

1
2
3
4
5
6
7
8
9
10
11
12
13
14 Performance on the processing portion of complex working memory span tasks is related to
15
16
17 working memory capacity estimates
18
19

20 Lauren L. Richmond¹, Lois K. Burnett¹, Alexandra B. Morrison², B. Hunter Ball³
21
22

23
24 ¹Department of Psychology, Stony Brook University, Stony Brook, NY
25

26 ²Department of Psychology, California State University, Sacramento, Sacramento, CA
27
28

29 ³ Department of Psychology, University of Texas at Arlington, Arlington, TX
30
31
32
33
34
35
36
37
38
39

40 Author Note:
41

42 Lauren L. Richmond and Lois K. Burnett, Department of Psychology, Stony Brook
43
44 University. Alexandra B. Morrison, Department of Psychology, California State University,
45
46 Sacramento. B. Hunter Ball, Department of Psychology, University of Texas at Arlington.
47
48

49 Correspondence concerning this article should be addressed to Lauren L. Richmond,
50
51 Department of Psychology, Stony Brook University; Psychology B Building, Stony Brook, NY
52
53 11794; email: lauren.richmond@stonybrook.edu.
54
55
56
57
58
59
60

Abstract

Individual differences in working memory capacity (WMC) have long been known to relate to performance in domains outside of WM, including attentional control, long-term memory, problem solving, and fluid intelligence to name a few. Complex span WM tasks, composed of a processing component and a storage component, are often used to index WMC in these types of investigations. Capacity estimates are derived from performance on the storage component only, while processing performance is often largely ignored. Here, we explore the relationship between processing performance and WMC in a large dataset for each of three complex span tasks to better characterize how the components of these tasks might be related. We provide evidence that enforcing an 85% or better accuracy criterion for the processing portion of the task results in the removal of a disproportionate number of individuals exhibiting lower WMC estimates. We also find broad support for differences in processing task performance, characterized according to both accuracy and reaction time metrics, as a function of WMC. We suggest that researchers may want to include processing task performance measures, in addition to capacity estimates, in studies using complex span tasks to index WMC. This approach may better characterize the relationships between complex span task performance and performance in disparate domains of cognition.

Performance on the processing portion of complex working memory span tasks is related to working memory capacity estimates

Working memory (WM) is best characterized as a multi-purpose mental workspace (Baddeley & Hitch, 1974; Cowan, 1999) that varies across individuals in terms of capacity (WMC; Kane et al., 2007). Individual differences in WMC have consistently been shown to relate to performance in a wide number of other cognitive domains (e.g., dichotic listening, Conway et al., 2001; attentional control, Kane et al., 2001; stroop interference, Kane & Engle, 2003; reasoning, Kyllonen & Christal, 1990; mind wandering, McVay & Kane, 2009; modes of attentional control, Richmond et al., 2015; fluid intelligence, Unsworth, Brewer, et al., 2009). In this type of research, complex span WM tasks are a popular method of assessing WMC (e.g., Conway et al., 2005; Foster et al., 2015; Redick et al., 2012; Unsworth et al., 2005). In contrast to simple span tasks (e.g., digit span) that involve only a storage component, complex span tasks contain both a processing component that involves the presentation of stimuli requiring decisions and a storage component involving the presentation of memory items (Conway et al., 2005; Unsworth et al., 2005). An important feature of the processing task is that it is typically thought to momentarily reduce access to or block rehearsal of memory items (Unsworth et al., 2005), and performance on the processing task is emphasized in complex span task instructions in order to better capture individual differences in the storage component of the task (Conway et al., 2005).

A Brief History of Complex Span Task Development

Complex span tasks have undergone a number of design implementations since their initial instantiation by Daneman and Carpenter (1980). The initial task structure introduced by

1
2
3 Daneman and Carpenter (1980) was developed to address the relationship between WMC and
4 reading comprehension that had been, to that point, only weakly observed. Previous studies
5
6 used mainly simple span tasks that only taxed memory storage, and Daneman and Carpenter
7
8 (1980) argued that the reason for the previously observed weak correlations was that simple
9
10 span tasks did not appropriately tap the multiple components of the working memory system.
11
12 In this view, simple span tasks are seen as indexing mainly short-term memory abilities. Short-
13
14 term memory, in turn, may be thought of as just one sub-component of working memory,
15
16 which additionally involves attentional control (see Engle et al., 1999 for an extended
17
18 discussion) and other mechanisms such as a controlled search through long-term memory
19
20 (Unsworth & Engle, 2007). Further, the authors argued that this complex span task would be
21
22 both a better measure of WMC and a better measure of individual differences in WMC across
23
24 participants (Daneman & Carpenter, 1980).
25
26
27
28
29
30
31

32
33 Based on the rationale outlined above, Daneman and Carpenter (1980) developed the
34
35 earliest implementation of the now-classic complex span task (originally dubbed “Reading
36
37 Span” and now commonly referred to as “Sentence Span”, as the term “Reading Span” is now
38
39 more typically used to describe the sentence verification/letter memory task; Redick et al.,
40
41 2012). This task required participants to judge the veracity of sentences presented by an
42
43 experimenter on index cards as the processing component, and to remember the final word of
44
45 each sentence as the storage component (Daneman & Carpenter, 1980). Later, Turner and
46
47 Engle (1989) replaced the sentences in Daneman and Carpenter’s task with mathematical
48
49 operations, thus creating the Operation Span task. Turner and Engle (1989) used this task to
50
51
52
53
54
55
56
57
58
59
60

1
2
3 demonstrate that reading comprehension could be predicted with a task that did not require
4
5 reading as the processing component.
6
7

8 In a later publication that influenced the ubiquity of complex span tasks, Unsworth and
9
10 colleagues (2005) developed an automated version of the operation span task (and later other
11
12 complex span tasks) that allowed for automated computerized administration with minimal
13
14 experimenter interaction. These tasks are now widely used in the literature, thanks to both the
15
16 ease of administration afforded by the automated tasks and the ability of researchers to access
17
18 the computerized versions of these tasks that have been made available for download by the
19
20 Engle lab (see <https://englelab.gatech.edu/taskdownloads>). Together, these two features have
21
22 contributed to cementing these tasks as mainstays in the working memory and cognitive
23
24 individual differences literatures.
25
26
27
28

29
30 In recent years, shortened versions of the automated span tasks (Foster et al., 2015;
31
32 Oswald et al., 2015), as well as advanced versions testing performance at larger set sizes
33
34 (Draheim et al., 2018) have been introduced. These new task versions are expected to further
35
36 increase the popularity and utility of complex span WM tasks. The scope of the present work
37
38 will be limited to the “standard” automated versions of the Operation Span, Symmetry Span,
39
40 and Reading Span tasks (Unsworth et al., 2005).
41
42
43
44

45 **The Relationship between Processing and Storage**

46
47 Although these complex span tasks necessarily contain both processing and storage,
48
49 typical methods of calculating WMC estimates rely exclusively on performance in the storage
50
51 component of the task. Adherence to an 85% processing accuracy criterion for inclusion in
52
53 analyses has been recommended to ensure that participants sufficiently engage the processing
54
55
56
57
58
59
60

1
2
3 component (Conway et al., 2005). However, ensuring that participants' processing performance
4
5
6 meets or exceeds this criterion is often the only consideration of this task component for
7
8
9 scoring. As noted above, early methods for testing WMC using complex span tasks were not
10
11 computerized and were completed in the presence of an experimenter. In these contexts,
12
13 experimenters could ensure adequate engagement with the processing task and could help to
14
15 correct participants' erroneous understanding regarding correct completion of the processing
16
17 task. With the advent of more automated computerized methods, these tasks are often
18
19 completed without such stringent oversight of the experimenter and therefore the 85%
20
21 criterion was suggested to ensure appropriate levels of participant engagement with the
22
23 processing portion of the task.
24
25
26

27
28 There is a relatively small body of literature focusing on investigating the relationship
29
30 between processing and storage performance in complex span tasks directly. Such
31
32 investigations have been conducted in healthy young adults (e.g., Engle et al., 1992; Friedman &
33
34 Miyake, 2004; St Clair-Thompson, 2007a, 2007b; Towse et al., 2000; Unsworth et al., 2005;
35
36 Waters & Caplan, 1996) and in typically developing children (e.g., Barrouillet & Camos, 2001;
37
38 Hitch et al., 2001; St Clair-Thompson, 2007b; Towse et al., 1998). Additional work has examined
39
40 the relationship between processing task accuracy and storage performance (Daneman &
41
42 Tardif, 1987; Engle et al., 1992; Lépine et al., 2005; Salthouse et al., 2008; Shah & Miyake, 1996;
43
44 Towse et al., 2000; Turner & Engle, 1989; Waters & Caplan, 1996). In both contexts, results
45
46 have been somewhat mixed, with some reports finding evidence for a relationship between
47
48 better processing performance indices (lower RT, higher accuracy) and better storage
49
50 performance (e.g., St Clair-Thompson, 2007a, 2007b; Waters & Caplan, 1996) in young adult
51
52
53
54
55
56
57
58
59
60

1
2
3 samples. However, evidence for the opposite pattern in terms of processing RT has been
4
5 observed in children (e.g., Towse et al., 1998)¹, as well as findings suggesting no relationship
6
7 between processing and storage performance in younger adults (e.g., Engle et al., 1992; Shah &
8
9 Miyake, 1996; Towse et al., 2000) and in children (Lépine et al., 2005). In sum, the extant
10
11 literature does not yet offer a clear picture of the way(s) in which performance indices on
12
13 processing and storage tasks interrelate. The present work aims to better characterize this
14
15 relationship.
16
17
18

19
20 When introducing the automated version of the Operation Span task (OSpan), Unsworth
21
22 and colleagues (2005) first tested the relationship between facets of processing time and
23
24 storage performance. In the OSpan task, the average RT for processing stimulus presentation
25
26 was found to be negatively related to storage accuracy (Unsworth et al., 2005). Therefore,
27
28 those who were faster on the processing task exhibited higher WMC. Importantly, the sample
29
30 in this study was restricted to participants exhibiting 85% processing accuracy or better,
31
32 resulting in approximately 15% data loss (Unsworth et al., 2005). Following this initial
33
34 examination, Unsworth, Redick, and colleagues (2009) tested the relationship between
35
36 processing and storage performance without enforcing any processing accuracy criterion, as
37
38 well as the relation of these variables to fluid intelligence scores. Unsworth and colleagues
39
40 (2009) found that processing accuracy and processing RT factors were non-redundant with one
41
42 another, and that each of these variables accounted for significant variance in fluid intelligence
43
44 scores over and above that predicted by WMC storage scores. More central to the interests of
45
46
47
48
49
50
51
52
53

54
55 ¹ Previous work has noted that relationships between processing RT and storage performance are more
56
57 consistently observed in children compared to adult samples (Towse et al., 2010).
58
59
60

1
2
3 the current paper, higher processing accuracy and faster processing RTs were associated with
4 better storage performance (i.e., WMC; Unsworth, Redick, et al., 2009). Building on this line of
5 work, Unsworth, Fukuda, and colleagues (2014) tested the relationships between processing
6 and storage performance and the relation of these factors to capacity, secondary memory,
7 attentional control, and fluid intelligence. In this analysis, the negative relationship between
8 processing time and storage accuracy that has been observed in prior work (Unsworth, Redick,
9 et al., 2009) was replicated. Moreover, capacity, secondary memory, and attentional control
10 were shown to fully account for the relationship between the WM indices (processing, storage)
11 and fluid intelligence (Unsworth et al., 2014). Such findings provide initial evidence that
12 consideration of both processing and storage together may be a worthwhile approach to
13 characterize task performance. Importantly, more recent work using complex span tasks to
14 characterize WMC appear to have abandoned strict adherence to an 85% processing accuracy
15 criterion (see, for example, Ellis et al., 2020; McVay & Kane, 2009; Redick et al., 2011; Richmond
16 et al., 2015; Unsworth et al., 2013). Đokić and colleagues (2018) recently suggested that
17 eliminating the 85% accuracy criterion does not impact the psychometric properties of the
18 tasks, and Unsworth, Redick, and colleagues (2009) suggested that enforcing the 85% accuracy
19 criterion is unnecessary. However, previous research has not strongly recommended against
20 enforcing this criterion, nor has a systematic analysis been undertaken to characterize the
21 impact of enforcing this criterion on WMC estimates retained for inclusion in the final sample.

