



Information Compression, Multiple Alignment, and the Representation and Processing of Knowledge in the Brain

J. Gerard Wolff*

CognitionResearch.org, Menai Bridge, UK

OPEN ACCESS

Edited by:

Asim Roy,
Arizona State University, USA

Reviewed by:

Jonathan C. W. Edwards,
University College London, UK
Luis C. Lamb,
Federal University of Rio Grande do
Sul, Brazil

*Correspondence:

J. Gerard Wolff
jgw@cognitionresearch.org

Specialty section:

This article was submitted to
Cognition,
a section of the journal
Frontiers in Psychology

Received: 18 May 2016

Accepted: 29 September 2016

Published: 03 November 2016

Citation:

Wolff JG (2016) Information
Compression, Multiple Alignment, and
the Representation and Processing of
Knowledge in the Brain.
Front. Psychol. 7:1584.
doi: 10.3389/fpsyg.2016.01584

The *SP theory of intelligence*, with its realization in the *SP computer model*, aims to simplify and integrate observations and concepts across artificial intelligence, mainstream computing, mathematics, and human perception and cognition, with information compression as a unifying theme. This paper describes how abstract structures and processes in the theory may be realized in terms of neurons, their interconnections, and the transmission of signals between neurons. This part of the SP theory—*SP-neural*—is a tentative and partial model for the representation and processing of knowledge in the brain. Empirical support for the SP theory—outlined in the paper—provides indirect support for *SP-neural*. In the abstract part of the SP theory (*SP-abstract*), all kinds of knowledge are represented with *patterns*, where a pattern is an array of atomic symbols in one or two dimensions. In *SP-neural*, the concept of a “pattern” is realized as an array of neurons called a *pattern assembly*, similar to Hebb’s concept of a “cell assembly” but with important differences. Central to the processing of information in *SP-abstract* is information compression via the matching and unification of patterns (ICMUP) and, more specifically, information compression via the powerful concept of *multiple alignment*, borrowed and adapted from bioinformatics. Processes such as pattern recognition, reasoning and problem solving are achieved via the building of multiple alignments, while unsupervised learning is achieved by creating patterns from sensory information and also by creating patterns from multiple alignments in which there is a partial match between one pattern and another. It is envisaged that, in *SP-neural*, short-lived neural structures equivalent to multiple alignments will be created via an inter-play of excitatory and inhibitory neural signals. It is also envisaged that unsupervised learning will be achieved by the creation of pattern assemblies from sensory information and from the neural equivalents of multiple alignments, much as in the non-neural SP theory—and significantly different from the “Hebbian” kinds of learning which are widely used in the kinds of artificial neural network that are popular in computer science. The paper discusses several associated issues, with relevant empirical evidence.

Keywords: multiple alignment, cell assembly, information compression, unsupervised learning, artificial intelligence

1. INTRODUCTION

The *SP theory of intelligence*, and its realization in the *SP computer model*, is a unique attempt to simplify and integrate observations and concepts across artificial intelligence, mainstream computing, mathematics, and human perception and cognition. The name “SP” derives from the central importance in the theory of information compression, something that may be seen as a process of maximizing the *Simplicity* of a body of information, by removing information that is repeated, whilst retaining as much as possible of its non-repeated expressive *Power*. Also, the theory itself may be seen to compress empirical information by combining simplicity in the theory with wide-ranging explanatory and descriptive power.

This paper, which draws on Wolff (2006, chapter 11) with revisions and updates, describes how abstract structures and processes in the SP theory may be realized in terms of neurons, their interconnections, and the transmission of impulses between neurons. This part of the SP theory—called *SP-neural*—may be seen as a tentative and partial theory of the representation and processing of knowledge in the brain. As such, it may prove useful as a source of ideas for theoretical and empirical investigations in the future. For the sake of clarity, the abstract parts of the theory, excluding SP-neural, will be referred to as *SP-abstract*.

It is envisaged that SP-neural will be further developed in the form of a computer model. As with the existing computer model of SP-abstract (which, unless otherwise stated, will be referred to as “the SP computer model”), the development of the new computer model of SP-neural will help to guard against vagueness in the theory, it will serve as a means of testing ideas to see whether or not they work as anticipated, and it will be a means of demonstrating what the model can do, and validating it against empirical data.

The next section says something about the theoretical orientation of this research. Then SP-abstract will be described briefly as a foundation for the several sections that follow, which describe aspects of SP-neural and associated issues.

2. THEORETICAL ORIENTATION

Cosmologist John Barrow has written that “Science is, at root, just the search for compression in the world” (Barrow, 1992, p. 247), an idea which may be seen to be equivalent to Occam’s Razor because, in accordance with the remarks above about the name “SP” and the theory itself, a good theory should combine conceptual *Simplicity* with descriptive or explanatory *Power*.

This works best when the range of phenomena to be described or explained is large. But this has not always been observed in practice: Newell (1973, p. 303) urged researchers in psychology to address “a genuine slab of human behavior”; and McCorduck (2004, pp. 417, 424) has described how research in artificial intelligence became fragmented into many narrow sub-fields.

In the light of these observations, and in the spirit of research on “unified theories of cognition” (Newell, 1990)

and “artificial general intelligence¹,” the SP programme of research has attempted to simplify and integrate observations and concepts across a broad canvass, resisting the temptation to concentrate only on one or two narrow areas.

3. SP-ABSTRACT IN BRIEF

As a basis for the description of SP-neural, this section provides a brief informal account of SP-abstract. The theory is described most fully in Wolff (2006) and quite fully but more briefly in Wolff (2013). Details of other publications in the SP programme, most of them with download links, are shown on (<http://www.cognitionresearch.org/sp.htm>).

