error of the mean.

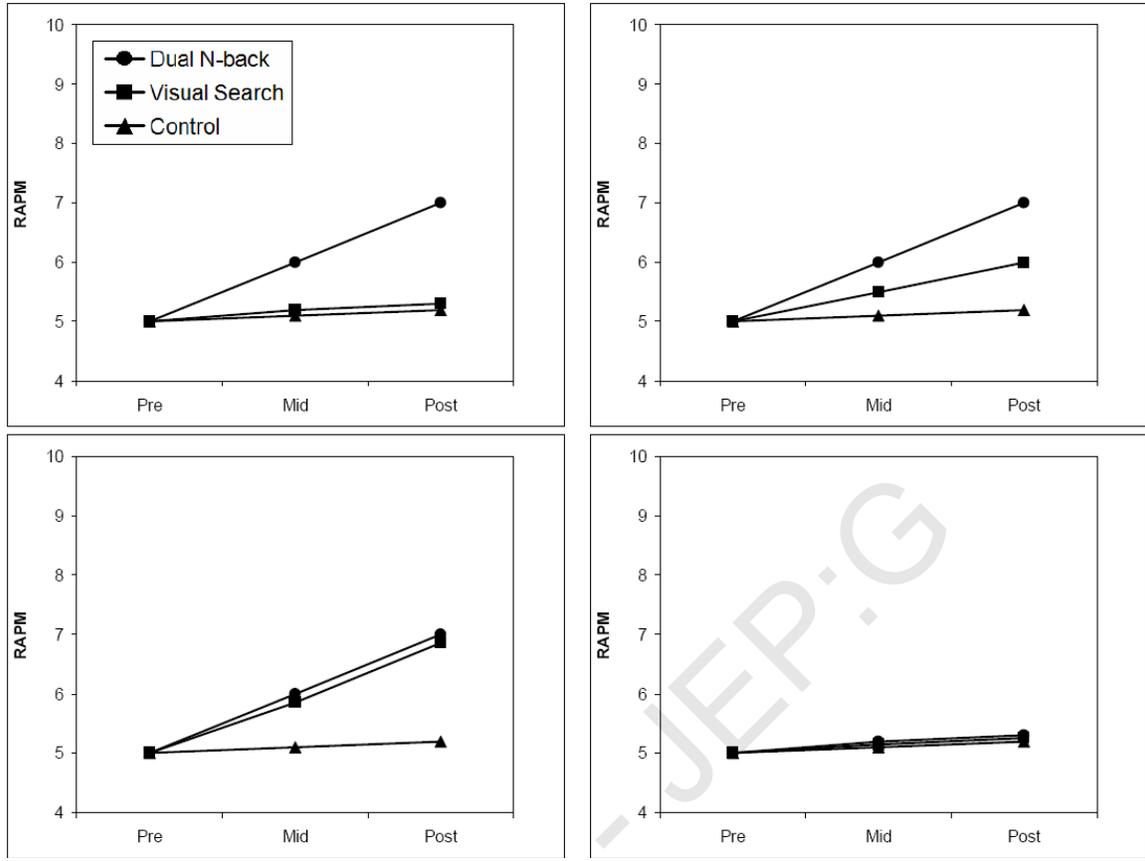


Figure 4. Four possible outcomes of current study (see text).