22 **Advances in Investigating Task Reaction Time and Accuracy**

23 The examination of trial-level variation in RTs, rather than characterizing RTs according
24 to their mean, has recently gained traction. One such approach involves application of the ex-

1
2
3 Gaussian model to RT distributions. The ex-Gaussian model convolves the Gaussian and
4
5 exponential distributions together, and are described by the parameters μ , σ , and τ . The μ
6
7 parameter approximates the mode of the Gaussian distribution and the σ parameter
8
9 approximates the standard deviation of the Gaussian distribution, whereas the τ parameter
10
11 reflects the mean and standard deviation of the exponential component of the distribution
12
13 (Balota & Yap, 2011). Because the sum of μ and τ are roughly equal to the mean RT, any
14
15 variable that results in an increase in τ accompanied by a decrease in μ (or vice versa) would
16
17 result in a null effect at the level of mean RT, but can be easily observed with ex-Gaussian RT
18
19 characterization (e.g., Ball & Brewer, 2018; Balota et al., 2008; Spieler et al., 1996). This work
20
21 suggests that the application of the ex-Gaussian model to RT data can therefore reveal effects
22
23 that would be masked by simply characterizing RT according to mean performance.
24
25
26
27
28
29

30 Similarly, characterization of task performance according to error types, rather than or
31
32 in addition to overall accuracy, has gained popularity in recent years (e.g., Giovannetti et al.,
33
34 2008; Scullin et al., 2012, 2020). For example, recent work by Giovannetti and colleagues in the
35
36 domain of naturalistic action execution has shown that omissions in the context of a
37
38 performance-based measure of everyday action are closely related to performance on tests of
39
40 episodic memory, whereas commission errors are more closely associated with deficits in
41
42 executive functioning (Devlin et al., 2014). This suggests that characterizing performance
43
44 according to error types may provide a more detailed analysis of participant performance than
45
46 simple accuracy measures alone, and this profile may map meaningfully to other domains of
47
48 cognition.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 The current work follows these two recent trends to provide a more detailed
4
5 examination of the relationship between processing and storage performance in complex span
6
7 tasks. Here, we consider four novel research questions in each of three complex span tasks.
8
9 First, we address the appropriateness of the aforementioned inclusion criterion by asking (1)
10
11 Does enforcing the recommended 85% processing accuracy criterion for inclusion result in the
12
13 removal of more participants with low WMC estimates compared to higher WMC estimates?
14
15 The next three questions investigate whether and how specific aspects of processing
16
17 performance relate to WMC. (2) Do RT means and standard deviations of RTs from the
18
19 processing practice portion of the task differ as a function of WMC?; (3) Do task-derived
20
21 measures of processing RT, including mean RT and ex-Gaussian parameters, and processing task
22
23 error profiles relate significantly to WMC?; and (4) Does modeling RT-based and error-based
24
25 processing profiles together explain more variance in WMC than consideration of either alone?
26
27
28
29
30
31

32 Method

33
34
35 **Operation Span** (OSpan) and **Symmetry Span** (SymSpan) data were collected between
36
37 2011-2019 at four large state universities: Arizona State University (ASU), California State
38
39 University, Sacramento (CSUS), Stony Brook University (SBU), and Temple University (TU).
40
41
42 **Reading Span** (RSpan) data were collected at ASU only. All data were collected in the context of
43
44 task batteries for large-scale projects. Study procedures were approved by the Institutional
45
46 Review Board of each institution.
47
48

49 Operation Span

50
51
52 In OSpan, participants alternated between solving simple math problems as the
53
54 processing component and remembering letters as the storage component. Participants started
55
56
57
58
59
60

1
2
3 out practicing each portion of the task separately - first practicing the letter memory (storage
4 portion) task, then practicing the math (processing) portion of the task, and last practicing
5 alternating between solving math problems and remembering letters (comparable to test
6 trials). See Figure 1 panel A for a task schematic.
7
8
9
10
11

12
13 -----
14
15 Insert Figure 1 about here
16
17
18 -----
19

20 Letters were displayed for 1000 ms each in all phases of the experiment. At the recall
21 phase, participants were shown a grid displaying 12 possible letters with a box beside each
22 letter. Participants were told to recall the letters in the order they were presented; the chosen
23 letters were displayed at the bottom of the screen. Participants were instructed to use the
24 'clear' button displayed on the screen if they made a mistake and wanted to start over. The
25 blank button was displayed on the screen to mark the position of a forgotten letter, and
26 participants were instructed to click the enter button displayed on the screen when they were
27 ready to submit their response. Participants were given as long as they needed to complete the
28 recall phase in all trials.
29
30
31
32
33
34
35
36
37
38
39
40

41
42 For the math problems, a simple arithmetic problem such as " $(6*0) + 1 = ?$ " was
43 displayed, and participants were told to solve the problem as quickly as possible without
44 sacrificing accuracy. Once participants had an answer in mind, they were instructed to click to
45 advance to the next screen. On this screen, a number is displayed at the top of the screen;
46 displayed below this number is a box marked 'true' and another marked 'false'. If the number
47 shown was the correct response to the math problem, the participant was instructed to choose
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 the 'true' button; else, they were instructed to click 'false'. During the processing practice
4
5 portion, mean and standard deviation (SD) of the RTs for the problem display screen were
6
7 calculated, and then an upper limit bound was set for solving math problems in the test context
8
9 by taking each individual's average RT and adding 2.5 standard deviations to that number
10
11 (Unsworth et al., 2005). This serves as the maximum allowable response time for clicking to
12
13 advance from the problem display to the true/false screen. On trials for which participants did
14
15 not click before the maximum time was reached, the trial was marked as a 'time out' error.
16
17
18 Regardless of RT on the problem screen, participants were given unlimited time to respond on
19
20 the true/false screen and the accuracy of the response was recorded.
21
22
23
24

25 The OSpan task consisted of 15 trials, with three trials each at set size ranging from 3-7.
26
27 Set sizes were presented randomly for each participant. In total, for the test phase, participants
28
29 solved 75 math problems and were shown 75 letters. The capacity score for this task was the
30
31 number of letters recalled in the correct position (with 75 being the maximum possible score).
32
33
34

35 **Reading Span**

36
37 The RSpan task is similar to OSpan, save for differences in processing task demands. In
38
39 RSpan, participants alternated between reading sentences and judging whether they made
40
41 sense as the processing component and remembering letters as the storage component. The
42
43 practice phase proceeded as described above in OSpan. See Figure 1 panel B for a task
44
45 schematic.
46
47
48

49 The storage component of the task was exactly the same as in OSpan, described above.
50
51 The RSpan processing task involved sentence verification. Here, a simple sentence ranging in
52
53 length from 10 to 15 words was displayed, and participants were told to make a judgment as
54
55
56
57
58
59
60

1
2
3 quickly as possible, without sacrificing accuracy, regarding whether the sentence made sense or
4
5 not. “Nonsense” sentences were created by replacing one word in the sentence (e.g. “The
6
7 young pencil kept his eyes closed until he was told to look.”). Once participants had an answer
8
9 in mind, they were instructed to click to advance to the next screen. On this screen, participants
10
11 saw a box marked ‘true’ and another marked ‘false’. If the sentence displayed on the previous
12
13 screen made sense, the participant was instructed that they should choose the ‘true’ button;
14
15 else, they should click ‘false’. During the processing practice portion, mean and SD for the RTs
16
17 on the sentence verification screen were calculated, and then an upper limit bound was set for
18
19 making sentence judgments in the test context by taking each individual’s average RT and
20
21 adding 2.5 standard deviations to that number. This served as the maximum allowable
22
23 response time for clicking to advance from the sentence display to the true/false screen. On
24
25 trials for which participants did not click before the maximum time was reached, the trial was
26
27 terminated and counted as a time out error. Participants were given unlimited time to respond
28
29 on the true/false screen.
30
31
32
33
34
35
36

37 The RSpan task consisted of 15 trials, with three trials each at set sizes ranging from 3-7.
38
39 Set sizes were presented randomly for each participant. In total, for the test phase, participants
40
41 completed 75 sentence judgments and were displayed 75 letters. The capacity score for this
42
43 task was the number of letters recalled in the correct position (with 75 being the maximum
44
45 possible score).
46
47
48

49 **Symmetry Span**

50
51
52 In SymSpan, participants alternated between making symmetry judgments as the
53
54 processing component and remembering highlighted locations as the storage component. The
55
56
57
58
59
60

1
2
3 practice portion of the task proceeded as described above. See Figure 1 panel C for a task
4
5 schematic.

6
7
8 Locations were shown in a 4x4 grid with one square of the grid highlighted in red.
9
10 Locations were displayed for 650 ms each in all phases of the experiment. At the recall phase,
11
12 participants were displayed a blank 4x4 grid. Participants were told to recall the locations in the
13
14 order they were presented by clicking on each location; the chosen locations were numbered in
15
16 the grid. Again, participants had access to 'clear', 'blank', and 'enter' buttons and were given as
17
18 much time as needed to complete the recall phase.
19
20
21

22
23 For the symmetry judgments, participants were displayed an 8x8 black and white grid
24
25 and were asked to make a symmetry judgment about the vertical axis. Participants were
26
27 instructed to solve the symmetry problem as quickly as possible without sacrificing accuracy.
28
29 Once participants had an answer in mind, they were instructed to click to advance to the next
30
31 screen. On this screen, participants were instructed to respond 'true' if the grid displayed on
32
33 the previous screen was symmetrical and 'false' if it was not. The processing task response
34
35 deadline was computed as described above (RT mean + 2.5 SDs) for the processing screen, and
36
37 unlimited time was allowed on the true/false screen. When participants did not submit a
38
39 response on the processing before the maximum time was reached, the trial was terminated
40
41 and counted as a time out error.
42
43
44
45

46
47 The SymSpan task consisted of 12 trials, with three trials each at set size ranging from 2-
48
49 5. Set sizes were presented randomly for each participant. In total, for the test phase,
50
51 participants solved 42 symmetry problems and were shown 42 to-be-remembered locations.
52
53
54
55
56
57
58
59
60

1
2
3 The capacity score for this task was the number of locations recalled in the correct order (with
4
5 42 being the maximum possible score).
6
7

8 **Participants**

9

10 Data were collected primarily from participant pools at each respective institution,
11
12 consisting of undergraduate students enrolled in psychology courses who were participating in
13
14 experiments for course credit and/or payment. Paid participants were also recruited through
15
16 flyer advertisements and word of mouth (i.e., paid participants were not required to be
17
18 registered with the university subject pool in order to participate)².
19
20
21
22

23 Although these data were collected across a number of different sites, study designs
24
25 were relatively similar across sites. In all studies, participants completed sessions that were
26
27 between 1 to 2 hours in length, and task batteries were completed over 1 or 2 sessions. All data
28
29 reported here were collected in the context of larger studies that included a variety of other
30
31 tasks in addition to the WMC measures. Participants were aged at least 18 years and provided
32
33 informed consent for their participation in each study. Sample sizes and descriptive statistics for
34
35 the capacity estimates in each task are displayed in Table 1, separated by site. Descriptive
36
37 statistics for storage performance, processing accuracy, and processing RT are displayed in
38
39 Table 2. Cronbach's alphas for each task for both processing and storage components are
40
41 displayed in Table 3.
42
43
44
45
46
47
48
49
50
51
52
53
54

Insert Tables 1, 2, & 3 about here

55 ² All analyses were also conducted without the inclusion of paid participants and the pattern of results was found
56 to be the same as the results of the full sample reported here.
57
58
59
60

Procedure & Statistical Approach

Significance criterion for all statistical tests was set to the $p \leq .01$ level. Statistical analyses were conducted in R (R Core Team, 2008) using the 'stats' package. Cronbach's alphas were computed using the alpha function from the 'psych' package (Revelle, 2018). Ex-gaussian RT distributional components were calculated using QMPE software (Heathcote et al., 2004) and were imported into R for analysis. Where appropriate, Cohen's d effect sizes were computed using the 'lsr' package (Navarro, 2015). Plots were created in R using the 'ggplot' package (Wickham, 2016). For variables where evidence of non-normality was observed, nonparametric statistical tests were conducted. For brevity, only parametric results are reported, with differences from the non-parametric test results footnoted. Tests of normality and non-parametric test results are provided in full in the Supplemental Materials.