3.1. Origins and Foundations of the SP Theory

The origins of SP theory are mainly in a body of research by Attneave (1954) and Barlow (1959, 1969) and others suggesting that much of the workings of brains and nervous systems may be understood as compression of information, and my own research on language learning (summarized in Wolff, 1988) suggesting that, to a large extent, the learning of language may be understood in the same terms. There is more about the foundations of the theory in Wolff (2014d).

3.2. Elements of SP-Abstract

In SP-abstract, all kinds of knowledge are represented with *patterns*, where a pattern is an array of atomic *symbols* in one or two dimensions. At present, the SP computer model² works only with 1D patterns but it is envisaged that the model will be generalized to work with 2D patterns. In this connection, a “symbol” is simply a “mark” that can make a yes/no match with any other symbol—no other result is permitted.

In most of the examples shown in this paper, symbols are shown as alphanumeric characters or short strings of characters but, when the SP system is used to model biological structures and processes, such representations may be interpreted as low-level elements of perception such as formants or formant ratios in the case of speech or lines and junctions between lines in the case of vision (see also Section 4.2).

To help cut through mathematical complexities associated with information compression, the SP system—SP-abstract and its realization in the SP computer model—is founded on a simple, “primitive” idea: that information may be compressed by finding full or partial matches between patterns and merging or “unifying” the parts that are the same. This principle—“Information Compression via the Matching and Unification of Patterns” (ICMUP)—provides the foundation for a powerful concept of *multiple alignment*, borrowed and adapted from bioinformatics. The multiple alignment concept, outlined in Section 3.5, below, is itself central in the workings of SP-abstract

¹See, for example, “Artificial General Intelligence”, *Wikipedia*, <http://bit.ly/1ZxCQP0>, retrieved 2016-01-19.

²The current version of the SP computer model is SP71, the source code for which may be downloaded via a link near the bottom of www.cognitionresearch.org/sp.htm. This version of the computer model is very similar to SP70, described in Wolff (2006, Sections 3.9.2, 9.2).

and is the key to versatility and adaptability in the SP system. It has the potential to be as significant for the understanding of “intelligence” in a broad sense as is DNA for biological sciences.

3.3. SP Patterns, Multiple Alignment, and the Representation and Processing of Knowledge

In themselves, SP patterns are not very expressive. But in the multiple alignment framework (Section 3.5) they become a very versatile medium for the representation of diverse forms of knowledge. And the building of multiple alignments, together with processes for unsupervised learning (Sections 3.4, 3.7), has proved to be a powerful means of modeling diverse aspects of intelligence.

The two things together—SP patterns and multiple alignment—have the potential to be a “Universal Framework for the Representation and Processing of Diverse Kinds of Knowledge” (UFK), as discussed in Wolff (2014c, Section III).

An implication of these ideas is that there would not, for example, be any difference between the representation and processing of non-syntactic cognitive knowledge and the representation and processing of the syntactic forms of natural language. A framework that can accommodate both kinds of knowledge is likely to facilitate their seamless integration, as discussed in Section 3.8.2.

3.4. Early Stages of Learning

The SP theory is conceived as a brain-like system that receives *New* patterns via its “senses” and stores some or all of them, in compressed form, as *Old* patterns. In broad terms, this is how the system learns.

In the SP system, all learning is “unsupervised³,” meaning that it does not depend on assistance by a “teacher,” the grading of learning materials from simple to complex, or the provision of “negative” examples of concepts to be learned—meaning examples that are marked as “wrong” (*cf.* Gold, 1967). Notwithstanding the importance of schools and colleges, it appears that most human learning is unsupervised. Other kinds of learning, such as “supervised” learning (learning from labeled examples)⁴, or “reinforcement” learning (learning with carrots and sticks)⁵, may be seen as special cases of unsupervised learning (Wolff, 2014b, Section V).

At the beginning of processing by the system, when the repository of Old patterns is empty⁶, New patterns are stored as they are received but with the addition of system-generated “ID” symbols at the beginning and end. For example, a New pattern like “t h e b i g h o u s e” would be stored as an Old pattern like “A 1 t h e b i g h o u s e #A.” Here,

³See “Unsupervised learning,” *Wikipedia*, bit.ly/22nEPL2, retrieved 2016-03-17.

⁴See “Supervised learning,” *Wikipedia*, bit.ly/1nR4ybK, retrieved 2016-03-17.

⁵See “Reinforcement learning,” *Wikipedia*, bit.ly/1R0RoDv, retrieved 2016-03-17.

⁶Although it is likely that, contrary to what Noam Chomsky and others have suggested, a newborn child does *not* have any kind of detailed knowledge of the structure of natural language, it is likely he or she does have inborn knowledge such as how to suck milk from a breast. In this respect (and others), the SP theory, insofar it is seen as a model of human cognition, is not entirely accurate.

the lower-case letters are atomic symbols that may represent actual letters but could represent basic elements of speech (such as formant ratios or formant transitions), or basic elements of vision (such as lines or corners), and likewise with other sensory data.

Later, when some Old patterns have been stored, the system may start to recognize full or partial matches between New and Old patterns. If a New pattern is exactly the same as an Old pattern (excluding the ID-symbols), then frequency measures for that pattern and its constituent symbols are incremented. These measures, which are continually updated at all stages of processing, have an important role to play in calculating probabilities of structures and inferences and in guiding the processes of building multiple alignments (Section 3.5) and unsupervised learning.

With partial matches, the system will form multiple alignments like the one shown in **Figure 1**, with a New pattern in row 0 and an Old pattern in row 1.

