# Supplemental Information

## Methods & Materials

*Details of DISCERN implementation*

A fundamental design choice of DISCERN is to incorporate multiple simple recurrent networks (SRNs, ref. 1) into a chain of hierarchically organized modules (Figure 1, main text). SRNs are backpropagation networks in which a copy of the previous hidden layer activation is saved at each computational step, and then used as a re-entrant input during the next computational step, thereby providing a sequential memory capacity.  SRNs have been successfully used to simulate aspects of normal sentence processing (1), sequential learning (2), and cognitive development (3-5) as well as hallucinated speech in schizophrenia (simulated as working memory (WM) disconnection, refs. 5,6), and therefore constitute a natural starting point for studying disruptions in sequential language behavior in patients with schizophrenia. With SRNs it was possible to model several illness mechanisms suggested by the research literature in a natural fashion (Figure 1, main text).

DISCERN learned a lexicon consisting of 159 words, including 10 specific agents or characters in the stories (e.g. "Stacy," "I," "Vito") and other more general agent references (e.g., "police," "mafia," "boss").  The story corpus consisted of two sets of 14 stories: autobiographical and crime-related.  The two story sets each had 5 specific agents that were entirely non-overlapping, with some overlap between other lexical elements.  For example, "wedding" occurred exclusively in the former, "bombing" occurred exclusive in the latter, while "boss," "meeting" and "police" occurred in both.

The original DISCERN system (7,8) was modified so that storage and retrieval of stories were associated with a positive/negative emotional valence code having five possible values (--, -, -/+, +, ++). Seven of the autobiographical stories had a positive emotional valence, and seven

had a negative emotional valence.  In contrast, seven of the crime stories had neutral emotional valence, and seven had a negative emotional valence.

The modules in DISCERN communicate using distributed representations of word meanings, i.e. fixed-size patterns of neuron activations, stored in a central lexicon.  These representations are learned based on how the words are used in the example stories, using the FGREP algorithm ("Forming Global Representations with Extended Backpropagation;" ref. 7,8), a modified version of backpropagation that treats input representations as an additional layer of adjustable weights.  During the memory storage phase (Figure 1A, left side of loop, main text), the input text is first translated into input activation patterns by the lexicon, then presented to the sentence parser SRN one word at a time.  The sentence parser builds a static representation of each sentence step-by-step as input words are received.  Individual sentences are represented through case roles corresponding to agent, predicate, indirect object, modifier and direct object. Each case is represented by a slot filled by one of the word representations.   A sequence of words in a sentence, for example, "`Vince entered the LA airport,`" would be turned into a static representation, [agent: Vince, predicate: entered, indirect object:__ modifier: LA, direct object: airport],  where each plain-text word represents a pattern of neural activations, and the underscore ("__") represents the "blank" pattern consisting of all zeroes.  At the end of each sentence, the sentence representation is passed on to the story parser.  This SRN in turn transforms sequences of sentences into static representations of scripts.  A script representation consists of the name of the script and the sequence of concepts filling its slots. The classic restaurant script, for instance, is comprised of a sequence of sentences expressing entering a restaurant, sitting down at a counter or table, ordering food, liking or disliking the food, receiving the bill, paying the bill and leaving a large or small tip.  This script includes slots for customer name, restaurant name, and various foods, which are functionally similar to case-role representations for sentences.  A script, in other words, corresponds to a sequence of sentences comprising a standardized schema.  In our study, the same script was often

incorporated into both autobiographical and crime stories with different words/concepts "filling" the script's "slots."

The original DISCERN system was expanded to include a memory encoder module that processes stories as a sequence of multiple scripts. This modification generated naturally occurring "breakpoints" in stories (corresponding to transitions between scripts) that could facilitate expression of derailment under various conditions. In brief, the memory encoder associates each script with a memory cue that is later used by the story generator to recall it. A script instance paired with its memory cue is called an *episodic memory trace*, i.e., an occurrence of a script that is stored in episodic memory. The memory encoder is a Recursive Auto-Associative Memory (RAAM; ref. 9), a neural network architecture that forms compact distributed representations of recursive data structures such as lists. RAAM networks are feedforward networks trained to reproduce their own input, forcing them to form compressed representations of inputs in their hidden layer. These compressed representations can then be re-used as part of the input to the RAAM, recursively building representations of arbitrarily long lists. In DISCERN, RAAM representations of sequences of scripts are used as cues to address episodic memory by content. Figure S1 shows a RAAM network that is being used to create a memory cue. The network uses the current cue (a compressed partial story) as part of the input to form the next cue in its hidden layer. In this way, the network steps backwards through a story, producing a compressed representation of the rest of the story at each step, and associating each new cue with the script used to create it. Shared emotional valence facilitates transitions from one script to the next within the stories. Figure S2 illustrates this process.