In press - JEP.G

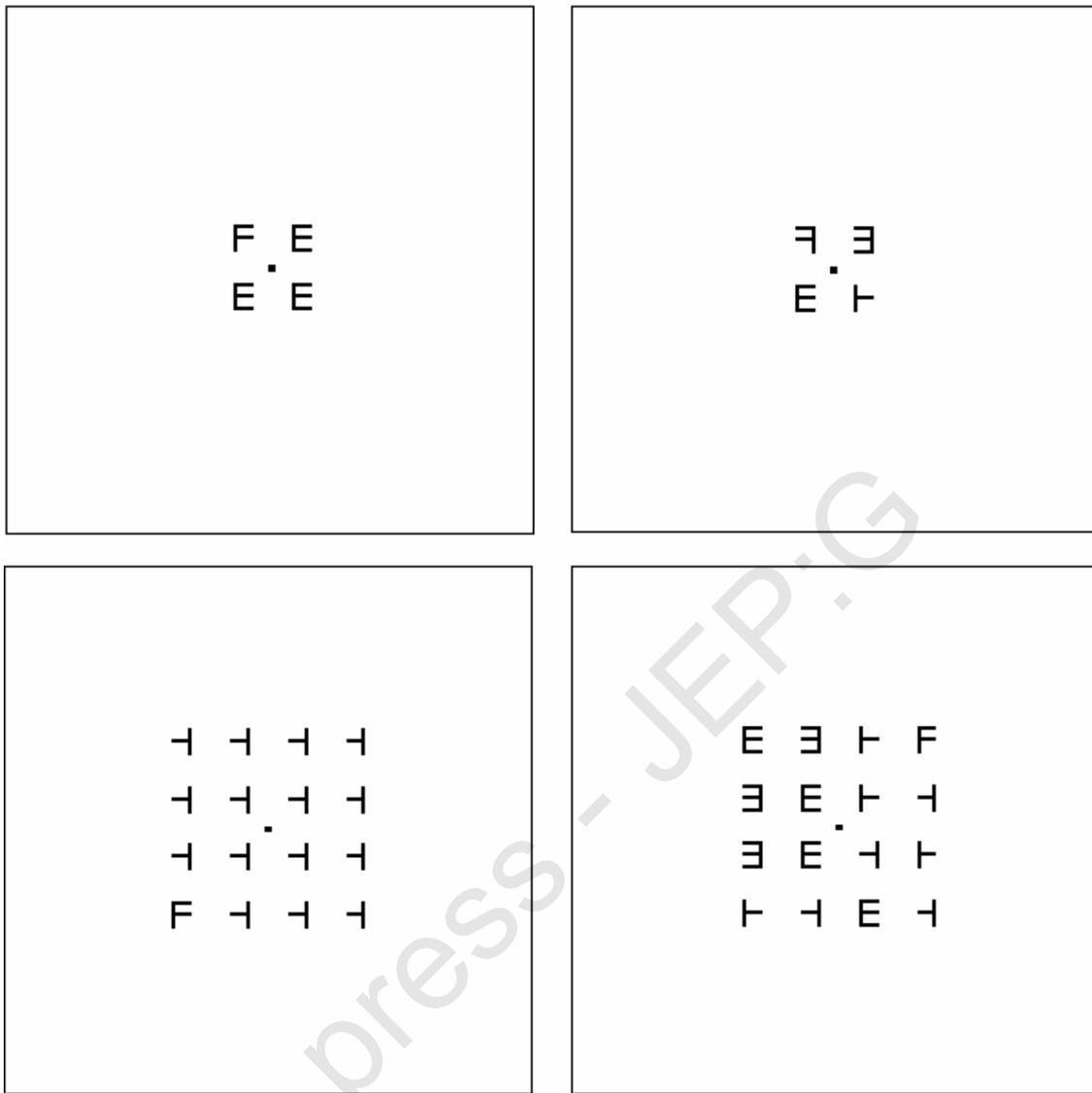