From a partial match like this, the system creates Old patterns from the parts that match each other and from the parts that don’t. Each newly-created Old pattern will be given system-generated ID-symbols. The result in this case would be patterns like these: “B 1 t h e #B,” “C 1 h o u s e #C,” “D 1 s m a l l #D,” and “D 2 b i g #D.” In addition, the system forms an abstract pattern like this: “E 1 B #B D #D C #C #E” which records the sequence [“B 1 t h e #B,” (“D 1 s m a l l #D” or “D 2 b i g #D”), “C 1 h o u s e #C”] in terms the ID-symbols of the constituent patterns.

Notice how “s m a l l” and “b i g” have both been given the ID-symbol “D” at their beginnings and the ID-symbol “#D” at their ends. These additions, coupled with the use of the same two ID-symbols in the abstract pattern “E 1 B #B D #D C #C #E” has the effect of assigning “s m a l l” and “b i g” to the same syntactic category, which looks like the beginnings of the “adjective” part of speech.

The overall result in this example is a collection of SP patterns that functions as a simple grammar to describe the phrases *the small house* and *the big house*.

In practice, the SP computer model may form many other multiple alignments, patterns and grammars which are much less tidy than the ones shown. But, as outlined in Sections 3.5, 3.7, the system is able to home in on structures that are “good” in terms of information compression.

As we shall (see Sections 3.5, 3.8.1, and 6), SP patterns, within the SP system, are remarkably versatile and expressive, with at least the power of context-sensitive grammars (Wolff, 2006, Chapter 5).

| | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|----|
| 0 | t | h | e | | s | m | a | l | l | h | o | u | s | e | 0 |
| | | | | | | | | | | | | | | | |
| 1 | A | 1 | t | h | e | b | i | g | | h | o | u | s | e | #A |

FIGURE 1 | A multiple alignment produced by the SP computer model showing a partial match between a New pattern (in row 0) and an Old pattern (in row 1).

directly with little modification except for the addition of system-generated ID-symbols. Later, when there are more Old patterns in store, the system begins to create Old patterns from partial matches between New and Old patterns. Part of this process is the creation of abstract patterns that describe sequences of lower-level patterns.

As the system begins to create abstract patterns, it will also begin to form multiple alignments like the one shown in **Figure 2**. And, as it begins to form multiple alignments like that, it will also begin to form code patterns, as described in Section 3.6.

At all stages of learning, but most prominent in the later stages, is a process of inferring one or more *grammars* that are “good” in terms of their ability to encode economically all the New patterns that have been presented to the system. Here, a “grammar” is simply a collection of SP patterns⁸.

Inferring grammars that are good in terms of information compression is, like the building multiple alignments, a stage-by-stage process of heuristic search through the vast abstract space of alternatives, discarding “bad” alternatives at each stage, and retaining a few that are “good.” As with the building of multiple alignments, the search aims to find solutions that are “good enough,” and not necessarily perfect. These kinds of heuristic search may be performed by means of genetic algorithms, simulated annealing, and other heuristic techniques.

It is envisaged that the SP computer model will be developed so that, in this later phase of learning, learning processes will be applied to code patterns as well as to New patterns. It is anticipated that this may overcome two weaknesses in the SP computer model as it is now: that, while it forms abstract patterns at the highest level, it does not form abstract patterns at intermediate levels; and that it does not recognize discontinuous dependencies in knowledge (Wolff, 2013, Section 3.3).

In Wolff (2006, Chapter 9), there is a much fuller account of unsupervised learning in the SP computer model.

3.8. Evaluation of SP-Abstract

The SP theory in its abstract form may be evaluated in terms of “simplicity” and “power” of the theory itself (discussed in Section 3.8.1 next), in terms of its potential to promote simplification and integration of structures and functions in natural or artificial systems that conform to the theory (Section 3.8.2 below), and in comparison with other AI-related systems.

3.8.1. Simplicity and Power

In terms of the principles outlined in Section 2, the SP system, with multiple alignment center stage, scores well. One relatively simple framework has strengths and potential in the representation of several different kinds of knowledge, in several different aspects of AI, and it has several potential benefits and applications:

- *Representation and processing of diverse kinds of knowledge.* The SP system (SP-abstract) has strengths and potential in the representation and processing of: class hierarchies

⁸The term “grammar” has been adopted partly because of the origins of the SP system in research on the learning of natural language (Wolff, 1988) and partly because the term has come to be used in areas outside computational linguistics, such as pattern recognition.

and heterarchies, part-whole hierarchies and heterarchies, networks and trees, relational knowledge, rules used in several kinds of reasoning, patterns with pattern recognition, images with the processing of images (Wolff, 2014a), structures in planning and problem solving, structures in three dimensions (Wolff, 2014a, Section 6), knowledge of sequential and parallel procedures (Wolff, 2014b, Section IV-H). It may also provide an interpretive framework for structures and processes in mathematics (Wolff, 2014d, Section 10).

There is a fuller summary in Wolff (2014c, Section III-B) and much more detail in Wolff (2006, 2013).

- *Strengths and potential in AI.* The SP theory has things to say about several different aspects of AI, as described most fully in Wolff (2006) and more briefly in Wolff (2013). In addition to its capabilities in the parsing of natural language, described above, the SP system has strengths and potential in the production of natural language, the representation and processing of diverse kinds of semantic structures, the integration of syntax and semantics, fuzzy pattern recognition, recognition at multiple levels of abstraction, computer vision and modeling aspects of natural vision (Wolff, 2014a), information retrieval, planning, problem solving, and several kinds of reasoning (one-step “deductive” reasoning; abductive reasoning; reasoning with probabilistic decision networks and decision trees; reasoning with “rules”; nonmonotonic reasoning and reasoning with default values; reasoning in Bayesian networks with “explaining away”; causal diagnosis; reasoning which is not supported by evidence; and inheritance of attributes in an object-oriented class hierarchy or heterarchy). There is also potential for spatial reasoning (Wolff, 2014b, Section IV-F.1) and what-if reasoning (Wolff, 2014b, Section IV-F.2). The system also has strengths and potential in unsupervised learning (Wolff, 2006, Chapter 9).
- *Many potential benefits and applications.* Potential benefits and applications of the SP system include: helping to solve nine problems associated with big data (Wolff, 2014c); the development of intelligence in autonomous robots, with potential for gains in computational efficiency (Wolff, 2014b); the development of computer vision (Wolff, 2014a); it may serve as a versatile database management system, with intelligence (Wolff, 2007); it may serve as an aid in medical diagnosis (Wolff, 2006); and there are several other potential benefits and applications, some of which are described in Wolff (2014e).