With the memory traces in place, the system is ready to recall the stories (Figure 1A, main text, right side). The story generator module (Figure 1B, main text) is cued with the first script in a story. At its output, it produces the representation of each sentence in the story one at a time, until a special end-of-story representation. Together with each sentence representation, it produces a memory cue that can be thought of as the system's discourse

plan.  The cue is used to retrieve the next memory trace from episodic memory, thus determining the story generator's own next input.  In this way, the story generator steps through each sentence of a story, and accesses each memory trace encoding it.  Note that as long as the story generator produces sentences belonging to the same script, the memory cue does not change.  However, when the story generator produces the last sentence of the script, the cue does change, and the input is replaced by the memory of the next script.  Figure S3 shows an example of such a switch from one script to the next.  Based on evidence of an editor function in human speakers (10), an output filter was attached to the story generator module to block sentence-level outputs falling below a quality threshold.  The filter tended to eliminate virtually all ungrammatical constructions and reduce many other errors produced by illness mechanisms at the cost of reducing successful recall.  Finally, the sentence generator, last in the chain, takes the sentence representations produced by the story generator and turns them back into a sequence of individual word representations.  The system then outputs plain text translations of these word representations as provided by the lexicon.

The sentence parser and the sentence generator were trained initially for 5000 iterations (or epochs) of the entire corpus, using FGREP to develop the word representations.  Each word representation consisted of a fixed-size pattern of 12 neuron activations.  With the word representations in place, 30 different DISCERN systems or exemplars were then trained starting with different random connection weights.  The hidden layer of the memory encoder had 48 neurons, while the story generator had 150 hidden neurons.  Sentence parsers and generators had 250 hidden neurons and the story parser had 225 hidden neurons.

Modules were trained in a chain, with the output from one module used as the input for the next.  Starting with the sentence module, new modules were added to the chain as meaningful input became available during the course of learning.  The learning rate for each module was always set to 0.4 times the average output error of the module during the last training epoch.  Thus, as the output error decreased during training, the learning rate decreased

4

automatically to allow fine-tuning of network response.  A total of 70,000 backpropagation learning epochs were employed overall for each DISCERN exemplar distributed across the different modules.  To provide a clean starting point for assessing performance after backpropagation training was completed, each exemplar's episodic memory was cleared of all individual story instances, and the 28 stories were once again read and stored in the episodic memory module.  After training, DISCERN was able to reproduce all 28 stories almost perfectly with the percentage of sentences correctly reproduced averaging 96.2% across the 30 DISCERN systems.

*Details of human story recall study*

In the human study, the story recall performance of 21 normal subjects and 43 subjects with schizophrenia or schizoaffective disorder was compared.  Patients were symptomatically stable outpatients.  Healthy control subjects were recruited by advertisement and word of mouth.  Psychiatric diagnoses of patients were based on DSM-IV criteria established using the Comprehensive Assessment of Symptoms and History (CASH, ref. 11).  Patients were prospectively divided into two subgroups: those who definitely demonstrated evidence of fixed delusions with a plot-like or narrative scheme, and those who produced questionable or no evidence of these delusions.  Typical examples included God choosing the patient to eliminate racial oppression, and the patient being trailed by Homeland Security agents due to allegations of terrorist activities.  The absence of psychiatric diagnosis in healthy controls was confirmed using the non-patient version of the Structured Clinical Interview for the DSM-IV (SCID, ref. 12). Antipsychotic drug level was quantified as chlorpromazine equivalents (13-15).  Out of this group of subjects, 20 healthy controls and 37 patients provided story recall data at seven days. To estimate verbal abilities, the Wechsler Adult Intelligence Scale-III vocabulary test (16) was administered.