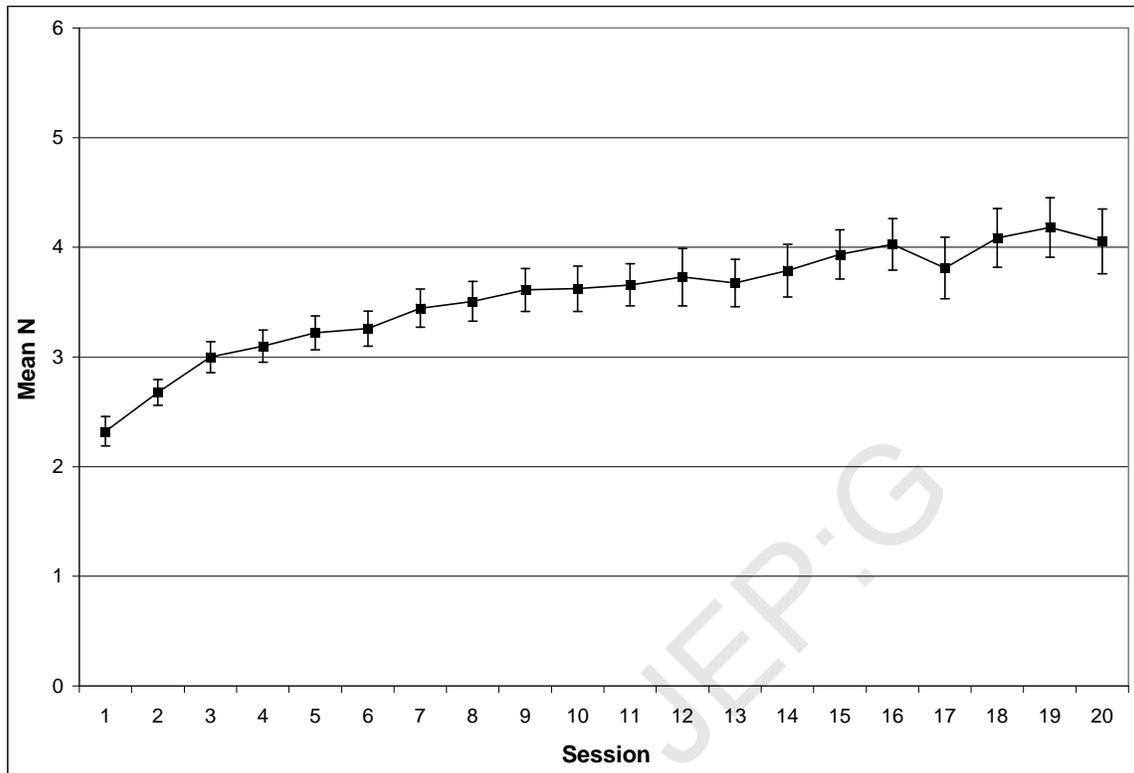