In short, the SP theory, in accordance with Occam’s Razor, demonstrates a favorable combination of simplicity and power across a broad canvass. As in other areas of science, this should increase our confidence in the validity and generality of the theory.

3.8.2. Simplification and Integration

Closely related to simplicity and power in the SP theory are two potential benefits arising from the use of one simple format (SP patterns) for all kinds of knowledge and one relatively simple framework (chiefly multiple alignment) for the processing of all kinds of knowledge:

- *Simplification.* Those two features (one simple format for knowledge and one simple framework for processing it) can mean substantial simplification of natural systems (brains) and artificial systems (computers) for processing information. The general idea is that one relatively simple system can serve many different functions. In natural systems, there is a potential advantage in terms of natural selection, and in artificial systems there are potential advantages in terms of costs.
- *Integration.* The same two features are likely to facilitate the seamless integration of diverse kinds of knowledge and diverse aspects of intelligence—pattern recognition, several kinds of reasoning, unsupervised learning, and so on—in any combination, in both natural and artificial systems. It appears that that kind of seamless integration is a key part of the versatility and adaptability of human intelligence and that it will be essential if we are to achieve human-like versatility and adaptability of intelligence in artificial systems.

With regard to the seamless integration of diverse kinds of knowledge, this is clearly needed in the understanding and production of natural language. To understand what someone is saying or writing, we obviously need to be able to connect words and syntactic structures with their non-syntactic meanings, and likewise, in reverse, when we write or speak to convey some meaning.

This has not yet been explored in any depth with the SP-abstract conceptual framework but preliminary trials with the SP computer model suggest that it is indeed possible to define syntactic-semantic structures in a set of SP patterns and then, with those patterns playing the role of Old patterns, to analyse a sample sentence and to derive its meanings (Wolff, 2006, Section 5.7, Figure 5.18), and, in a separate exercise with the same set of Old patterns, to derive the same sentence from a representation of its meanings (Wolff, 2006, Figure 5.19).

3.8.3. Distinctive Features and Advantages of the SP System Compared with Other AI-Related Systems

In several publications, such as Wolff (2006, 2007, 2014e), potential benefits and applications of the SP system have been described.

More recently, it has seemed appropriate to say what distinguishes the SP system from other AI-related systems and, more importantly, to describe advantages of the SP system compared AI-related alternatives. Those points have now been set out in some detail in *The SP theory of intelligence: its distinctive features and advantages* (Wolff, 2016). Of particular relevance to this paper are the several advantages of the SP system compared with systems for deep learning in artificial neural networks (Wolff, 2016, Section V).

Since many AI-related systems may also be seen as models of cognitive structures and processes in brains, Wolff (2016) may also be seen to demonstrate the relative strength of the SP system in modeling aspects of human perception and cognition.

In this connection, the SP system appears to have some advantages compared with concepts developed in research in

“neural-symbolic computation,” described in d’Avila Garcez et al. (2015), de Penning et al. (2011), d’Avila Garcez et al. (2009), Komendantskaya et al. (2007), and d’Avila Garcez (2007) amongst other publications. The main apparent advantages are:

- *The AI scope of the SP system.* The scope of SP-abstract in AI, meaning the range of AI-related capabilities where it has strengths and potential (summarized in Section 3.8.1), appears to be greater than the range of AI-related capabilities considered in research on neural-symbolic computation. There is potential for SP-neural to inherit that same wide scope.
- *Problems with deep learning in artificial neural networks, and potential SP solutions.* As mentioned above, the SP system has the potential to overcome several problems with deep learning in artificial neural networks (Wolff, 2016, Section V).

4. INTRODUCTION TO SP-NEURAL

As we have seen in Section 3, SP-abstract is a relatively simple system with descriptive and explanatory power across a wide range of observation and phenomena in artificial intelligence and related areas. How can such a system have anything useful to say about the extraordinary complexity of brains and nervous systems, both in their structure and in their workings?

An answer in brief is that SP-neural—a realization of SP-abstract in terms of neurons, their interconnections, and the transmission of impulses between neurons—may help us to interpret neural structures and processes in terms of the relatively simple concepts in SP-abstract. To the extent that this is successful, it may—like any good theory in any field—help us to understand empirical phenomena in our area of interest, it may help us to make predictions, and it may suggest lines of investigation.

It is anticipated that SP-neural will work in broadly the same way as SP-abstract, but the characteristics of neurons and their interconnections raise some issues that do not arise in SP-abstract and its realization in the SP computer model. These issues will be discussed at appropriate points in this and subsequent sections.