The story recall task consisted of three prerecorded stories presented binaurally on headphones as described in the main text.  All three stories involved a gift.  Two of the stories

("The Gift" (17) and "The Hitchhiker" (custom-written for this study)) shared other content, involving a travel theme and a specific character reference ("wispy old man"). Stories were presented in random order. Immediate recall, recall at 45 minutes after exposure to stories and recall after seven days were tape-recorded and transcribed for analysis by a rater not involved in data collection who was blind to group, presence/absence of fixed delusions, and subject identification. Seven-day recall was by surprise to prevent preparatory rehearsal during the intervening period, and comprised the narrative language behavior against which alternative DISCERN models were assessed. Below are the verbatim instructions used for each subject prior to presenting stories, and at 45-minute and 7-day recall:

> *Subject instructions prior to playing the stories were as follows:*
> "I am going to play a tape, and you will hear a man's voice reading a story. The story is short – about 5 or 6 sentences. The idea is for you to listen carefully, and then, when it is finished playing, I will ask you to recall as much of the story as possible. Don't worry about "passing" or "failing" – there is no such thing on this task. Just do the best you can. This procedure will be repeated for two additional stories."

> *For the 45 minute rehearsal, the instructions were:*
> "A little while ago, I had you listen to three stories played back on a tape recorder and then recall them. Now I want you to recall them again, you can recall them in any order. Just do your best to recall as much of each story as possible."

> *For 7-day recall, the instructions were:*
> "We now want you to try to recall as completely as possible the three stories you heard on headphones last week. This may be somewhat of a surprise, but we didn't want you to rehearse in your mind the content of the stories over the last week. We wished this to be a test of story memory capacity that occurs naturally without practice or rehearsing. Please take your time and try to repeat as much of the story as you can recall. The words of the story needn't be exact and no one is able to recall these stories completely. Just try as hard as you can."

> *If the subject cannot recall any detail of a particular story, then (s)he is provided with the corresponding prompt:*
> 1. There was a story about flowers.
> 2. There was a story about a hitchhiker.
> 3. There was a story about a robbery.

*Comparing DISCERN and human language performance*

Several variables reflecting story recall deviance by human subjects and DISCERN simulations were developed and tested. Five of them were eventually discarded:

- Within-story accretions.  These were words and phrases derived from the target story that were grouped together in ways that misrepresented meaning.  Limited flexibility of DISCERN's output language precluded these types of errors.

- Pronoun reference failures.  This version of DISCERN used did not generate pronouns (outside of the first person pronoun, "I") so this variable could not be used to evaluate illness mechanisms.

- Word/phrase insertions.  These were words or phrases that were inserted into otherwise grammatical responses that derived from outside of the story; the more rigid sentence syntax of DISCERN did not permit these text insertions.

- Ungrammatical constructions.  Surprisingly, such errors were virtually non-detectable in human data but were quite prevalent in all DISCERN illness mechanisms prior to filtering other than hyperlearning.  However, this form of error was virtually eliminated in DISCERN when the output filter was adjusted to match language profiles of human patients.

- Between-story migrations.  These were errors clearly where text from one story intruded into the recall of another story.  Given that we had full knowledge of the entire story corpus for DISCERN, a large majority of simulated errors fell into this category.  For human story recall, on the other hand, errors could presumably derive from an extremely large number of narrative memories (or other sources) not involving the three target stories used in the experiment.  Therefore, human and DISCERN story migration counts were not comparable.

Four variables could, however, be scored comparably for both humans and DISCERN story recall while demonstrating sufficient variance to allow contrasts between human subject groups and between illness models.  Importantly, derailed clauses were not scored (see main text for definition) if inserted text was interpretable as any of the error types described above.  A manual

for scoring narrative memory distortions, including a breakdown of propositional structure of the three stories used for quantifying recall success, is available on request from the first author.

More extensive editing/filtering of language outputs by human subjects and DISCERN was projected to reduce errors at the cost of reducing successful recall. Therefore, when comparing human and DISCERN story recall performance, the three commission error variables described above (derailments, agent-slotting errors and lexical misfires) were re-calibrated as penetrance scores, where totals for each type of error were divided by recall success (scored as kernels successfully paraphrased across stories). This strategy also accommodated the fact that the number of propositions in the DISCERN story corpus was much greater than the human story corpus, which provided much greater opportunity for error for the former.