Figure 5. Example stimuli from different levels of the adaptive visual search task. Top-left: 2x2 homogeneous (level 1); Top-right: 2x2 heterogeneous (level 2); Bottom-left: 4x4 homogeneous (level 3); Bottom-right: 4x4 heterogeneous (level 4).

A



B

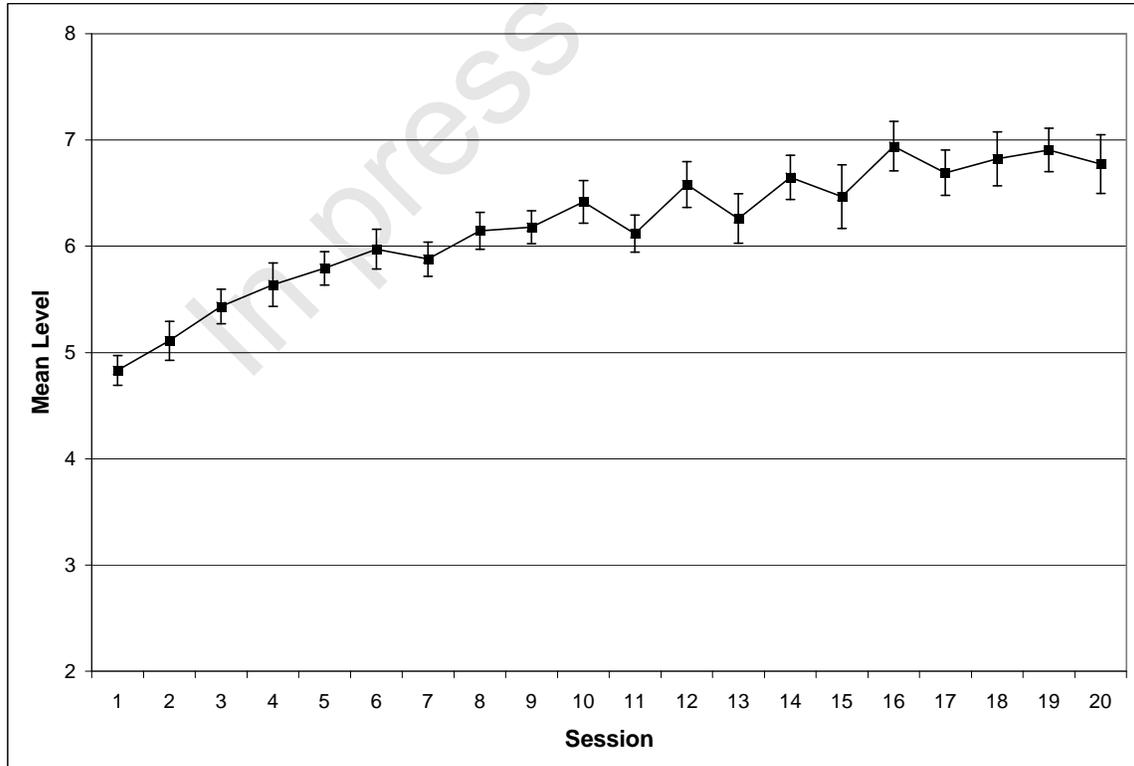


Figure 6. Practice data for the (A) dual n -back and (B) visual search tasks. Error bars ± 1 standard error of the mean.

In press - JEP.G

Supplemental Materials

Table S1

Item number from original tests used in each parallel version for transfer sessions