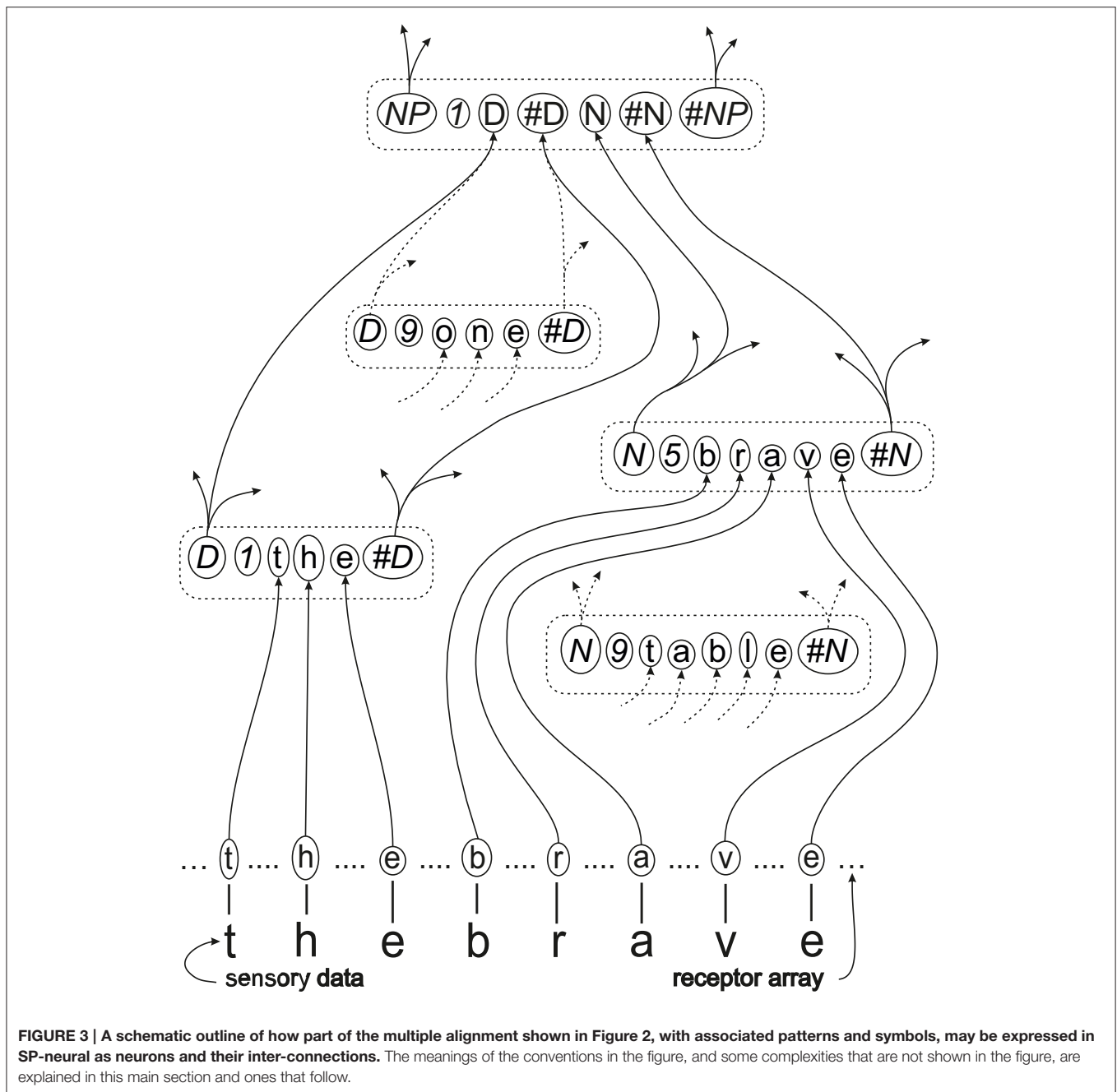
This section introduces SP-neural in outline, and sections that follow describe aspects of the theory in more detail, drawing where necessary on aspects of SP-abstract that have been omitted from or only sketched in Section 3.

4.1. Sensory Data and the Receptor Array

Figure 3, to be discussed in this and the following subsections, shows in outline how a portion of the multiple alignment shown in Figure 2, may be realized in SP-neural, with associated patterns and symbols.

In the figure, “sensory data” at the bottom means the visual, auditory or tactile data entering the system which, in this case, corresponds with the phrase “t h e b r a v e.” In a more realistic illustration, the sensory data would be some kind of analog signal. Here, the letters are intended to suggest the kinds of low-level perceptual primitives outlined below.

It is envisaged that, with most sensory modalities, the receptor array would be located in the primary sensory cortex. Of course,



a lot of processing goes on in the sense organs and elsewhere between the sense organs and the primary sensory cortices. But it seems that most of this early processing is concerned with the identification of the perceptual primitives just mentioned.

As with SP-abstract, it is anticipated that SP-neural will, at some stage, be generalized to accommodate patterns in two dimensions, such as visual images, and then the sensory data may be received in two dimensions, as in the human eye.

Between the sensory data and the *receptor array* (above it in the figure), there would be, first, cells that are specialized to receive particular kinds of input (auditory, visual, tactile

etc.). These send signals to neurons that encode the sensory data as *neural symbols*, the neural equivalents of “symbols” in SP-abstract. In the receptor array, each letter enclosed in a solid ellipse represents a neural symbol, expressed as a single neuron or, more likely, a small cluster of neurons. As we shall see Section 5.1, the reality is more complex, at least in some cases.

In vision, neural symbols in the receptor array would represent such low-level features as lines, corners, colors, and the like, while in speech perception, they would represent such things as formants, formant ratios and transitions, plosive and fricative sounds, and so on. Whether or how the SP concepts can

be applied in the discovery or identification of features like these is an open question (Wolff, 2013, Section 3.3). For now, we shall assume that they can be identified and can be used in the creation and use of higher-level structures.

4.2. Pattern Assemblies

In the rest of **Figure 3**, each broken-line rectangle with rounded corners represents a *pattern assembly*—corresponding to a “pattern” in SP-abstract. The word “assembly” has been adopted within the expression “pattern assembly” because the concept is quite similar to Hebb’s concept of a “cell assembly”—a cluster of neurons representing a concept or other coherent mental entity. Differences between Hebb’s concept of a cell assembly and the SP concept of a pattern assembly are described in the Appendix.