Importantly, popularly used automated complex span tasks (e.g., Unsworth et al., 2005) enforce a response deadline for the processing component of the task (mean RT from processing practice performance + 2.5 standard deviations from that mean). Therefore, we examine RT data from both the practice component (where no response deadline is enforced) in question 2 as well as from the task itself in questions 3 and 4. The response deadline applies only to the screen on which the processing task itself (math problem, symmetry grid, sentence reading) is presented, and not the following screen where a response is input. After participants have been alerted to the response deadline in the context of the task proper, it is possible for savvy participants to then 'game' the system by moving on from the processing task screen to the response screen and then lingering on this screen while continuing to think. Therefore, in

1
2
3 addressing research questions 3 and 4 we used total RTs for both the RT exhibited on the
4 response screen in the processing component as well as RT on the processing screen itself. All
5 analyses focusing on processing time are based on RTs derived from correct trials only.
6
7

8
9
10 To examine processing performance, we characterized errors according to two types.
11 For trials on which participants failed to move on to the response screen before their
12 individualized response deadline, these were counted as ‘time out’ errors (regardless of the
13 response rendered on the response screen). Overtly incorrect responses (i.e., for the operation
14 “ $2 \times 3 + 5$?” and a response screen displaying a value of “12,” choosing the “true” box would
15 count as an error) regardless of participant RTs were characterized as ‘incorrect’ errors.
16
17
18
19
20
21
22
23

24
25 Capacity estimates were examined by awarding credit for each to-be-remembered item
26 recalled in the correct position, summed over the entire task (Unsworth et al., 2005).
27
28
29

30 Results

31
32 The results for each of our substantive research questions are reported below,
33 separated by sub-headings.
34
35

36 37 **Are WMC estimates significantly lower for individuals who do not meet the 85% processing** 38 **criterion cutoff than for those that do?** 39 40

41
42 This question was tested with a two-tailed Welch’s two-sample t-test for unequal
43 variances given the different sample sizes for those who missed versus met/exceeded the
44 criterion cutoff. For OSpan, we observe a significant difference in WMC estimates by processing
45 performance ($t(206.73) = 10.74, p < .001, 95\% \text{ CI } [10.60, 15.37], d = 0.99$). This pattern
46 replicates in both SymSpan ($t(247.82) = 10.75, p < .001, 95\% \text{ CI } [5.67, 8.22], d = 0.94$) and
47 RSpan ($t(163.60) = 7.97, p < .001, 95\% \text{ CI } [9.62, 15.95], d = 0.89$; see Figure 2 for visualization of
48
49
50
51
52
53
54
55
56
57
58
59
60

these data). Together, these results strongly suggest that enforcing a processing accuracy criterion for inclusion in the final dataset results in the removal of a greater number of individuals with low WMC compared to those who achieve high WMC estimates.

 Insert Figure 2 about here

Do individuals who display higher mean RTs and/or more variable RTs (i.e., higher RT SDs) in the practice phase for the processing task also exhibit lower WMC estimates?

We tested these questions using two-tailed Pearson's correlations between (a) an individual's mean RT during the processing practice task and WMC, and (b) an individual's RT SD exhibited during the processing practice task and WMC. We observe a small relationship between mean RT during the processing practice task and WMC that was nonetheless consistent across OSpan and SymSpan. This effect was not observed for RSpan. For OSpan, there is a small but significant negative correlation between mean processing practice RT and WMC ($r(1683) = -0.09, p < .001, 95\% \text{ CI } [-0.13, -0.04]$), but not between RT SDs and WMC ($r(1683) = -0.03, p = .256, 95\% \text{ CI } [-0.08, 0.02]$). This pattern replicates in SymSpan (mean practice RT: $r(1051) = -0.13, p < .001, 95\% \text{ CI } [-0.18, -0.07]$; SD practice RT: $r(1051) = -0.05, p = .142, 95\% \text{ CI } [-0.11, 0.02]$)³. However, deviation from this pattern is observed in RSpan, where we observe non-significant relationships between RT indices and WMC estimates (mean practice RT: $r(1060) = -0.02, p = 0.436, 95\% \text{ CI } [-0.08, 0.04]$; SD practice RT: $r(1060) = -0.07, p = 0.015,$

³ The Spearman correlation for the relationship between SymSpan WMC and processing practice RT SDs reaches significance: $r_s(1051) = -0.10, p < .001$.

1
2
3 95% CI [-0.13, -0.01]). Overall, **as indicated by small correlations**, evidence for differences in RTs
4
5 during the processing practice phase as they relate to WMC is weak (Sawilowsky, 2009), with
6
7 slightly stronger support for mean RT differences by WMC estimates compared to RT SDs from
8
9 the practice portion of the task. **Given the small correlations and the large sample size,**
10
11 **however, the relationship between RTs on the practice processing task and WMC is of little**
12
13 **practical significance.**

14
15
16
17
18 **Do task-derived measures of processing RT, including mean RT and ex-Gaussian parameters,**
19
20 **and processing task error profiles relate significantly to WMC?**

21
22
23 To mirror the strategy employed by Unsworth and colleagues (Unsworth et al., 2005;
24
25 Unsworth, Redick, et al., 2009) in which measures of central tendency were used to
26
27 characterize processing RT, a regression model with mean task-derived processing RT entered
28
29 as a predictor and WMC as the outcome was built separately for each task (OSpan, SymSpan,
30
31 and RSpan). Significant models predicting WMC estimates with mean RT derived from the
32
33 processing task were observed **consistently across all three tasks, though variance explained**
34
35 **was small;** OSpan ($R^2 = .019$, R^2 adjusted = $.018$, $F(1, 1683) = 31.78$, $p < .001$), SymSpan ($R^2 =$
36
37 $.050$, R^2 adjusted = $.049$, $F(1, 1051) = 54.82$, $p < .001$), and RSpan ($R^2 = .007$, R^2 adjusted = $.006$, F
38
39 $(1, 1060) = 7.17$, $p = .008$).

40
41
42
43
44
45 Next, we characterized the distribution of task-derived RTs using an ex-Gaussian
46
47 approach and examined the significance of each predictor in these models as well as overall
48
49 model fit for each task. Vincentile plots were created for each task to examine the overlap
50
51 between predicted and observed values derived from the ex-Gaussian model. Vincentiles were
52
53 created by rank-ordering raw RTs from fastest to slowest for each individual and calculating the
54
55
56
57
58
59
60

1
2
3 mean RT for the first 20% of RTs, the next 20%, and so on. The substantial overlap between
4
5 predicted and observed values for the top third, middle third, and bottom third of participants
6
7 according to WMC depicted in Figure 3 suggests that the ex-Gaussian model provided a good fit
8
9 for our RT data. Vincentile plots for our entire sample and for only those participants who had
10
11 35 or more correct RTs (due to concerns over the appropriateness of the ex-Gaussian model for
12
13 fitting a small number of RTs) can be found in the Supplemental Materials.
14
15
16

17
18 -----
19
20 Insert Figure 3 about here
21
22 -----
23
24

25 Across all three tasks, we find significant overall models which explain between 2.8%
26
27 and 4.9% of the variance in WMC estimates. We also observe some inconsistency in terms of
28
29 the significance of individual ex-Gaussian predictors. In OSpan, we observe an overall significant
30
31 model ($R^2 = .031$, R^2 adjusted = $.030$, $F(3, 1681) = 18.10$, $p < .001$). The σ component emerged
32
33 as the only significant predictor in the model ($\beta = -0.188$, $t(1681) = -4.28$, $p < .001$; μ and τ p
34
35 values $> .028$). Similarly, an overall significant model is observed in RSpan ($R^2 = .031$,
36
37 R^2 adjusted = $.028$, $F(3, 1058) = 11.30$, $p < .001$), and the σ ($\beta = -0.188$, $t(1058) = -5.16$, $p < .001$)
38
39 component again emerged as the only significant predictor in the model (μ and τ p values $>$
40
41 $.10$). In SymSpan, we replicate the overall significance of the model ($R^2 = .052$, R^2 adjusted =
42
43 $.049$, $F(3, 1049) = 19.16$, $p < .001$), and the τ ($\beta = -0.104$, $t(1049) = -3.38$, $p < .001$) component
44
45 emerged as the only significant predictor in the model (μ and σ p values $> .10$). In general, we
46
47 find consistent support for the use of ex-Gaussian analyses to characterize task-derived RTs in
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 complex span tasks, but the contributions of individual components to overall predictive power
4
5 varies across tasks⁴.
6
7

8 To test whether the ex-Gaussian approach provides an advantage over characterizing
9
10 RTs by only mean performance, the models described above were tested against one another.
11
12 In OSpan ($F(2, 1681) = 11.075, p < .001, \Delta AIC = 18.06$) and RSpan ($F(2, 1058) = 13.28, p < .001,$
13 $\Delta AIC = 22.33$) we observed significant improvements in model fit for the ex-Gaussian models
14
15 over mean RT models⁵. There was no difference between models for the SymSpan task ($F(2,$
16 $1049) = 1.31, p = .270, \Delta AIC = 1.37$). **In summary, for two of the three tasks (OSpan, RSpan), we**
17
18 **observed improvements in model fit when using an ex-Gaussian approach compared to the**
19
20 **more common mean RT approach.** Together, these data provide support for characterizing RT
21
22 data in complex span tasks using an ex-Gaussian approach.
23
24
25
26
27
28
29

30 Last, we explored whether characterizing processing errors according to their type
31
32 would predict WMC in each task context by running a series of simultaneous multiple
33
34
35

36 ⁴ It has been suggested that ex-Gaussian approaches may produce poor model fits when the number of trials for
37 inclusion in the model are small (Heathcote et al., 2004). In order to ensure that the inclusion of participants with
38 few correct RT trials were not skewing our results, we assessed the significance of the ex-Gaussian model using
39 only participants with 35 or more correct RTs. Removal of those with very few RTs changed neither the
40 directionality nor the significance of results for OSpan nor RSpan. Following the removal of participants with few
41 RTs in the SymSpan task, the overall model remained significant, but in this model only the μ ($\beta = -0.003, t$
42 $(913) = -3.75, p < .001$) component emerged as a significant predictor. Vincentile plots for datasets
43 including those with only 35 or more correct RT trials can be found in the Supplemental Materials.
44

45 ⁵ To compare the non-nested models reported here, we report ΔAIC to provide information on the best fitting model.
46 AIC (Akaike, 1985, 1973) is a penalized likelihood model that is based on the number of estimated regression
47 parameters. Using the relative AIC values for the candidate models being compared, one can select the best model
48 from the set and determine whether the others provide good estimates of the observed data. Burnham and
49 Anderson (2002) provide some general rules of thumb for estimating the level of empirical support for competing
50 models (compared to the best-fitting model) on the basis of ΔAIC , with ΔAIC between 0 and 2 indicating substantial
51 support, between 4 and 10 with less support, and greater than 10 indicating little to no empirical support. The
52 authors note that these rules of thumb are generally applied to nested models and that guidelines values may be
53 larger for non-nested models. Although there is support for the use of AIC to compare non-nested models (Burnham
54 & Anderson, 2002), we also note that some argue against the use of AIC for selecting between non-nested models
55 (Ripley, 2004).
56
57
58
59
60

1
2
3 regression models entering error types (time out, incorrect) as predictors and WMC as the
4
5 outcome variable. Results were consistent across all three complex span tasks for both overall
6
7 model fit and significance of individual predictors (time out errors and incorrect response
8
9 errors). The overall model was significant for OSpan ($R^2 = .139$, R^2 adjusted = $.138$, $F(2, 1682) =$
10
11 135.40 , $p < .001$), and both error types emerged as significant predictors (time out: $\beta = -0.151$, t
12
13 $(1682) = -6.63$, $p < .001$; incorrect: $\beta = -0.323$, $t(1682) = -14.17$, $p < .001$) in this model. Similar
14
15 results were observed for SymSpan (overall model: $R^2 = .153$, R^2 adjusted = $.151$, $F(2, 1050) =$
16
17 94.69 , $p < .001$; time out: $\beta = -0.191$, $t(1050) = -6.66$, $p < .001$; incorrect: $\beta = -0.313$, $t(1050) = -$
18
19 10.88 , $p < .001$) and RSpan (overall model: $R^2 = .119$, R^2 adjusted = $.118$, $F(2, 1059) = 71.61$, $p <$
20
21 $.001$; time out: $\beta = -0.196$, $t(1059) = -6.75$, $p < .001$; incorrect: $\beta = -0.262$, $t(1059) = -9.04$, $p <$
22
23 $.001$). Across all task contexts, results indicate moderate associations between processing error
24
25 profiles and WMC estimates.
26
27
28
29
30
31