To measure how well a DISCERN exemplar matched the human data, a mean square deviation metric was used (18), and represented in the following form:

$$GOF^{C/P}(D,m,f) = \sum_{i=1}^{4} \frac{(\overline{V}_i^{C/P} - V_i(D,m,f))^2}{SE(V_i^{C/P})}$$

where $GOF^{C/P}(D,m,f)$ is the goodness-of-fit of a given DISCERN exemplar, *D,* with mechanism parameter, *m,* and filter parameter, *f,* calculated relative to either the group of human healthy controls (*C*) or human patients with schizophrenia (*P*), $\overline{V}_i^{C/P}$ is the mean value of the story-recall variable, *i* (recall success, derailment penetrance, lexical misfire penetrance, agent-slotting-error penetrance, see Table 4, main text) calculated for the corresponding subject groups, *C* and *P*, $SE(V_i^{C/P})$ is the standard error for variable, *i*, and $V_i(D,m,f)$ is the score for that variable ascertained for DISCERN exemplar, *D*, with mechanism parameter, *m*, and filter setting, *f*. Incorporating standard errors of human variables as divisors of sums of square differences has the effect of adjusting contributions to *GOF* of individual variables in terms of their respective spans of dispersion in the human population. A lower *GOF* indicates better fit. For each DISCERN exemplar, parameters were selected separately to produce the smallest *GOF* relative

to the four story recall criterion variables determined for the healthy control group and the

patient group.  In order to ascertain the closest match for two-dimensional illness mechanisms,

the parameter space was searched across a 100 (illness mechanism setting) x 1000 (output

filter) grid to minimize *GOF* relative to mean healthy control data.  This optimization was then

repeated relative to mean patient data.  Note that in this analysis, it was not necessary to make

assumptions about how model behaviors were distributed for each mechanism since best-fit

*GOF* was assessed relative to a family of 30 independently generated DISCERN exemplars.

These findings are shown in Figure 2A (*GOF* for the 30 DISCERN exemplars relative to healthy

controls) and Figure 2B (*GOF* for the 30 DISCERN exemplars relative to patients) of the main

text.  For three-dimensional simulations, closest matches to patient data were sought by

searching a 40 (illness mechanism setting 1) x 40 (illness mechanism setting 2) x 1000 (output

filter) grid.  These findings are shown in Figure 2C of the main text.  The best-fitting variance-

covariance structure according to BIC (Schwartz-Bayesian Information Criterion) was compound

symmetry heterogeneous (CSH). This structure assumes equal covariances within cluster and

allows for unequal variance per mechanism/group.

The two-dimensional hyperlearning mechanism applied to the memory encoder module

and to the two-dimensional WM disconnection mechanism produced a robustly better fit to

patient story recall performance than the other six mechanisms, but these two mechanisms

were not significantly different from each other (Table 5, main text).

The goal of the second set of simulations was to determine whether adding a second

model-fitting parameter to the two best two-dimensional mechanisms (WM disconnection and

memory-encoder hyperlearning) resulted in a significantly better-fit to the patient story-recall

performance profile, and whether either of these three-dimensional mechanisms proved to be a

significantly better-fit to the patient story-recall relative to the other.  To accomplish these

objectives, the best *GOF* to the story-recall disturbance profile of patients was used as the

dependent variable and simulation was used as the clustering factor (again reflecting the 30

DISCERN exemplars), while dimensionality (2D vs. 3D) and model (hyperlearning vs. WM disconnection) were treated as within-subject factors.   As expected, 3D expansions of both mechanisms produced further improvements in GOF.  Most importantly, 3D hyperlearning proved to be a significantly better match to the story-recall profile of patients relative to 3D disconnection, thereby providing a new illness model that should be tested in future clinical studies.

In order to model fixed, self-referential delusions, DISCERN's agent-slotting errors need to be systematic, i.e. the same confusion of a personal-story and a crime-story agent needs to recur in the output stories.  This systematicity was assessed using a randomization test.  This test generated cross-context errors randomly using the same base rate of cross-context agent-slotting errors exhibited by each of the 30 DISCERN exemplars, and counting how many of the errors generated by random within-context agent selection turned out to repeat earlier systematic cross-context errors (in the same or opposite direction, see results section for examples).  A count of these random-generated systematic errors was repeated 10,000 times for different random selections of agents within-context, and compared to that observed for the 30 DISCERN simulations.

## *Results*

### *Additional findings*

Pooling data across both groups of human subjects, there was no significant correlation between any of the performance variables, and age, parental education level, or WAIS-scaled vocabulary, assuming an uncorrected cut-off of $\alpha = 0.05$.  Within just the patient group, number of hospitalizations and antipsychotic dose (scored as chlorpromazine equivalents) were also not significantly correlated with any of the performance variables using the same cut-off.