Within each pattern assembly, as represented in the figure, each character or group of characters enclosed in a solid-line ellipse represents a *neural symbol* which, as already mentioned, corresponds to a “symbol” in SP-abstract. As with neural symbols in the receptor array, it is envisaged that each neural symbol would comprise a single neuron or, more likely, a small cluster of neurons.

It is supposed that, within each pattern assembly, there are lateral connections between neural symbols—but these are not shown in the figure.

It is envisaged that most pattern assemblies would represent knowledge that is learned and not inborn, and would be located mainly outside the primary sensory areas of the cortex, in other parts of the sensory cortices. Pattern assemblies that integrate two or more sensory modalities may be located in “association” areas of the cortex.

Research with fMRI recordings from volunteers (Huth et al., 2016) has revealed “semantic maps” that “show that semantic information is represented in rich patterns that are distributed across several broad regions of cortex. Furthermore, each of these regions contains many distinct areas that are selective for particular types of semantic information, such as people, numbers, visual properties, or places. We also found that these cortical maps are quite similar across people, even down to relatively small details”⁹. Of course, this research says nothing about whether or not the knowledge is represented with pattern assemblies and their interconnections. But it does apparently confirm that knowledge is stored in several regions of the cortex and throws light on how it is organized.

Although most parts of the mammalian cerebral cortex has six layers and many convolutions, it may be seen, topologically, as a sheet which is very much broader and wider than it is thick. Correspondingly, it is envisaged that 1D and 2D pattern assemblies will be largely “flat” structures, rather like writing or pictures on a sheet of paper. That said, it is quite possible, indeed likely, that pattern assemblies would take advantage of two or more layers of the cortex, not just one.

⁹From the website of the Gallant Lab at UC Berkeley, retrieved 2016-05-02, <http://bit.ly/1WvVlhX>. See also “Brain ‘atlas’ of words revealed,” *BBC News*, 2016-04-27, bbc.in/1SGESLz.

Incidentally, since 2D SP patterns may provide a basis for 3D models, as described in Wolff (2014a, Sections 6.1, 6.2), flat neural structures in the cortex may serve to represent 3D concepts.

4.3. Connections between Pattern Assemblies

In **Figure 3**, the solid or broken lines that connect with neural symbols represent axons, with arrows representing the direction of travel of neural impulses. Where two or more connections converge on a neural symbol, we may suppose that, contrary to the simplified way in which the convergence is shown in the figure, there would be a separate dendrite for each connection.

Axons represented with solid lines are ones that would be active when the multiple alignment in **Figure 2** is in the process of being identified. Broken-line connections show a few of the many other possible connections.

As mentioned in Section 4.2, it is envisaged that there would be one or more neural connections between neighboring neural symbols within each pattern assembly but these are not marked in the figure.

Compared with what is shown in the figure, it likely that, in reality, there would be more “levels” between basic neural symbols in the receptor array and ID-neural-symbols representing pattern assemblies for relatively complex entities like the words “one,” “brave,” “the,” and “table,” as shown in the figure.

In this connection, it is perhaps worth emphasizing that, as with the modeling of hierarchical structures in multiple alignments (Section 3.5), while pattern assemblies may form “strict” hierarchies, this is not an essential feature of the concept, and it is likely that many neural structures formed from pattern assemblies may be only loosely hierarchical or not hierarchical at all.

4.4. SP-Neural, Quantities of Knowledge, and the Size of the Brain

Given the foregoing account of how knowledge may be represented in the brain, a question that arises is “Are there enough neurons in the brain to store what a typical person knows?” This is a difficult question to answer with any precision but an attempt at an answer, described in Wolff (2006, Section 11.4.9), reaches the tentative conclusion that there are. In brief:

- Given that estimates of the size of the human brain range from 10^{10} up to 10^{11} neurons,¹⁰ we may estimate, via calculations given in Wolff (2006, Section 11.4.9), that the “raw” storage capacity of the brain is between approximately 1000 and 10,000 MB.
- Given a conservative estimate that, using SP compression mechanisms, compression by a factor of 3 may be achieved across all kinds of knowledge, our estimates of the storage capacity of the brain will range from about 3000 MB up to about 30,000 MB.

¹⁰This is consistent with another estimate, not quoted in Wolff (2006, Section 11.4.9), that there may be as many as 86 billion neurons in the human brain (Herculano-Houzel, 2012).

- Assuming: (1) That the average person knows only a relatively small proportion of what is contained in the *Encyclopedia Britannica* (EB); (2) That the average person knows lots of “everyday” things that are *not* in the EB; (3) That the “everyday” things that we *do* know are roughly equal to the things in the EB that we *do not* know; Then (4), we may conclude that the size of the EB provides a rough estimate of the volume of information that the average person knows.
- The EB can be stored on two CDs in compressed form. Assuming that most of the space is filled, this equates to 1300 MB of compressed information or approximately 4000 MB of information in uncompressed form.
- This 4000 MB estimate of what the average person knows is the same order of magnitude as our range of estimates (3000 to 30,000 MB) of what the human brain can store.
- Even if the brain stores two or three copies of its compressed knowledge—to guard against the risk of losing it, or to speed up processing, or both—our estimate of what needs to be stored (lets say $4000 \times 3 = 12,000$ MB) is still within the 3000 to 30,000 MB range of estimates of what the brain can store.