32 Does modeling processing RT and processing errors together explain more variance in WMC 33 than either alone? 34 35

36 RT components are derived only from correct trials, so error types add a non-
37
38 overlapping piece of information about processing performance. Here, we compared model fits
39
40 for regression models including RT components (μ , σ , τ) only (hereafter called the RT
41
42 distribution model) and a model including error types as predictors (hereafter referred to as the
43
44 error model) to a model containing both RT distributional components and error types (the RT
45
46 distribution + error model). In OSpan, the RT distribution + error model explained significantly
47
48 more variance in WMC estimates compared to the RT distribution model (model comparison: F
49
50 $(2, 1679) = 130.49$, $p < .001$, $\Delta R^2 = 0.130$) and the error model (model comparison: $F(3, 1679) =$
51
52
53
54
55
56
57
58
59
60

1
2
3 15.34, $p < .001$, $\Delta R^2 = 0.023$). This pattern is replicated in SymSpan (comparison to the RT
4 distribution model: ($F(2, 1047) = 95.24$, $p < .001$, $\Delta R^2 = 0.146$); comparison to the error
5 model: ($F(3, 1047) = 19.61$, $p < .001$, $\Delta R^2 = 0.045$), and RSpan (comparison to the RT distribution
6 model: ($F(2, 1056) = 59.91$, $p < .001$, $\Delta R^2 = 0.099$); comparison to the error model: ($F(3, 1056) =$
7 4.31 , $p = .005$, $\Delta R^2 = 0.011$)⁶. While the increase in explanatory power as a result of using the RT
8 distribution + error model varied across tasks, using the combined ex-Gaussian parameters and
9 error types as predictors provided notable improvements in terms of variance explained.
10 Compared to the RT distribution (only) model, the RT distribution + error model explained
11 between 4.5% and 14.6% of the variance in WMC estimates. Compared to the error (only)
12 model, the RT distribution + error model explained between 1% and 2.3% of the variance in
13 WMC estimates. Across all task contexts, the full model provided a significantly improved
14 model fit compared to models using either the ex-Gaussian RT distribution components or error
15 profiles alone. Taken together, these data provide support for the joint use of RT-based and
16 error-based metrics to characterize processing task performance in relation to WMC estimates.

Discussion

17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40 For our substantive research questions, we observed a fair degree of consistency across
41 the three complex span tasks, as well as some points of departure. Below, we discuss each in
42 terms of consistency across tasks and size of the observed effects.

43
44
45
46
47 First, we see strong evidence in all three tasks that WMC estimates are significantly
48 lower for those that miss the 85% processing task accuracy cutoff than those who meet or
49

50
51
52
53
54 ⁶ The comparison for the quantile regression RT distribution + error model to quantile regression error-only model
55 in RSpan does not meet our $p < .01$ significance criterion (model comparison significance: $p = .035$).
56
57
58
59
60

1
2
3 exceed the cutoff. Notably, the intention of this cutoff is to remove those who do not
4
5 sufficiently engage the processing task (Unsworth et al., 2005), and was not devised to
6
7 eliminate truly low performers. In the current datasets, enforcement of this criterion results in
8
9 the removal of 10-17% of the full sample, and in Unsworth et al., 2005 approximately 15% data
10
11 loss is reported when this criterion is enforced. Our findings suggest that enforcing this criterion
12
13 will result in the removal of a disproportionate number of individuals exhibiting lower WMC
14
15 estimates.
16
17
18

19
20 We observed some inconsistency across tasks for our second question. There were no
21
22 significant associations observed between processing task RT SDs and WMC. Significant
23
24 negative correlations between mean processing practice RT and WMC were observed in two of
25
26 the three tasks (OSpan, SymSpan). Inconsistency by task is perhaps not surprising, given the
27
28 small number of trials and the relatively simple nature of processing task practice. Moreover,
29
30 these findings should be interpreted with caution as, in all cases, the strength of the correlation
31
32 value was low.
33
34
35

36
37 Turning to our third question, in general we find significant relationships between RT
38
39 (characterized by mean RT and ex-Gaussian components) and WMC. For both OSpan and
40
41 RSpan, the ex-Gaussian model fit significantly better than the mean RT model. There were no
42
43 differences for SymSpan. This suggests ex-Gaussian analyses are no worse than measures of
44
45 central tendency for characterizing RTs, and are actually more informative for two of the three
46
47 tasks. Under the ex-Gaussian analyses, the specific components that were significant differed
48
49 across tasks (σ for OSpan and RSpan, τ for SymSpan when all participants were included for
50
51 analysis and μ when only participants with 35 or more correct RTs were retained for modeling).
52
53
54
55
56
57
58
59
60

1
2
3 It is somewhat surprising that the τ component was only found to be a significant predictor in
4
5 one of the three tasks. However, we note that trials reaching the response deadline were
6
7 characterized as time-out errors, so the extent to which effects of τ could be observed may
8
9 have been limited in this context. In support of this view, we saw that the significant
10
11 contribution of τ to the overall model for SymSpan was eliminated when errors were modeled
12
13 together with ex-Gaussian parameters, suggesting that τ and error profiles may be accounting
14
15 for overlapping variance in WMC. It is possible that eliminating or extending the response
16
17 deadline for each trial may lead to more consistency in terms of the contribution of the τ
18
19 component of the ex-Gaussian parameters by capturing responses that exceed the cut-off to be
20
21 classified as time-out errors. However, researchers may be reluctant to alter the standard RT
22
23 deadline of 2.5 times the mean as it is intended to prevent participants from rehearsing items
24
25 from the storage portion of the task when they should be completing the processing
26
27 component (Unsworth et al., 2005).
28
29
30
31
32
33
34

35 Future work using other analytic techniques for characterizing RTs, such as diffusion
36
37 modeling (Ratcliff, 1978), may further elucidate these relationships by offering a
38
39 straightforward link between parameter estimates and cognitive processes. Unfortunately,
40
41 concerns over an inadequate number of RTs available in the current tasks precluded the
42
43 inclusion of this analysis here (Lerche et al., 2017). Researchers who wish to pursue such a
44
45 characterization in future studies are encouraged to explore methods by which to increase the
46
47 number of RTs available in the context of WMC tasks for the application of this model, perhaps
48
49 by utilizing the advanced complex span tasks with larger set sizes (Draheim et al., 2018),
50
51 increasing the number of cycles through each set size, or by loosening or removing the response
52
53
54
55
56
57
58
59
60

1
2
3 deadline for the processing task in order to obtain usable RT data for long response trials (that are
4 currently captured instead as time-out errors under standard task conditions). Nonetheless, across
5
6 all tasks, we observed predictive power for processing error profiles, with both time-out errors
7
8 and incorrect errors emerging as significant predictors of WMC.
9
10

11
12 Last, in all tasks, results provide consistent evidence that the RT distribution + error
13
14 models provided better fits than models containing only error information and models
15
16 containing only RT distributional information. In terms of variance explained, the RT distribution
17
18 + error model increased variance explained by 4.5% to 14.6% compared to the RT distribution
19
20 model. Increases in variance explained compared to the error model were smaller (between ~1-
21
22 2%). While the observed effects are small for some of the tasks, inclusion of both ex-Gaussian
23
24 RT components and error profiles together consistently improves the explanatory power of the
25
26 models. Future research may consider the specific relationships between each processing
27
28 component (RT, accuracy) and other tasks sensitive to cognitive individual differences. For
29
30 instance, it is possible that ex-Gaussian RT components may be more related to tasks that
31
32 emphasize speeded responding whereas processing accuracy may be more related to tasks that
33
34 don't require speeded responses (Unsworth, Redick, et al., 2009).
35
36
37
38
39
40
41

42 The current work dovetails with prior work regarding individual differences in WM and
43
44 extends this work in some important ways. Previous work has sought to characterize individual
45
46 differences in WMC and the relation of those differences to other cognitive factors of interest,
47
48 including processing speed (Conway et al., 2002). WMC and processing speed are generally
49
50 found to be only weakly related to one another in samples of healthy adults (see e.g., Conway
51
52 et al., 2002). WMC is often observed as a better predictor of higher-order cognitive functions
53
54
55
56
57
58
59
60

1
2
3 such as reasoning (Kyllonen & Christal, 1990) and fluid intelligence (Conway et al., 2002)
4
5 compared to processing speed measures in healthy adult samples, despite in some cases
6
7 observing strong relationships between processing speed and WMC (Kyllonen & Christal, 1990,
8
9 but see Conway et al., 2002). The relationship between processing speed and WMC is more
10
11 strongly observed in early development (Kail, 2007) and in late life (Brown et al., 2012).
12
13
14 Nonetheless, consideration of RT-based metrics that can be derived from complex span tasks
15
16 themselves may provide a more fruitful way to characterize relationships between WMC and
17
18 speed. In this regard, future work to compare the strength of the relationships observed
19
20 between WMC-storage measures derived from complex span tasks and RT metrics derived from
21
22 traditional processing speed tasks and from the processing portion of the complex span WM
23
24 task itself would be useful.
25
26
27
28
29

30 Reflecting on the present findings, we make some recommendations for future WM
31
32 research using these tasks. First, we caution against use of an 85% processing accuracy criterion
33
34 as it may inadvertently bias WMC estimates against lower capacity participants, and if
35
36 adherence to the cutoff is maintained researchers should be aware that this is likely to result in
37
38 skewed estimates of WMC in their samples. In considering whether to forgo enforcing this
39
40 criterion, researchers should assess whether their sample sizes are sufficiently powered to
41
42 tolerate approximately 10-17% data loss. Researchers may instead choose to adopt a criterion
43
44 closer to 50% (where performance worse than 50% likely represents misunderstanding of or
45
46 insufficient engagement with the processing task) or to embed attention checks in the task in
47
48 order to justify the removal of participants' data. **If researchers wanted to move beyond the
49
50 simple 50% accuracy criterion for processing task performance and instead include a cutoff for
51
52
53
54
55
56
57
58
59
60**

1
2
3 ensuring that participants were included in the final were 95% or 99% likely to be above
4
5 guessing probability, this could be easily achieved and data below the cut score could be
6
7 discarded on either a task-wise or trial-wise basis. At the task level, to achieve 95% confidence
8
9 that participants weren't guessing on a task with 42 processing trials (SymSpan) the criterion
10
11 should be set to 61.9% overall processing accuracy. For a task with 75 processing steps (Ospan,
12
13 RSpan), a cut score below 69% accuracy should be adopted under the 95% confidence
14
15
16
17
18 criterion⁷.

19
20 Alternatively, researchers may choose to adopt a data-driven approach to set a
21
22 processing task performance threshold for inclusion in the final dataset based on their own
23
24 sample. For example, participants who are found to exhibit processing task performance 2 or
25
26 2.5 SDs below the mean processing performance in that sample could be excluded from the
27
28 final dataset. In cases where participants complete more than 1 complex span task, z-scored
29
30 processing task performance could be computed for each participant across tasks, similar to the
31
32 typical approach for combining WMC scores across tasks (see Morrison & Richmond, 2020;
33
34 Redick et al., 2011; Richmond et al., 2015; Shipstead & Broadway, 2013 for examples of this
35
36 approach). From here, participants exhibiting z-scores equal to or less than -2 or -2.5 could be
37
38 removed from the final sample. This approach could perhaps be used in combination with a
39
40 criterion for acceptable lower-bound RTs displayed on the processing task screen itself (e.g.,
41
42 RTs shorter than 200 ms). These approaches may be particularly useful in samples that exhibit
43
44 higher average performance on complex span tasks (see for example Redick et al., 2012 Table 5
45
46 showing differences in WMC estimates by data collection site).