## *Discussion*

### *Additional comments*

It is noteworthy that Friston (19) has proposed a disconnection hypothesis for schizophrenia, where illness arises from altered synaptic efficacy in neural systems responsible for emotional learning and memory. Hyperlearning recalls this hypothesis insofar as the functional "site" of pathophysiology is memory consolidation; story memories in DISCERN all had an emotional valance, and the net result of hyperlearning is altered synaptic efficacy within modules.

One other limitation of the DISCERN findings is that best-fit hyperlearning simulations also confused agents drawn from the same story. For instance, `Stacy`, the girlfriend, was often exchanged with `Mary`, the fiancée of the boss, `Joe`. Exchanging the identity of two persons, both personally known to the patient, does not commonly occur in schizophrenia. It is possible, however, that humans overlearn deployment of personally known agents within autobiographical stories, thereby reducing likelihood that these agents are exchanged across these stories. A future computational study will test this hypothesis.
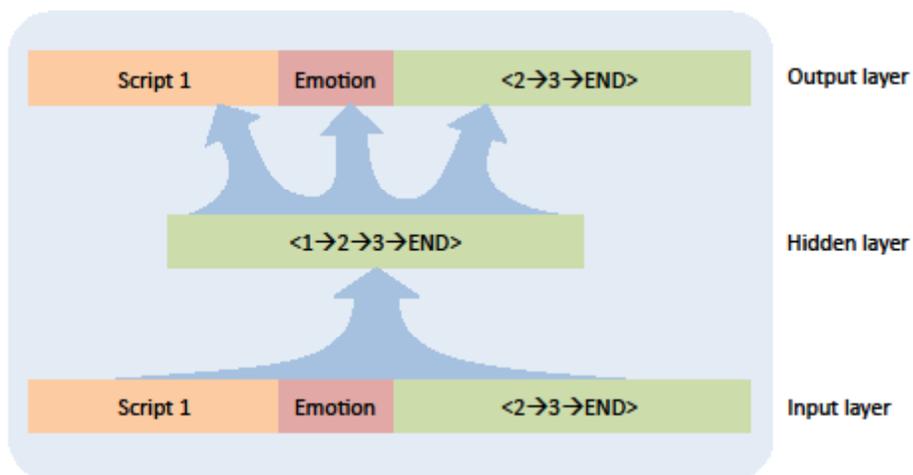
**Figure S1.** The memory encoder in DISCERN is a Recursive Auto-Associative Memory (RAAM, ref. 9), a neural network that is trained to reproduce its input in the output layer, forcing the input information to be compressed in the smaller hidden layer. The figure illustrates how the network creates a compressed representation of an entire story (consisting of scripts 1, 2, 3, and an end-of-story representation, denoted as <1→2→3→END>) in its hidden layer, given as input the slot filler representations and emotion code of the first script, as well as a compressed representation of the second and third script (<2→3→END>). Using its own previous output in this way, the memory encoder creates compressed representations of stories that serve as memory cues during story recall. This figure depicts step 3 of the encoding process illustrated in Figure S2.
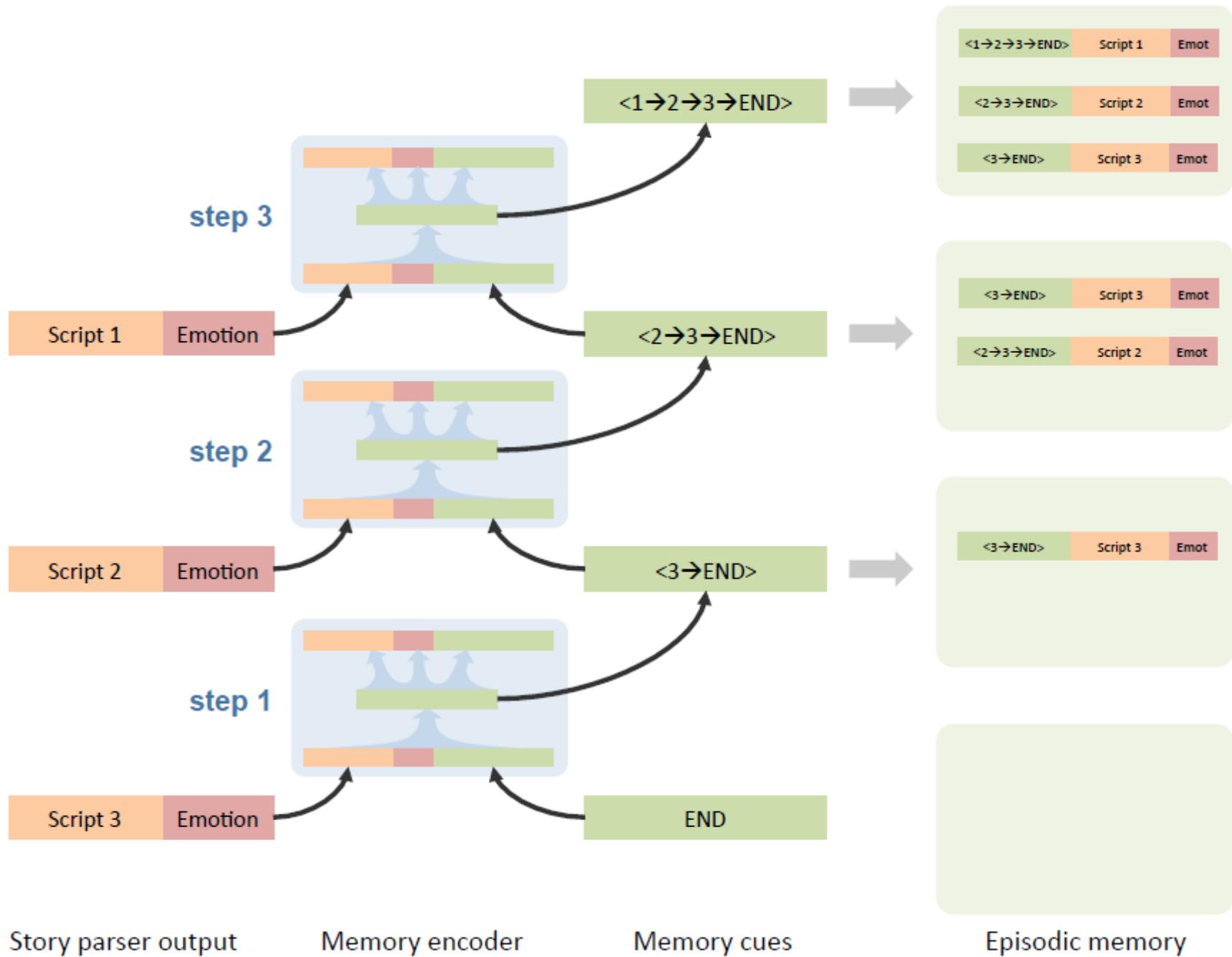
Story parser output          Memory encoder          Memory cues          Episodic memory