4.5. Neural Processing

In broad terms, it is envisaged that, for a task like the parsing of natural language or pattern recognition:

1. SP-neural will work firstly by receiving sensory data and interpreting it as neural symbols in the receptor array—with excitation of the neural symbols that have been identified:
 - Excitatory signals would be sent from those excited neural symbols to pattern assemblies that can receive signals from them directly. In **Figure 3**, these would be all the pattern assemblies except the topmost pattern assembly.
 - Within each pattern assembly, excitatory signals will spread laterally via the connections between neighboring neural symbols.
 - Pattern assemblies would become excited, roughly in proportion to the number of excitatory signals they receive.
2. At this stage, there would be a process of selecting amongst pattern assemblies to identify one or two that are most excited.
3. From those pattern assemblies—more specifically, the neural ID-symbols at the beginnings and ends of those pattern assemblies—excitatory signals would be sent onwards to other pattern assemblies that may receive them. In **Figure 3**, this would be the topmost pattern assembly (that would be reached immediately after the first pass through stages 2 and 3).

As in stage 1, the level of excitation of any pattern assembly would depend on the number of excitatory signals it receives, but building up from stage to stage so that the highest-level pattern assemblies are likely to be most excited.
4. Repeat stages 2 and 3 until there are no more pattern assemblies that can be sent excitatory signals.

The “winning” pattern assembly or pattern assemblies, together with the structures below them that have, directly or indirectly, sent excitatory signals to them, may be seen as neural analogs of multiple alignments (NAMAs), and we may guess that they

provide the best interpretations of a given portion of the sensory data.

If the whole sentence, “*f o r t u n e f a v o u r s t h e b r a v e*,” is processed by SP-neural with pattern assemblies that are analogs of the SP patterns provided for the example shown in **Figure 2**, we may anticipate that the overall result would be a pattern of neural excitation that is an analog of the multiple alignment shown in that figure.

When a neural symbol or pattern assembly has been “recognized” by participating in a winning (neural) multiple alignment, we may suppose that some biochemical or physiological aspect of that structure is increased as an at least approximate measure of the frequency of occurrence of the structure, in accordance with the way in which SP-abstract keeps track of the frequency of occurrence of symbols and patterns (Section 3.4).

Some further possibilities are discussed in Sections 5, 9.

5. SOME MORE DETAIL

The bare-bones description of SP-neural in Section 4 is probably inaccurate in some respects and is certainly too simple to work effectively. This section and the ones that follow describe some other features which are likely to figure in a mature version of SP-neural, drawing on relevant empirical evidence where it is available.

5.1. Encoding of Information in the Receptor Array

With regard to the encoding of information in the receptor array, it seems that the main possibilities are these:

1. *Explicit alternatives*. For the receptor array to work as described in Section 4, it should be possible to encode sensory inputs with an “alphabet” of alternative values at each location in the array, in much the same way that each binary digit (bit) in a conventional computer may be set to have the value 0 or 1, or how a typist may enter any one of an alphabet of characters at any one location on the page. At each location in the receptor array, each option may be provided in the form of a neuron or small cluster of neurons. Here, there seem to be two main options:
 - a. *Horizontal distribution of alternatives*. The several alternatives may be distributed “horizontally,” in a plane that is parallel to the surface of the cortex.
 - b. *Vertical distribution of alternatives*. The several alternatives may be distributed “vertically” between the outer and inner surfaces of the cortex, and perpendicular to those surfaces.
2. *Implicit alternatives*. At each location there may be a neuron or small cluster of neurons that, via some kind of biochemical or neurophysiological process, may be “set” to any one of the alphabet of alternative values.
3. *Rate codes*. Something like the intensity of a stimulus may be encoded via “an interaction between [the] firing rates and the number of neurons [that are] activated by [the] stimulus.” (Squire et al., 2013, p. 503).

4. *Temporal codes.* A stimulus that varies with time may be encoded via “the time-varying pattern of activity in small groups of receptors and central neurons.” (Squire et al., 2013).

In support of option 1.a, there is evidence that neurons in the visual cortex (of cats) are arranged in columns perpendicular to the surface of the cortex, where, for example, all the neurons in a given column respond most strongly to a line at one particular angle in the field of view, that—within a “hypercolumn” containing several columns—the preferred angle increases progressively from column to column, and that there are many hypercolumns across the primary visual cortex (Barlow, 1982). “Hubel and Wiesel point out that the organization their results reveal means that each small region, about 1 mm^2 at the surface, contains a complete sequence of ocular dominance and a complete sequence of orientation preference.” (Barlow, 1982, pp. 148–149).

Leaving out the results for ocular dominance, these observations are summarized schematically in **Figure 4**. In terms of this scheme, the way in which the receptor array is shown in **Figure 3**, is a considerable simplification—each neural symbol in the receptor array in that figure should really be replaced by a hypercolumn.

With something like the intensity of a stimulus, it seems that, at least in some cases: “... activity in one particular population of somatosensory neurons ... leads the CNS to interpret it as painful stimulus” (Squire et al., 2013, p. 503), while “An entirely separate population of neurons ... would signal light pressure.” (Squire et al., 2013). Since it is likely that relevant receptors appear repeatedly across one’s skin, this appears to be another example of option 1.a.

There seems to be little evidence of encoding via option 1.b. Indeed, since the concept of a cortical column is, in effect, defined by the fact that all the neurons in any one column have the same kind of receptive field, this seems to rule out the 1.b option (see also Section 5.2).