47
48
49
50
51
52
53
54
55
56 ⁷ Thank you to an anonymous reviewer for suggesting this approach.
57
58
59
60

1
2
3 At the same time, we acknowledge that simply reducing or abandoning a criterion for
4
5 inclusion on the processing task is not expected to fully eliminate issues regarding inclusion of
6
7 problematic data. Instead, this approach is expected to improve the retention of data for
8
9 engaged participants at the lower end of the WMC spectrum. In other words, it is possible that
10
11 enforcing a less stringent processing accuracy criterion level could allow for the inclusion of a
12
13 small number of participants who fail to adequately engage with the processing task (i.e., faux
14
15 lows) in the final dataset. More importantly, forgoing strict adherence to the 85% accuracy
16
17 criterion is expected to allow for the inclusion of an interesting and important sub-set of
18
19 individuals with lower working memory capacity estimates. Approaches such as those described
20
21 above (enforcing an accuracy criterion at chance levels, using a data driven approach to setting
22
23 a criterion for inclusion, setting a lower-bound of acceptable processing RTs for inclusion) are
24
25 expected to minimize the number of participants included who truly failed to engage with the
26
27 processing task, but also maximize the number of participants retained for analysis who were
28
29 engaged in the task and simply struggled with both processing and storage components of the
30
31 task at hand. Moreover, we note that the traditional criterion for inclusion only applies to the
32
33 processing portion of the task, and there is no lower bound accuracy metric enforced for
34
35 performance on the storage portion. To separate out truly disengaged participants from
36
37 engaged participants with lower WMC, setting a criterion for acceptable storage scores (e.g., at
38
39 least one correct trial at the lowest set size), in addition to adopting some or all of the practices
40
41 we recommend above, could prove useful in future studies.

42
43 Another important avenue for future research will be to examine if the current findings
44
45 extend to different complex span task configurations and design choices. For example, in the
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 complex span tasks used in the current research, the processing task itself was designed to
4
5 disrupt active maintenance of information in working memory through rehearsal (Conway et
6
7 al., 2005). Previous work examining complex span tasks including semantically related
8
9 information for both the processing and storage components (Towse et al., 2010) suggests that
10
11 the reinstatement of context at processing could serve to boost, rather than disrupt, storage
12
13 performance (see e.g., Delaney & Sahakyan, 2007; Wahlheim et al., 2016, 2017; Wahlheim &
14
15 Huff, 2015; Wingfield & Kahana, 2002 for further discussions of contextual reinstatement and
16
17 memory performance). It is as yet unknown whether the patterns observed in the current study
18
19 would extend to conditions where the processing task is designed to support, rather than
20
21 disrupt, access to the to-be-remembered information. In considering other task configurations,
22
23 it has recently been suggested that the order of the storage and processing components can
24
25 impact estimates of WMC, with the processing-storage sequence resulting in higher estimates
26
27 of WMC compared to storage-processing (Debraise et al., 2020). Future work may explore
28
29 whether the relationships between processing and storage task components observed in the
30
31 current work (using the standard processing-storage sequence) would extend to complex span
32
33 tasks with a storage-processing sequence.
34
35
36
37
38
39
40

41
42 Finally, we also encourage those using complex span tasks in their own work to examine
43
44 processing performance more thoroughly. The advantages of this approach are twofold: (1)
45
46 more efficient use of collected data, and (2) task length could potentially be reduced while still
47
48 obtaining stable estimates of WMC. Large-scale individual differences studies should consider
49
50 processing RT, processing error types, and storage together in relation to other cognitive
51
52 constructs (see also Unsworth, Redick, et al., 2009). Inclusion of both storage and processing
53
54
55
56
57
58
59
60

1
2
3 components is expected to be informative and may reveal subtle relationships between
4
5 component processes embedded in complex span tasks and other cognitive domains. Overall,
6
7 the present work provides strong support for careful consideration of processing performance
8
9 indices, in addition to storage performance, in the context of complex span tasks.
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only

Open Practices Statement

Analysis scripts, aggregate data files, PDF outputs, and Supplemental Materials are available at https://osf.io/85yfe/?view_only=204c6d7d4ba24b1d80ce145675960f21. [*Please note: This link will be replaced with a publicly accessible link after acceptance.*] The analyses reported in this paper were not preregistered.

For Review Only

References

- 1
2
3
4
5
6 Akaike, H. (1985). Prediction and Entropy. In E. Parzen, K. Tanabe, & G. Kitagawa (Eds.), *Selected*
7
8 *Papers of Hirotugu Akaike* (pp. 387–410). Springer. <https://doi.org/10.1007/978-1-4612->
9
10 1694-0_30
11
12
13 Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In
14
15 B. N. Petrov & F. Csaki (Eds.), *Proceedings of the 2nd International Symposium on*
16
17 *Information Theory* (pp. 267–281). Akadémiai Kiadó.
18
19
20 Baddeley, A., & Hitch, G. (1974). Working memory. In G. Bower (Ed.), *Recent advances in*
21
22 *learning and motivation* (pp. 47–89). Academic Press.
23
24
25 Ball, B. H., & Brewer, G. A. (2018). Proactive control processes in event-based prospective
26
27 memory: Evidence from intraindividual variability and ex-Gaussian analyses. *Journal of*
28
29 *Experimental Psychology: Learning, Memory, and Cognition*, 44(5), 793–811.
30
31 <https://doi.org/10.1037/xlm0000489>
32
33
34
35 Balota, D. A., & Yap, M. J. (2011). Moving beyond the mean in studies of mental chronometry:
36
37 The power of response time distributional analyses. *Current Directions in Psychological*
38
39 *Science*, 20(3), 160–166. <https://doi.org/10.1177/0963721411408885>
40
41
42 Balota, D. A., Yap, M. J., Cortese, M. J., & Watson, J. M. (2008). Beyond mean response latency:
43
44 Response time distributional analyses of semantic priming. *Journal of Memory and*
45
46 *Language*, 59(4), 495–523. <https://doi.org/10.1016/j.jml.2007.10.004>
47
48
49 Barrouillet, P., & Camos, V. (2001). Developmental increase in working memory span: Resource
50
51 sharing or temporal decay? *Journal of Memory and Language*, 45(1), 1–20.
52
53 <https://doi.org/10.1006/jmla.2001.2767>
54
55
56
57
58
59
60

- 1
2
3 Brown, L. A., Brockmole, J. R., Gow, A. J., & Deary, I. J. (2012). Processing speed and visuospatial
4 executive function predict visual working memory ability in older adults. *Experimental*
5
6 *Aging Research, 38*(1), 1–19. <https://doi.org/10.1080/0361073X.2012.636722>
7
8
9
- 10 Burnham, K. P., & Anderson, D. R. (2002). A practical information-theoretic approach. *Model*
11
12 *Selection and Multimodel Inference, 2*.
13
- 14
15 Conway, A., Cowan, N., & Bunting, M. (2001). The cocktail party phenomenon revisited: The
16
17 importance of working memory capacity. *Psychonomic Bulletin & Review, 8*, 331–335.
18
19
- 20 Conway, A., Cowan, N., Bunting, M. F., Theriault, D. J., & Minkoff, S. R. (2002). A latent variable
21
22 analysis of working memory capacity, short-term memory capacity, processing speed,
23
24 and general fluid intelligence. *Intelligence, 30*(2), 163–183.
25
26
- 27
28 Conway, A., Kane, M. J., Bunting, M., Hambrick, D., Wilhelm, O., & Engle, R. (2005). Working
29
30 memory span tasks: A methodological review and user's guide. *Psychonomic Bulletin &*
31
32 *Review, 12*(5), 769–786.
33
34
- 35
36 Cowan, N. (1999). An embedded process model of working memory. In A. Miyake & P. Shah
37
38 (Eds.), *Models of working memory: Mechanisms of active maintenance and executive*
39
40 *control* (pp. 62–101). Cambridge University Press.
41
- 42
43 Daneman, M., & Carpenter, P. (1980). Individual differences in working memory and reading.
44
45 *Verbal Learning and Verbal Behavior, 19*(4), 450–466.
46
- 47
48 Daneman, M., & Tardif, T. (1987). Working memory and reading skill re-examined. In *Attention*
49
50 *and performance 12: The psychology of reading*. (pp. 491–508). Lawrence Erlbaum
51
52 Associates, Inc.
53
54
55
56
57
58
59
60

1
2
3 Delaney, P., & Sahakyan, L. (2007). Unexpected costs of high working memory capacity
4 following directed forgetting and contextual change manipulations. *Memory &*
5
6 *Cognition*, 35(5), 1074–1082.
7
8
9

10 Devlin, K. N., Giovannetti, T., Kessler, R. K., & Fanning, M. J. (2014). Commissions and omissions
11 are dissociable aspects of everyday action impairment in schizophrenia. *Journal of the*
12
13 *International Neuropsychological Society*, 20(08), 812–821.
14
15

16 <https://doi.org/10.1017/S1355617714000654>
17
18

19 Đokić, R., Koso-Drljević, M., & Đapo, N. (2018). Working memory span tasks: Group
20 administration and omitting accuracy criterion do not change metric characteristics.
21
22 *PLOS ONE*, 13(10), e0205169. <https://doi.org/10.1371/journal.pone.0205169>
23
24
25

26 Draheim, C., Harrison, T. L., Embretson, S. E., & Engle, R. W. (2018). What item response theory
27 can tell us about the complex span tasks. *Psychological Assessment*, 30(1), 116–129.
28
29

30 <https://doi.org/10.1037/pas0000444>
31
32
33

34 Ellis, D. M., Ball, B. H., Kimpton, N., & Brewer, G. A. (2020). The role of working memory
35 capacity in analytic and multiply-constrained problem-solving in demanding situations.
36
37 *Quarterly Journal of Experimental Psychology*, 73(6), 920–928.
38
39

40 <https://doi.org/10.1177/1747021820909703>
41
42
43

44 Engle, R. W., Cantor, J., & Carullo, J. J. (1992). Individual differences in working memory and
45 comprehension: A test of four hypotheses. *Journal of Experimental Psychology:*
46
47

48 *Learning, Memory, and Cognition*, 18(5), 972–992. [https://doi.org/10.1037/0278-](https://doi.org/10.1037/0278-7393.18.5.972)
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 Engle, R. W., Tuholski, S. W., Laughlin, J. E., & Conway, A. (1999). Working memory, short-term
4 memory, and general fluid intelligence: A latent-variable approach. *Journal of*
5
6 *Experimental Psychology: General*, 128(3), 309.
7
8
9
- 10 Foster, J. L., Shipstead, Z., Harrison, T. L., Hicks, K. L., Redick, T. S., & Engle, R. W. (2015).
11
12 Shortened complex span tasks can reliably measure working memory capacity. *Memory*
13
14 & *Cognition*, 43(2), 226–236. <https://doi.org/10.3758/s13421-014-0461-7>
15
16
17
- 18 Friedman, N. P., & Miyake, A. (2004). The reading span test and its predictive power for reading
19
20 comprehension ability. *Journal of Memory and Language*, 51(1), 136–158.
21
22
23 <https://doi.org/10.1016/j.jml.2004.03.008>
24
- 25 Giovannetti, T., Bettcher, B. M., Brennan, L., Libron, D. J., Kessler, R. K., & Duey, K. (2008).
26
27 Coffee with jelly or unbuttered toast: Commissions and omissions are dissociable
28
29 aspects of everyday action impairment in Alzheimer's disease. *Neuropsychology*, 22(2),
30
31
32 235–245. <https://doi.org/10.1037/0894-4105.22.2.235>
33
34
- 35 Heathcote, A., Brown, S., & Cousineau, D. (2004). QMPE: Estimating Lognormal, Wald, and
36
37 Weibull RT distributions with a parameter-dependent lower bound. *Behavior Research*
38
39 *Methods, Instruments, & Computers*, 36(2), 277–290.
40
41
42 <https://doi.org/10.3758/BF03195574>
43
44
- 45 Hitch, G. J., Towse, J. N., & Hutton, U. (2001). What limits children's working memory span?
46
47 Theoretical accounts and applications for scholastic development. *Journal of*
48
49 *Experimental Psychology: General*, 130(2), 184–198. <https://doi.org/10.1037//0096->
50
51
52 3445.130.2.184
53
54
55
56
57
58
59
60