**Figure S2.** During the memory encoding process, each script of a story is paired with a memory cue, transforming the output of the story parser (left column) into content-addressable episodic memory traces (right column).  Each script's memory cue is a compressed version of the remaining story, and represents DISCERN's discourse plan at that point (e.g. the cue for script 2 is the compressed version of scripts 2 and 3, denoted by <2→3→END>).  The memory encoder builds these cues by stepping backwards (from bottom to top) through the scripts of a story, at each step creating a memory cue by combining a script with the memory cue produced previously.  In this manner, stories of variable length can be compressed into a single distributed memory representation, leading to cognitive capacities typical of connectionist models.
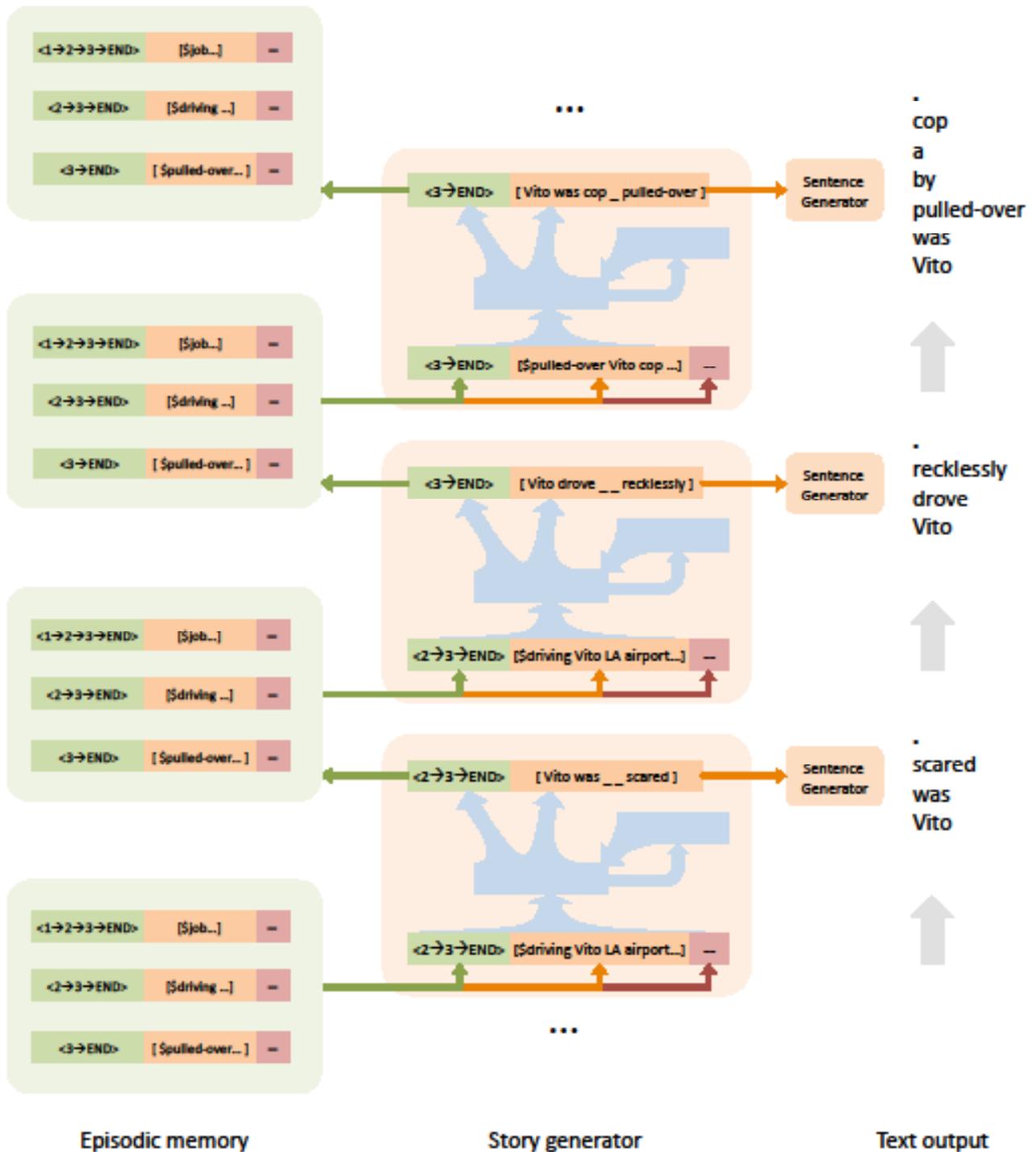
**Figure S3.** During story recall, the story generator steps through each sentence of the story, and accesses each memory trace encoding it. Three consecutive snapshots of the story generator's input and output are shown during the switch from the second ($driving) to the