Raven Advanced Progressive Matrices

A - 1, 6, 7, 12, 13, 18, 19, 24, 25, 30, 31, 36

B - 2, 5, 8, 11, 14, 17, 20, 23, 26, 29, 32, 35

C - 3, 4, 9, 10, 15, 16, 21, 22, 27, 28, 33, 34

Raven Standard Progressive Matrices

A - A1, A6, A7, A12, B2, B5, B8, B11, C3, C4, C9, C10, D1, D6, D7, D12, E2, E5, E8, E11

B - A2, A5, A8, A11, B3, B4, B9, B10, C1, C6, C7, C12, D2, D5, D8, D11, E3, E4, E9, E10

C - A3, A4, A9, A10, B1, B6, B7, B12, C2, C5, C8, C11, D3, D4, D9, D10, E1, E6, E7, E12

Cattell's Culture-Fair

Series Completion

A - P2, 4, 5, 10, 11

B - P3, 3, 6, 9, 12

C - 1, 2, 7, 8, 13

Odd Elements

A - P2, 5, 6, 11, 12

B - 1, 4, 7, 10, 13

C - 2, 3, 8, 9, 14

Matrix Completion

A - P2, 4, 5, 10, 11

B - P3, 3, 6, 9, 12

C - 1, 2, 7, 8, 13

Dot Task

A - P2, 4, 5, 10

B - P3, 3, 6, 9

C - 1, 2, 7, 8

Paper Folding

A - 1, 6, 7, 12, 13, 18

B - 2, 5, 8, 11, 14, 17

C - 3, 4, 9, 10, 15, 16

Letter Sets

A - 1, 6, 7, 12, 13, 18, 19, 24, 25, 30

B - 2, 5, 8, 11, 14, 17, 20, 23, 26, 29

C - 3, 4, 9, 10, 15, 16, 21, 22, 27, 28

Number Series

A - 1, 6, 7, 12, 13

B - 2, 5, 8, 11, 14

C - 3, 4, 9, 10, 15

Inferences

A - 3, 8, 9, 14, 15, 20

B - 4, 7, 10, 13, 16, 19

C - 5, 6, 11, 12, 17, 18

Analogies

A - 1, 6, 7, 12, 13, 18, 19, 24

B - 2, 5, 8, 11, 14, 17, 20, 23

C - 3, 4, 9, 10, 15, 16, 21, 22

General Knowledge

A - 1, 6, 7, 12, 13, 18, 19, 24, 25, 30

B - 2, 5, 8, 11, 14, 17, 20, 23, 26, 29

C - 3, 4, 9, 10, 15, 16, 21, 22, 27, 28

Vocabulary

A - 2, 7, 8, 13, 14, 19, 20, 25, 26, 31, 32, 37, 38

B - 3, 6, 9, 12, 15, 18, 21, 24, 27, 30, 33, 36, 39

C - 4, 5, 10, 11, 16, 17, 22, 23, 28, 29, 34, 35, 40

Table S2
Test orders used for transfer sessions

	A	B	C
1)	Inferences	SynWin	Paper Folding
2)	ControlTower	Number Series	ATClab
3)	Cattell	Letter Comparison	Analogies
4)	Running Letter Span	Number Comparison	General Knowledge
5)	Letter Sets	Raven Standard	Raven Advanced
6)	General Knowledge	Vocabulary	Letter Comparison
7)	Raven Advanced	Analogies	Number Comparison
8)	Letter Comparison	Symmetry Span	Inferences
9)	Number Comparison	Cattell	Running Letter Span
10)	Number Series	ControlTower	Raven Standard
11)	ATClab	Letter Sets	SynWin
12)	Paper Folding	General Knowledge	Letter Sets
13)	Symmetry Span	Raven Advanced	Vocabulary
14)	Analogies	Running Letter Span	Cattell
15)	Vocabulary	Inferences	Symmetry Span
16)	Raven Standard	ATClab	Number Series
17)	SynWin	Paper Folding	Control Tower

Note. All subjects completed the demographic questionnaire before the first test in the pre-test session and the survey questionnaire after the last test in the post-test session.

SynWin [1.0.03]: UNKNOWN, Session 2

Click the List Box to view list

List Box --> **APJZIU**

Probe Letter --> **[]**

YES **NO**

Click YES if Probe letter is in the list, NO if it isn't

Digit Adjustment controls -->

$$\begin{array}{r} 715 \\ + 551 \\ \hline 1266 \end{array}$$

Adjust digits until answer is correct then click DONE

DONE

14

Fuel

50

0 100

Click to refuel

ALERT

Click the ALERT button when you hear the HIGH tone

Figure S1. Example screenshot from SynWin multitask. Probe-recognition, arithmetic, visual monitoring, and auditory monitoring sub-tasks, along with the subjects' total current score, are shown.