- 1
2
3 Kail, R. V. (2007). Longitudinal evidence that increases in processing speed and working
4
5 memory enhance children's reasoning. *Psychological Science*, *18*(4), 312.
6
7
- 8 Kane, M. J., Bleckley, M., Conway, A., & Engle, R. (2001). A controlled-attention view of working
9
10 memory capacity. *Journal of Experimental Psychology: General*, *130*(2), 169–183.
11
12
- 13 Kane, M. J., Conway, A., Hambrick, D., & Engle, R. (2007). Variation in working memory capacity
14
15 as variation in executive attention and control. In C. Jarrold & A. Conway (Eds.),
16
17 *Variation in Working Memory* (pp. 21–48). Oxford University Press.
18
19
- 20 Kane, M. J., & Engle, R. (2003). Working-memory capacity and the control of attention: The
21
22 contributions of goal maintenance, response competition, and task set to Stroop
23
24 interference. *Journal of Experimental Psychology: General*, *132*, 47–70.
25
26
- 27 Kyllonen, P., & Christal, R. (1990). Reasoning ability is (little more than) working-memory
28
29 capacity?! *Intelligence*, *14*, 383–433.
30
31
- 32 L epine, R. Ile, Parrouillet, P., & Camos, V. (2005). What makes working memory spans so
33
34 predictive of high-level cognition? *Psychonomic Bulletin & Review*, *12*(1), 165–170.
35
36
37 <https://doi.org/10.3758/BF03196363>
38
39
- 40 Lerche, V., Voss, A., & Nagler, M. (2017). How many trials are required for parameter
41
42 estimation in diffusion modeling? A comparison of different optimization criteria.
43
44
45 *Behavior Research Methods*, *49*(2), 513–537. [https://doi.org/10.3758/s13428-016-0740-](https://doi.org/10.3758/s13428-016-0740-2)
46
47 2
48
- 49 McVay, J., & Kane, M. J. (2009). Conducting the train of thought: Working memory capacity,
50
51 goal neglect and mind wandering in an executive-control task. *Journal of Experimental*
52
53 *Psychology: Learning, Memory and Cognition*, *35*(1), 196–294.
54
55
56
57
58
59
60

- 1
2
3 Morrison, A. B., & Richmond, L. L. (2020). Offloading items from memory: Individual differences
4
5 in cognitive offloading in a short-term memory task. *Cognitive Research: Principles and*
6
7 *Implications*, 5(1). <https://doi.org/10.1186/s41235-019-0201-4>
8
9
- 10 Navarro, D. (2015). *Learning statistics with R: A tutorial for psychology students and other*
11
12 *beginners. (Version 0.5)*. University of Adelaide. <http://ua.edu.au/ccs/teaching/lr>
13
14
- 15 Oswald, F. L., McAbee, S. T., Redick, T. S., & Hambrick, D. Z. (2015). The development of a short
16
17 domain-general measure of working memory capacity. *Behavior Research Methods*,
18
19 47(4), 1343–1355. <https://doi.org/10.3758/s13428-014-0543-2>
20
21
- 22 R Core Team. (2008). *R: A Language and Environment for Statistical Computing*. R Foundation
23
24 for Statistical Computing. <http://www.R-project.org>
25
26
- 27 Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, 85(2), 59.
28
29
- 30 Redick, T. S., Broadway, J. M., Meier, M. E., Kuriakose, P. S., Unsworth, N., Kane, M. J., & Engle,
31
32 R. W. (2012). Measuring working memory capacity with automated complex span tasks.
33
34 *European Journal of Psychological Assessment*, 28(3), 164–171.
35
36
- 37 Redick, T. S., Calvo, A., Gay, C. E., & Engle, R. W. (2011). Working memory capacity and go/no-
38
39 go task performance: Selective effects of updating, maintenance, and inhibition. *Journal*
40
41 *of Experimental Psychology: Learning, Memory and Cognition*, 37(2), 308–324.
42
43
- 44 Revelle, W. (2018). *psych: Procedures for Psychological, Psychometric, and Personality*
45
46 *Research*. Northwestern University. <https://CRAN.R-project.org/package=psych>
47
48
- 49 Richmond, L. L., Redick, T. S., & Braver, T. S. (2015). Remembering to prepare: The benefits (and
50
51 costs) of high working memory capacity. *Journal of Experimental Psychology: Learning,*
52
53 *Memory, and Cognition*, 41(6), 1764–1777. <https://doi.org/10.1037/xlm0000122>
54
55
56
57
58
59
60

- 1
2
3 Ripley, B. D. (2004). Selecting amongst large classes of models. In *Methods and Models in*
4
5 *Statistics: In Honour of Professor John Nelder, FRS* (pp. 155–170). World Scientific.
6
7
8 Salthouse, T. A., Pink, J. E., & Tucker-Drob, E. M. (2008). Contextual analysis of fluid intelligence.
9
10 *Intelligence, 36*(5), 464–486. <https://doi.org/10.1016/j.intell.2007.10.003>
11
12
13 Sawilowsky, S. S. (2009). New effect size rules of thumb. *Journal of Modern Applied Statistical*
14
15 *Methods, 8*(2), 26.
16
17
18 Scullin, M. K., Ball, B. H., & Bugg, J. M. (2020). Structural correlates of commission errors in
19
20 prospective memory. *Cortex, 124*, 44–53. <https://doi.org/10.1016/j.cortex.2019.10.013>
21
22
23 Scullin, M. K., Bugg, J. M., & McDaniel, M. A. (2012). Whoops, I did it again: Commission errors
24
25 in prospective memory. *Psychology and Aging, 27*(1), 46–53.
26
27 <https://doi.org/10.1037/a0026112>
28
29
30 Shah, P., & Miyake, A. (1996). The separability of working memory resources for spatial thinking
31
32 and language processing: An individual differences approach. *Journal of Experimental*
33
34 *Psychology: General, 125*(1), 4–27.
35
36
37 Shipstead, Z., & Broadway, J. M. (2013). Individual differences in working memory capacity and
38
39 the Stroop effect: Do high spans block the words? *Learning and Individual Differences,*
40
41 *26*, 191–195. <https://doi.org/10.1016/j.lindif.2012.04.003>
42
43
44
45 Spieler, D. H., Balota, D. A., & Faust, M. E. (1996). Stroop performance in healthy younger and
46
47 older adults and in individuals with dementia of the Alzheimer's type. *Journal of*
48
49 *Experimental Psychology: Human Perception and Performance, 22*(2), 461–479.
50
51 <https://doi.org/10.1037/0096-1523.22.2.461>
52
53
54
55
56
57
58
59
60

1
2
3 St Clair-Thompson, H. L. (2007a). The influence of strategies on relationships between working
4
5 memory and cognitive skills. *Memory*, 15(4), 353–365.

6
7
8 <https://doi.org/10.1080/09658210701261845>
9

10 St Clair-Thompson, H. L. (2007b). The effects of cognitive demand upon relationships between
11
12 working memory and cognitive skills. *Quarterly Journal of Experimental Psychology*,
13
14 60(10), 1378–1388. <https://doi.org/10.1080/17470210601025505>
15
16

17
18 Towse, J. N., Hitch, G. J., Horton, N., & Harvey, K. (2010). Synergies between processing and
19
20 memory in children’s reading span: Synergies between processing and memory in
21
22 reading span. *Developmental Science*, 13(5), 779–789. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-7687.2009.00929.x)
23
24 7687.2009.00929.x
25
26

27
28 Towse, J. N., Hitch, G. J., & Hutton, U. (1998). A reevaluation of working memory capacity in
29
30 children. *Journal of Memory and Language*, 39(2), 195–217.
31
32 <https://doi.org/10.1006/jmla.1998.2574>
33
34

35
36 Towse, J. N., Hitch, G. J., & Hutton, U. (2000). On the interpretation of working memory span in
37
38 adults. *Memory & Cognition*, 28(3), 341–348. <https://doi.org/10.3758/BF03198549>
39

40
41 Turner, M. L., & Engle, R. W. (1989). Is working memory capacity task dependent? *Journal of*
42
43 *Memory and Language*, 28(2), 127–154.
44

45
46 Unsworth, N., Brewer, G. A., & Spillers, G. J. (2009). There’s more to the working memory
47
48 capacity—Fluid intelligence relationship than just secondary memory. *Psychonomic*
49
50 *Bulletin & Review*, 16(5), 931–937. <https://doi.org/10.3758/PBR.16.5.931>
51
52
53
54
55
56
57
58
59
60

1
2
3 Unsworth, N., Brewer, G. A., & Spillers, G. J. (2013). Working memory capacity and retrieval
4 from long-term memory: The role of controlled search. *Memory & Cognition*, *41*(2),
5 242–254. <https://doi.org/10.3758/s13421-012-0261-x>
6
7
8
9

10 Unsworth, N., Fukuda, K., Awh, E., & Vogel, E. K. (2014). Working memory and fluid intelligence:
11 Capacity, attention control, and secondary memory retrieval. *Cognitive Psychology*, *71*,
12 1–26. <https://doi.org/10.1016/j.cogpsych.2014.01.003>
13
14
15
16
17

18 Unsworth, N., Heitz, R., Schrock, J., & Engle, R. (2005). An automated version of the operation
19 span task. *Behavior Research Methods*, *37*, 498–505.
20
21
22

23 Unsworth, N., Redick, T. S., Heitz, R. P., Broadway, J. M., & Engle, R. W. (2009). Complex
24 working memory span tasks and higher-order cognition: A latent-variable analysis of the
25 relationship between processing and storage. *Memory*, *17*(6), 635–654.
26
27
28
29
30
31
32

33 Wahlheim, C. N., Ball, B. H., & Richmond, L. L. (2017). Adult age differences in production and
34 monitoring in dual-list free recall. *Psychology and Aging*, *32*(4), 338–353.
35
36
37
38
39
40

41 Wahlheim, C. N., & Huff, M. J. (2015). Age differences in the focus of retrieval: Evidence from
42 dual-list free recall. *Psychology and Aging*, *30*(4), 768–780.
43
44
45
46
47

48 Wahlheim, C. N., Richmond, L. L., Huff, M. J., & Dobbins, I. G. (2016). Characterizing adult age
49 differences in the initiation and organization of retrieval: A further investigation of
50 retrieval dynamics in dual-list free recall. *Psychology and Aging*, *31*(7), 786–797.
51
52
53
54
55
56
57
58
59
60

1
2
3 Waters, G. S., & Caplan, D. (1996). The measurement of verbal working memory capacity and its
4
5 relation to reading comprehension. *The Quarterly Journal of Experimental Psychology*
6
7 *Section A, 49(1)*, 51–79. <https://doi.org/10.1080/713755607>
8
9

10 Wickham, H. (2016). *Ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
11
12 <https://ggplot2.tidyverse.org>
13
14

15 Wingfield, A., & Kahana, M. J. (2002). The dynamics of memory retrieval in older adulthood.
16
17 *Canadian Journal of Experimental Psychology, 56(3)*, 187–199.
18
19

20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only

Table 1***Sample Sizes and Descriptive Statistics for WMC estimates from each Task and Site***

| | Overall | | ASU | | CSUS/SBU | | TU | | Site difference |
|---------|----------|-------|----------|-------|----------|-------|---------|-------|-----------------|
| | M | SD | M | SD | M | SD | M | SD | |
| OSpan | 56.80 | 13.72 | 58.19 | 12.77 | 50.07 | 16.73 | 54.45 | 14.72 | ** |
| | n = 1685 | | n = 1180 | | n = 121 | | n = 384 | | |
| SymSpan | 28.61 | 7.77 | 29.63 | 7.79 | 29.40 | 7.14 | 26.56 | 7.75 | ** |
| | n = 1053 | | n = 585 | | n = 121 | | n = 347 | | |
| RSpan | -- | | 52.52 | 14.81 | -- | | -- | | N/A |
| | | | n = 1062 | | | | | | |

Note. ** indicates significant site differences observed at the $p \leq .001$ level.