third (`$pulled-over`) script of the story illustrated in the body of the report.  Time flows from

bottom to top. *Bottom:* DISCERN reproduces the sentence "`Vito is scared`" in the second

script of the story.  The story generator produces a representation of the sentence, which is then

passed on to the sentence generator (to the right).  Additionally, it produces a memory cue that

is used to retrieve the next input memory trace from episodic memory (on the left).  In this case,

the same memory trace as before is retrieved, since the script is not yet finished.

*Middle:* DISCERN produces the last sentence of the script, "`Vito drives recklessly`."

The memory cue changes, and the memory trace for the third ($pulled-over) script is retrieved.

*Top:* Using the retrieved memory trace, DISCERN now starts to reproduce the third script.  By

switching memory cues successively in this manner, the story generator can step through each

script in the correct order.  Scripts trigger subsequent scripts within a single story as is

commonly done in symbolic script-processing systems (20).  In DISCERN, this model of

narrative structure is given a subsymbolic connectionist implementation.

**Table S1.** Comparison of patients with and without fixed delusions with narrative organization completing seven day delayed story-recall.

| | Age[1] | Gender (M/F) | Parental education (grades)[1] | WAIS scaled vocabulary score[1] | SAPS positive thought disorder score[1,2] | Alogia score[1] | Antipsychotic drug treatment[3] |
|---|---|---|---|---|---|---|---|
| Patients with definite fixed narrative delusions (FND+; N = 27) | 41.2 (9.4) | (12/15) | 16.0 (8.7) | 9.6 (3.9) | 1.6 (1.2) | 1.2 (1.1) | 14/4/4/5 |
| Patients with questionable or absent evidence of fixed narrative delusions (FND-; N = 10) | 38.8 (10.4) | (4/6) | 12.7 (2.0) | 10.9 (6.2) | 1.7 (1.7) | 1.4 (1.2) | 7/0/0/3 |
| Significance test (two-tailed) | $t(55) = 0.71$ | $\chi^2 = 0.06$ | $t(53) = 1.12^4$ | $t(55) = .78$ | $t(55) = .30$ | $t(55) = .50$ | $\chi^2 = 3.8$, df = 3, $p = 0.28$ |

[1] mean (standard deviation).

[2] SAPS thought disorder of 1=questionable, 2=mild.

[3] first-generation (FG) /second generation (SG) /FG+SG/ SGx2.

[4] data missing for two subjects.

F, female; M, male; SAPS, Scale for the Assessment of Positive Symptoms; WAIS, Wechsler Adult Intelligence Scale.

## *References*

1. Elman JL (1990): Finding structure in time. *Cog Sci.* 14:179-211.

2. Spiegel R, McLaren IP (2006): Associative sequence learning in humans. *J Exp Psychology: Animal Behav Process* ;32:150-63.

3. Elman JL (1993): Learning and development in neural networks: the importance of starting small. *Cognition* 48:71-99.

4. Rohde DL, Plaut DC (1999): Language acquisition in the absence of explicit negative evidence: how important is starting small? *Cognition* 72:67-109.

5. Hoffman RE, McGlashan TH (1997): Synaptic elimination, neurodevelopment, and the mechanism of hallucinated "voices" in schizophrenia. *Am J Psychiatry* 154:1683-1689.

6. Valle-Lisboa JC, Reali F, Anastasia H, Mizraji E (2005): Elman topology with sigma-pi units: an application to the modeling of verbal hallucinations in schizophrenia. *Neural Networks* 18:863-77.

7. Miikkulainen R (1993): *Subsymbolic natural language processing.* MIT Press: Cambridge, MA.

8. Miikkulainen R, Dyer MG (1991): Natural language processing with modular PDP networks and distributed lexicon. *Cogn Sci* 13:343-399.

9. Pollack J (1990): Recursive distributed representations. *Art Intell.* 46:159–216.

10. Fox Tree JE (2000): In *Aspects of language production*, L. Wheeldon, Ed., Psychology Press: New York, NY, pp. 375-406.

11. Andreasen NC (1987): *Comprehensive assessment of symptoms and history*. University of Iowa: Iowa City, IA.

12. First MB, Spitzer RL, Gibbon M, Williams JBW (2002): *Structured Clinical Interview for DSM-IV-TR Axis I Disorders - Non-patient Edition*. New York Psychiatric Institute: New York, NY.

13. Davis JM (1974): Dose equivalence of the antipsychotic drugs. *J Psychiatric Res* 11:65-69.

14. Woods SW (2003): Chlorpromazine equivalent dosages for the newer atypical antipsychotics. *J Clin Psychiatry* 64:663-667.

15. Centorrino F, Eakin M, Bahk W-M, Kelleher JP, Goren J, Salvatore P, Egli S, Baldessarini RJ (2002): Inpatient antipsychotic drug use in 1998, 1993, and 1989. *Am J Psychiatry.* 159:1932-1935.

16. Wechsler D (1987): *Manual for the Wechsler memory scale – revised.* The Psychological Corporation: San Antonio, TX.

17. Cerf B (1993): The gift. In *Chicken Soup for the Soul.* J. Carfield, M.V. Hansen, Eds. Health Communications, Inc.: Deerfield Beach, FL.

18. Marchiori D, Warglien M (2008): Predicting human interactive learning by regret-driven neural networks. *Science* 319:1111-1113.

19. Friston KJ (1998): The disconnection hypothesis. *Schizophr Res* 30(2):115-25.

20. Schank RC (1999): *Dynamic memory revisited.* 2nd Edition. Cambridge University Press: New York, NY.