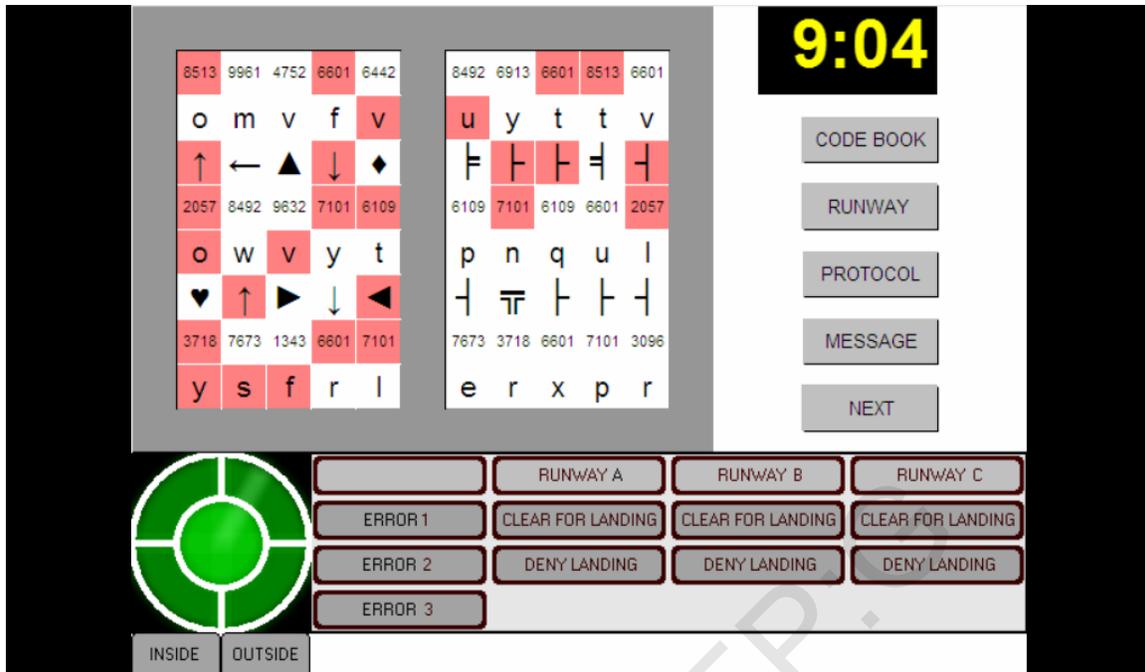


Figure S2. Example screenshot from ControlTower multitask. Primary task comparisons (number, letter, symbols) are shown in the side-by-side arrays. Subjects' responses are made in the right array. Other task components (radar, airplane, color, problem-solving) are considered distractor tasks.

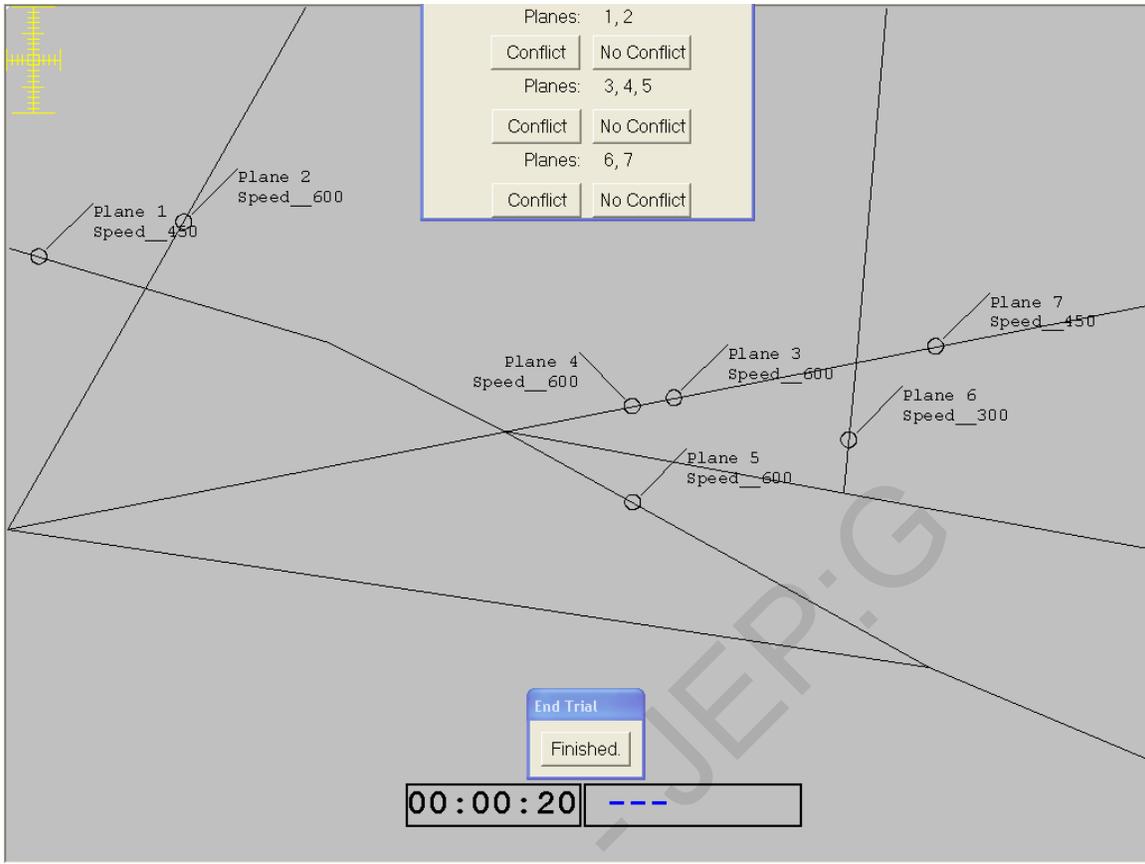


Figure S3. Example screenshot from ATClab multitask.