For Review Only

Table 2

Descriptive statistics for storage performance, processing accuracy, and processing RT measures

| Span Task | | M | SD | Skew | Kurtosis | Norm Violation? | |
|--------------------|----------------------|----------------------|---------|---------|----------|-----------------|---|
| OSpan | Proc Acc | 0.91 | 0.07 | -2.71 | 11.81 | Y | |
| | Mean Proc Prac RT | 3243.03 | 1385.79 | 1.51 | 3.68 | N | |
| | SD Proc Prac RT | 1768.66 | 1306.06 | 2.41 | 9.23 | Y | |
| | RT Dist, μ | 2546.80 | 761.69 | 1.63 | 5.40 | Y | |
| | RT Dist, σ | 489.49 | 360.66 | 2.41 | 10.03 | Y | |
| | RT Dist, τ | 1253.81 | 741.13 | 1.97 | 7.08 | Y | |
| | Time-Out Err | 1.43 | 1.89 | 3.37 | 21.70 | Y | |
| | Incorrect Err | 5.17 | 4.41 | 3.28 | 17.43 | Y | |
| | SymSpan | Proc Acc | 0.91 | 0.10 | -2.66 | 10.42 | Y |
| | | Mean Proc Prac RT | 2141.52 | 1041.46 | 1.58 | 4.29 | Y |
| SD Proc Prac RT | | 1190.12 | 831.42 | 2.86 | 14.30 | Y | |
| RT Dist, μ | | 1721.80 | 676.88 | 1.74 | 3.85 | Y | |
| RT Dist, σ | | 309.22 | 303.46 | 1.84 | 3.84 | Y | |
| RT Dist, τ | | 746.69 | 467.15 | 1.82 | 6.67 | Y | |
| Time-Out Err | | 0.80 | 1.28 | 2.56 | 9.58 | Y | |
| Incorrect Err | | 3.18 | 3.74 | 3.04 | 13.02 | Y | |
| RSpan | Proc Acc | 0.90 | 0.11 | -3.43 | 15.44 | Y | |
| | Mean Proc Prac RT | 3947.43 | 1276.23 | 1.35 | 3.85 | N | |
| | SD Proc Prac RT | 1456.58 | 749.58 | 2.45 | 10.37 | Y | |
| | RT Dist, μ | 3436.06 | 1054.71 | 0.97 | 4.77 | Y | |

| | | | | | |
|-------------------|--------|--------|------|--------|---|
| RT Dist, σ | 669.88 | 415.34 | 3.16 | 20.40 | Y |
| RT Dist, τ | 991.19 | 587.05 | 2.11 | 11.16 | Y |
| Time-Out Err | 1.61 | 2.28 | 8.03 | 125.49 | Y |
| Incorrect Err | 5.62 | 7.40 | 3.73 | 17.83 | Y |

Note. WMC: capacity estimate; Proc Acc: processing task accuracy; Mean Proc Prac RT: mean practice processing task RT; SD Proc Prac RT: standard deviation practice processing task RT; RT Dist, μ : mu component of the ex-Gaussian analysis for task-relevant RTs; RT Dist, σ : sigma component of the ex-Gaussian analysis for task-relevant RTs; RT Dist, τ : tau component of the ex-Gaussian analysis for task-relevant RTs; Time-Out Err: processing task errors due slow responding; Incorrect Err: processing task errors due to incorrect responding. “Y” under Norm Violation? column indicates that the assumption of normality was violated, defined as skew > |2| and/or kurtosis > |4|, whereas “N” indicates that skew and kurtosis values were found to be in the acceptable range. For rows marked “Y”, non-parametric statistical tests can be found in the Supplemental Materials, and in cases where non-parametric and parametric findings differed these are footnoted throughout the manuscript.

Table 3*Cronbach's alpha scores for processing and storage components of each task*

| | OSpan | SymSpan | RSpan |
|---------------------|-------|---------|-------|
| Processing Accuracy | .77 | .81 | .91 |
| Processing RT | .93 | .89 | .91 |
| Storage Accuracy | .92 | .83 | .92 |

For Review Only

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

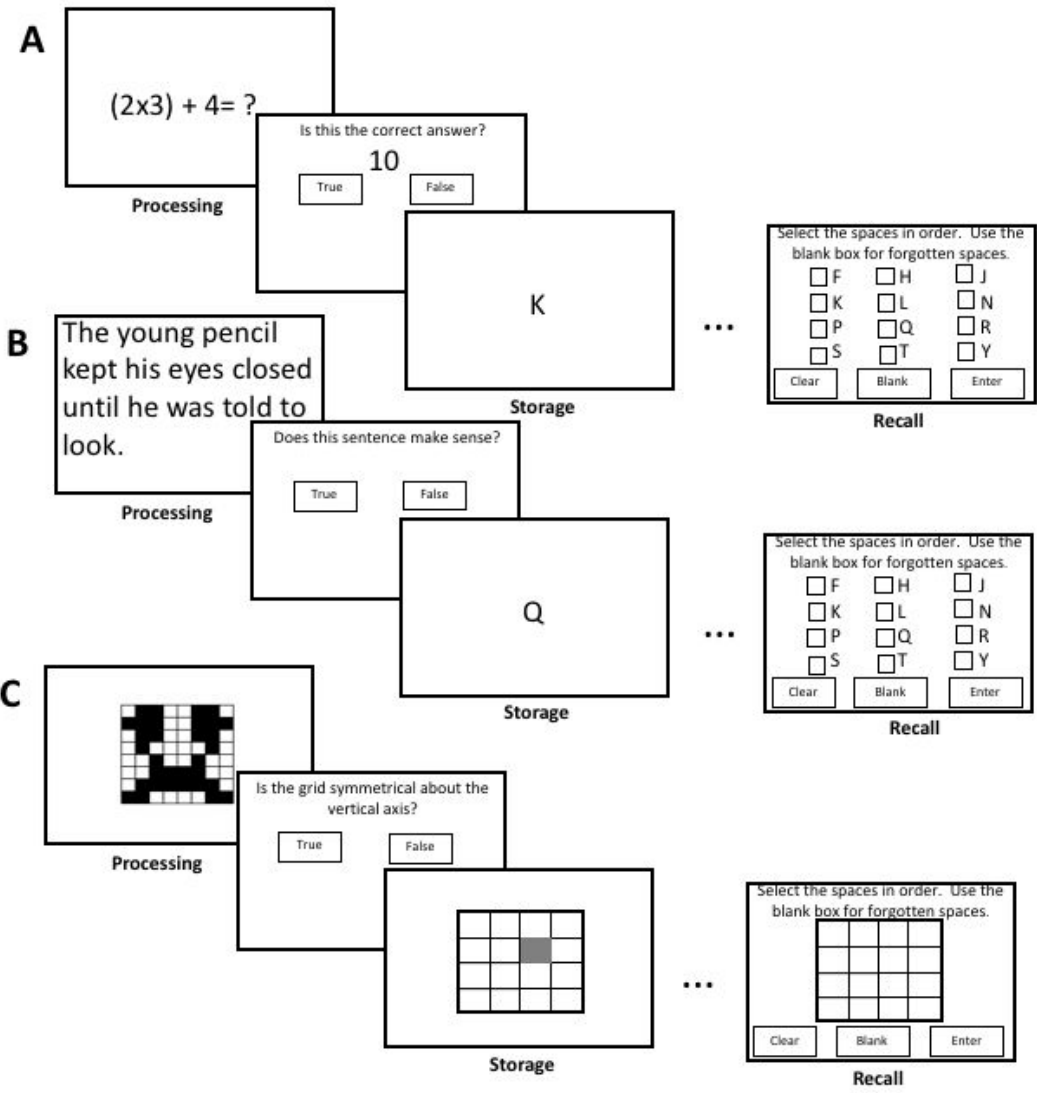


Figure 1. Complex Span Working Memory task schematics depicting Operation Span (panel A), Reading Span (panel B), and Symmetry Span (panel C).

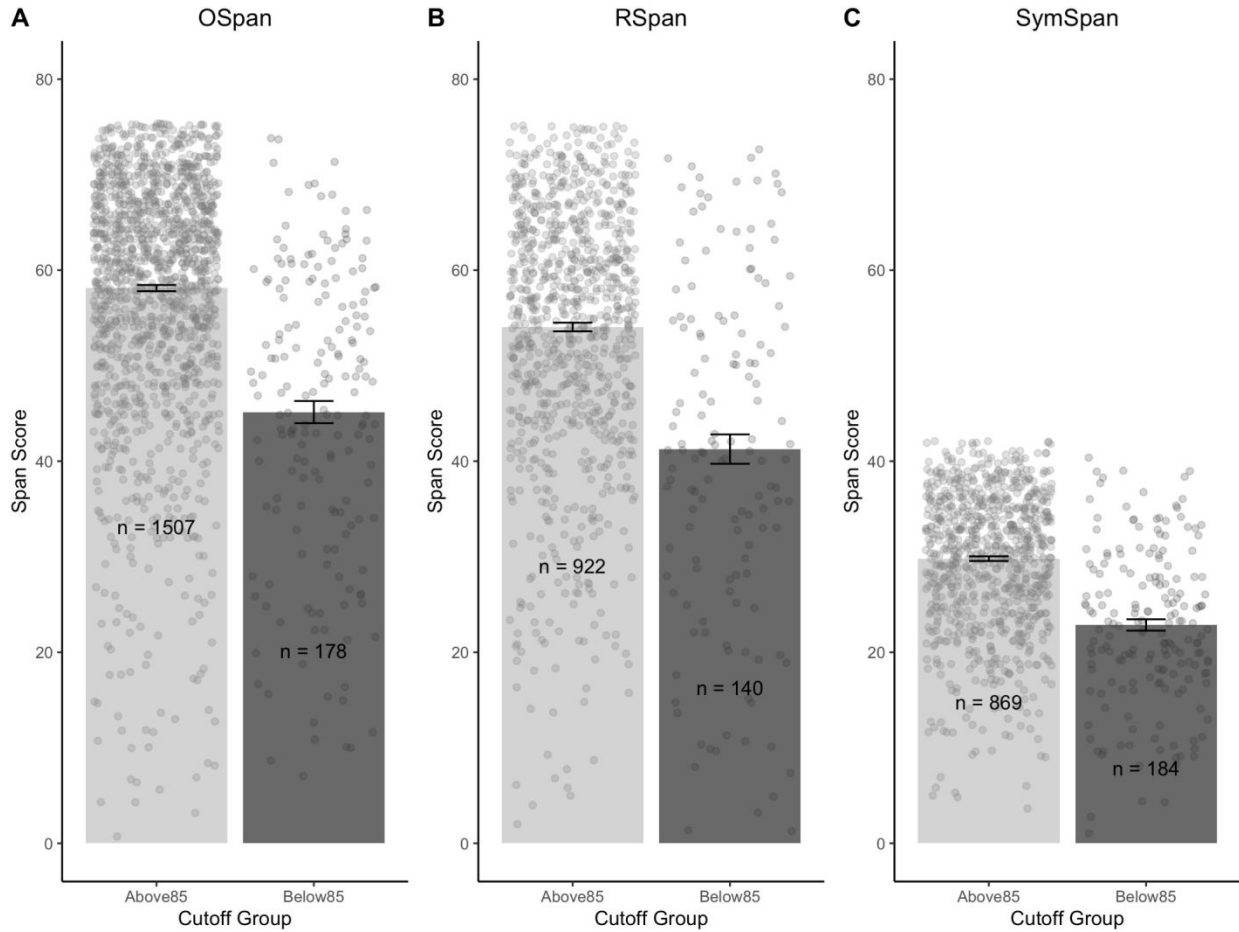


Figure 2. Participants who meet or exceed the 85% processing accuracy cutoff have significantly higher span scores across tasks compared to those who miss the cutoff. Individual data points are depicted by the circles and the error bar represents standard error.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

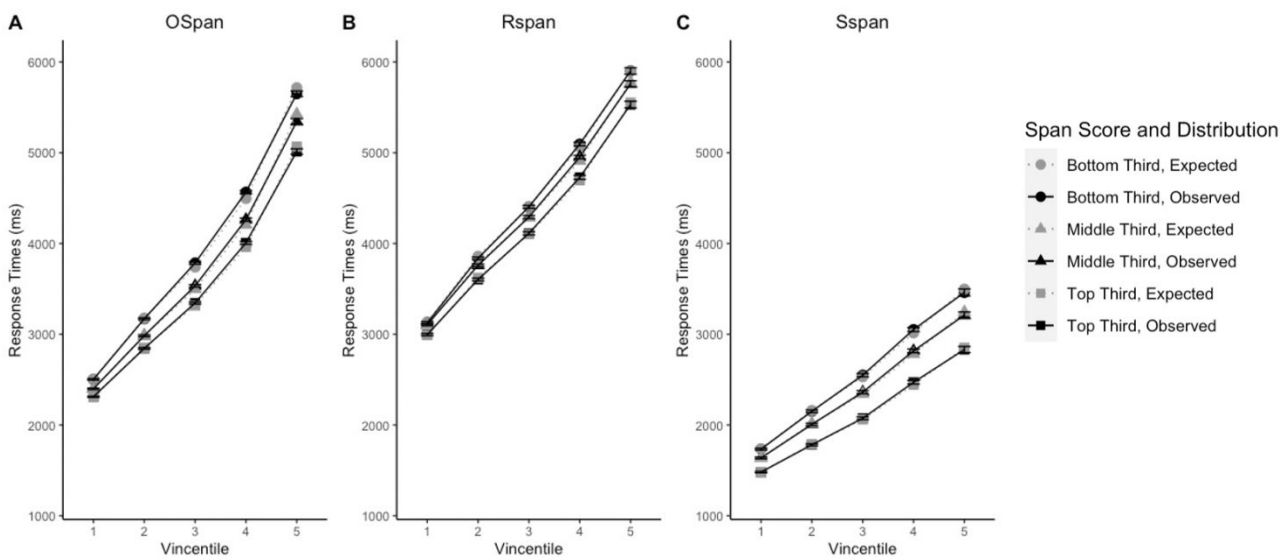


Figure 3. Vincentile plots by task (panel A: Operation Span, panel B: Reading Span, panel C: Symmetry Span) and by WMC partial span score (top third, middle third, and bottom third). This depiction is for illustrative purposes only; WMC was included in our models as a continuous variable. Expected values for each group are denoted by the grey dotted lines and the observed values are denoted by the black solid lines. The substantial overlap between predicted and observed values indicates that these data are fit well by the exGaussian function.

Preprint Only

Supplemental Materials

Non-Parametric Analyses: OSpan

The following nonparametric analyses were conducted in the OSpan dataset due to the observation of violations of normality for the factors included in the tests:

For our research question regarding differences in WMC as a function of processing accuracy (question 1), we conducted a Wilcoxon Rank Sum Test. Consistent with the Welch's two sample t-test for unequal variances, results here indicate a significant difference in WMC estimates by processing performance ($W = 203059$, $p < .001$, 95% CI [11.00, 15.00], $r = -0.35$) for those who met/exceeded the 85% criterion cutoff, consistent with our parametric analysis.

For the research question regarding the extent to which the SD of practice RTs are related to WMC (question 2b), a Spearman's Rho correlation revealed a non-significant correlation ($\rho(1683) = -0.03$, $p = 0.175$) between WMC and SDs of RTs during the practice portion of the processing task. This finding is consistent with the result of the parametric Pearson's Product Moment correlation.

For our question regarding the extent to which ex-Gaussian characterization of task-derived RTs predict WMC (question 3b), we entered the components μ , σ , and τ as simultaneous predictors into a quantile regression where WMC estimates served as the outcome variable. Consistent with the results of the linear regression, the σ component emerged as the only significant predictor in the model ($\beta = -0.21$, $t(1681) = -3.85$, $p < .001$; τ and μ p values $> .14$).

We also conducted a quantile regression for our question regarding the relation between processing errors and WMC (question 3c). Consistent with the linear regression model, both time out errors ($\beta = -0.17$, $t(1682) = -5.66$, $p < .001$) and incorrect errors ($\beta = -0.35$, $t(1682) = -9.24$, $p < .001$) emerged as significant predictors in the model.

Turning to our last research question regarding improvements in model fit for RT distribution + error predictor models over error alone and RT distributional components alone (question 4), we tested the quantile regression models described above (RT distribution only, error only) against a model that added error type (time-out errors, incorrect errors) to the ex-Gaussian components to predict WMC estimates. Consistent with the comparison of linear regression models, the RT distribution + error model explained significantly more variance in WMC estimates compared to the RT-distribution only model (model comparison: $F(2, 1679) = 71.90$, $p < .001$) and the error only model (model comparison: $F(3, 1679) = 9.18$, $p < .001$).

Non-Parametric Analyses: SymSpan

The following nonparametric analyses were conducted in the SymSpan dataset due to the observation of violations of normality for the factors included in the tests:

For our research question regarding differences in WMC as a function of processing accuracy (question 1), we conducted a Wilcoxon Rank Sum Test. Consistent with the Welch's two sample t-test for unequal variances, results here indicate a significant difference in WMC estimates by processing performance ($W = 118992$, $p < .001$, 95% CI [6.00, 9.00], $r = -0.32$) for those who met/exceeded the 85% criterion cutoff, consistent with our parametric analysis.

For the research questions regarding the extent to which mean and SD of practice RTs are related to WMC (questions 2a and 2b), we conducted Spearman's Rho correlations. The result for question 2a (mean practice RT) revealed a significant correlation ($\rho (1051) = -0.16$, $p < .001$) between WMC and mean RT during the practice portion of the processing task. This finding is consistent with the result of the parametric Pearson's Product Moment correlation. Inconsistent with the result of the parametric test for question 2b (SD of practice RT), we observed a significant Spearman's Rho correlation between RT SD during the practice portion of the processing task and WMC ($\rho (1051) = -0.10$, $p < .001$).

For our question regarding the extent to which ex-Gaussian characterization of task-derived RTs predict WMC (question 3b), components μ , σ , and τ were entered as simultaneous predictors into a quantile regression where WMC estimates served as the outcome variable. Consistent with the results of the linear regression, the τ component emerged as a significant predictor in the model ($\beta = -0.13$, $t (1049) = -2.85$, $p = .004$; μ and σ p values $> .04$).

We also conducted a quantile regression for our question regarding the relation between processing errors and WMC (question 3c). Consistent with the linear regression model, both time out errors ($\beta = -0.25$, $t (1050) = -4.57$, $p < .001$) and incorrect errors ($\beta = -0.38$, $t (1050) = -7.930$, $p < .001$) emerged as significant predictors in the model.

Turning to our last research question regarding improvements in model fit for RT distribution + error predictor models over error alone and RT distributional components alone (question 4), we tested the quantile regression models described above (RT distribution only, error only) against a model contained both elements (RT distribution + error). Consistent with the comparison of linear regression models, the RT distribution + error model explained significantly more variance in WMC estimates compared to the RT-distribution only model (model comparison: $F (2, 1047) = 48.56$, $p < .001$) and the error only model (model comparison: $F (3, 1047) = 27.84$, $p < .001$).

Non-Parametric Analyses: RSpan

The following nonparametric analyses were conducted in the RSpan dataset due to the observation of violations of normality for the factors included in the tests:

For our research question regarding differences in WMC as a function of processing accuracy (question 1), we conducted a Wilcoxon Rank Sum Test. Consistent with the Welch's two sample t-test for unequal variances, results here indicate a significant difference in WMC estimates by processing performance ($W = 91722$, $p < .001$, 95% CI [10.00, 16.00], $r = -0.25$) for those who met/exceeded the 85% criterion cutoff, consistent with the findings from our parametric analysis.

For the research question regarding the extent to which the SD of practice RTs are related to WMC (question 2b), we conducted a Spearman's Rho correlation which revealed a non-significant correlation ($\rho(1060) = -0.07$, $p = .021$), consistent with the findings from our parametric analysis.

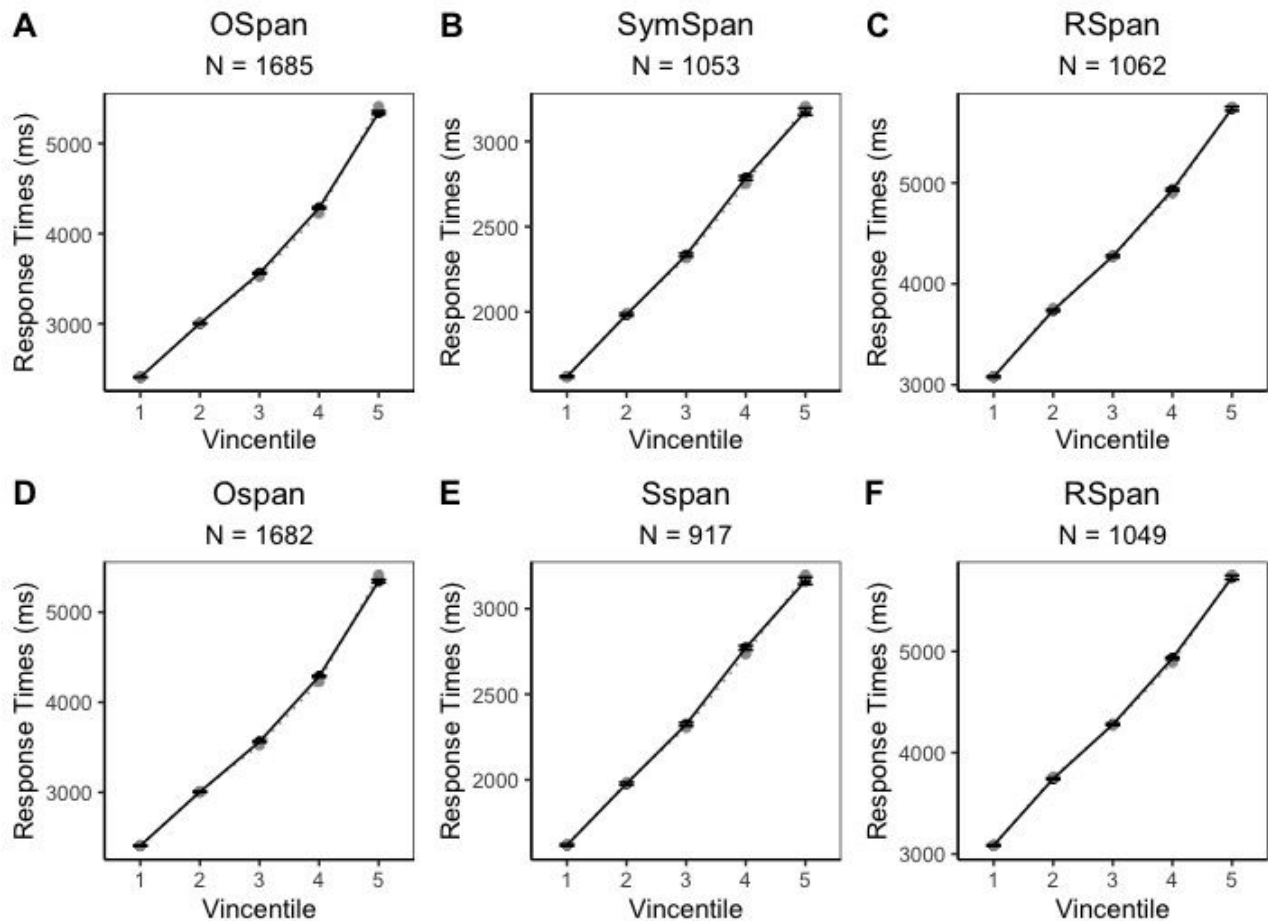
For our question regarding the extent to which ex-Gaussian characterization of task-derived RTs predict WMC (question 3b), components μ , σ , and τ were entered as simultaneous predictors into a quantile regression where WMC estimates served as the outcome variable. Consistent with the parametric analysis, the σ component emerged as the only significant predictor in the model ($\beta = -0.007$, $t(1058) = -2.85$, $p = .004$; μ and τ p values $> .71$).

We also conducted a quantile regression for our question regarding the relation between processing errors and WMC (question 3c). Consistent with the linear regression model, both time out errors ($\beta = -0.21$, $t(1059) = -4.33$, $p < .001$) and incorrect errors ($\beta = -0.35$, $t(1059) = -6.39$, $p < .001$) emerged as significant predictors in the model.

Turning to our last research question regarding improvements in model fit for RT distribution + error predictor models over error alone and RT distributional components alone (question 4), we tested the quantile regression models described above (RT distribution only, error only) against a model that included both RT distribution and error type as predictors. Consistent with the comparison of linear regression models, the RT distribution + error model explained significantly more variance in WMC estimates compared to the RT-distribution only model (model comparison: $F(2, 1056) = 38.75$, $p < .001$). Inconsistent with the comparison of linear regression models, the RT distribution + error model did not explain significantly more variance in WMC estimates compared to the error only model at the level of $p < .01$ (model comparison: $F(3, 1056) = 2.88$, $p = .035$).

Ex-Gaussian Model Fits

Below we present vincentile plots depicting model fits for the Ex-Gaussian characterization of RTs. Vincentile plots allow for examination of the raw RT distribution across conditions without making assumptions about the underlying shape of the distribution (Andrews & Heathcote, 2001). For each of the three span tasks, vincentiles were computed for by rank ordering raw RTs from shortest to longest for each participant, and calculating the mean of the first 20%, the second 20%, and so on. These plots present the best-fitting predicted vincentiles (dotted lines) superimposed on the observed vincentiles (solid lines). Across all tasks, the minimal divergence of the expected line from the observed line indicates that the data was well fit by the Ex-Gaussian function. Across all three complex span tasks, we present plots including all participants as well as plots excluding participants with fewer than 35 correct RTs.



Supplementary Figure 1. Panels A, B, and C depict observed and expected vincentiles for the complete samples. Panels D, E, and F depict observed and expected vincentiles when including only participants who contributed 35 or more correct RTs.

References

- 1
2
3
4
5
6 Andrews, S., & Heathcote, A. (2001). Distinguishing common and task-specific processes in word
7 identification: A matter of some moment? *Journal of Experimental Psychology: Learning,*
8 *Memory, and Cognition*, 27(2), 514–544. <https://doi.org/10.1037/0278-7393.27.2.514>
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only