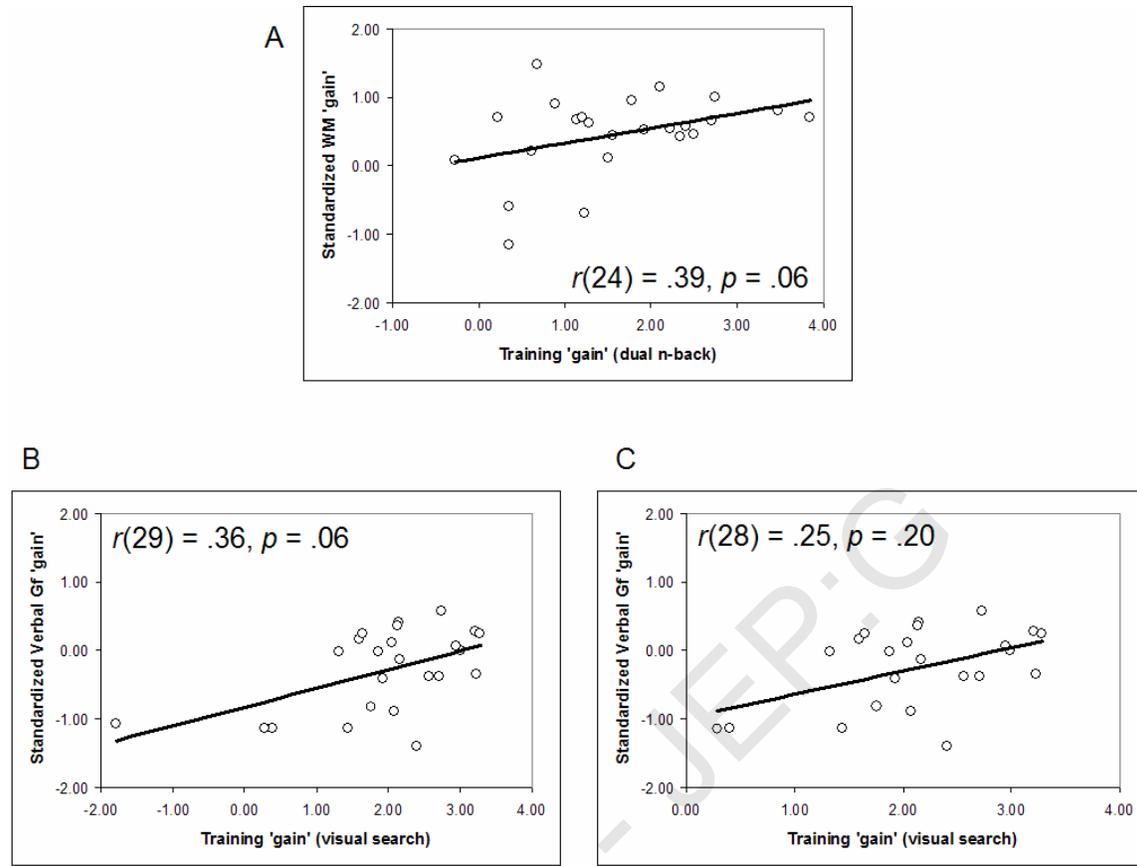


Figure S4. Scatterplots depicting the relationship between training gain and transfer. (A) Amount of dual *n*-back training gain and standardized WM capacity gain. (B) and (C) Amount of visual search training gain and standardized verbal fluid intelligence gain including and excluding one outlier subject, respectively.