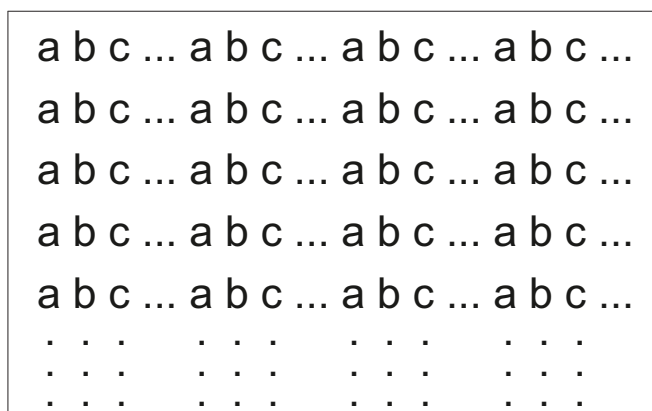


FIGURE 4 | Schematic representation of one hypercolumn in the receptor array in the cortex. Each letter represents a neural symbol that responds to a particular small pattern in the sensory data. The ellipsis, “...,” in each row and each column represents other neural symbols that would be shown in a more comprehensive representation of the given hypercolumn. Each vertical sequence of letters, all of one kind such as “a” or “b,” represents a simple column in the cortex.

But, with respect to option 2, it appears that in some cases, as noted above, the intensity of a stimulus may be encoded via the rates of firing of neurons, together with the numbers of neurons that are activated (option 3). And, since we can perceive and remember time-varying stimuli such as the stroking of a finger across one’s skin, or the rising or falling pitch of a note, some kind of temporal encoding must be available (option 4).

Here, it must be acknowledged that options 3 and 4 appear superficially to be outside the scope of the SP theory, in view of the emphasis in many examples on discrete atomic symbols. But, as we know from the success of digital recording, or indeed digital computing, any continuum may be encoded digitally, in keeping with the digital nature of the SP theory. How the SP theory may be applied to the digital encoding and processing of continua has been discussed elsewhere in relation to vision (Wolff, 2014a) and the development of autonomous robots (Wolff, 2014b).

5.2. Why Are There Multiple Neurons with the Same Receptive Fields in Columns in the Cortex?

As we have seen (Section 5.1), some aspects of vision are mediated via columns of neurons in the primary visual cortex in which each column contains many neurons with receptive fields that are all the same, all of them responding, for example, to a line in the visual field with a particular orientation.

Why, at each of several locations across the visual cortex, should there be many neurons with the same receptive field, not just one? There seem to be two possible answers to this question (and they are not necessarily mutually exclusive):

- *Encoding of sensory patterns.* If, in the receptor array, we wish to encode two or more patterns such as “m e t” and “h e m,” they need to be independent of each other, with repetition of the “e” neural symbol, otherwise there will be the possibly unwanted implication that such things as “m e m” or “h e t” are valid patterns.
- *Error-reducing redundancy.* At any given location in the receptor array, multiple instances of neurons representing a given neural symbol may help to guard against the problems that may arise if there is only one neuron at that location and if, for any reason, it becomes partially or fully disabled.

With regard to the first point, the receptor array may have a useful role to play, *inter alia*, as a short-term memory for many sensory patterns pending their longer-term storage (Section 11). In vision, for example, the receptor array may store many short glimpses of a scene, as outlined in Section 5.6, until such time as further processing may be applied to weld the many glimpses into a coherent structure (Wolff, 2014b) and to transfer that structure to longer-term memory.

5.3. The Labeled Line Principle

Section 4.5 suggests that normally, at some early stage in sensory processing, raw sensory data is encoded in terms of the excitation of neuronal symbols in a receptor array, then excited neural symbols send excitatory signals to appropriate neural symbols within pattern assemblies, and pattern assemblies that

are sufficiently excited send excitatory signals on to other pattern assemblies, and so on. As we shall see (Section 9), it is likely that, in this processing, there will also be a role for inhibitory processes.

At first sight, it may be thought that, in the same way that each location in the receptor array should provide an alphabet of alternative encodings (Section 5.1), the same should be true for the location of each neural symbol within each pattern assembly. But if a neural symbol in a pattern assembly (let's call it "NS1") receives signals only from neural symbols in the receptor array that represent a given feature, let us say, "a," then, in accordance with the "labeled line" principle (Squire et al., 2013, p. 503), NS1 also represents "a."

For most sensory modalities, this principle applies all the way from each sense organ, through the thalamus, to the corresponding part of the primary sensory cortex¹¹. It seems reasonable to suppose that the same principle will apply onwards from each primary sensory cortex into non-primary sensory cortices and non-sensory association areas.

5.4. How the Ordering or 2D Arrangement of Neural Symbols May Be Respected

In SP-neural, as in SP-abstract and the SP computer model, the process of matching one pattern with another should respect the orderings of symbols. For example, "A B C D" matched with "A B C D" should be rated more highly in terms of information compression than, for example, "A B C D" matched with "C A D B"¹².

It appears that this problem may be solved by the adoption, within SP-neural, of the following feature of natural sensory systems:

"Receptors within [the retina and skin surface] communicate with ganglion cells and those ganglion cells with central neurons in a strictly ordered fashion, such that relationships with neighbors are maintained throughout. This type of pattern, in which neurons positioned side by side in one region communicate with neurons positioned side-by-side in the next region, is called a *typographic pattern*." (Squire et al., 2013, p. 504) (emphasis in the original).

5.5. How to Accommodate the Variable Sizes of Sensory Patterns

A prominent feature of human visual perception is that we can recognize any given entity over a wide range of viewing distances, with correspondingly wide variations in the size, on the retina, of the image of that entity.

¹¹Thus, for example, "Even within one function, mappings of neurons [within the thalamus] are preserved so that there is separation of neurons providing touch information from the arm vs. from the leg and of neurons responding to low vs. high sound frequencies" (Squire et al., 2013, p. 507). Also, "Nuclei in the central pathways often contain multiple maps." but "The functional significance of multiple maps in general, however, remains to be clarified." (Squire et al., 2013).

¹²A possible exception is when one pattern is a mirror image or inversion of another, since Leonardo da Vinci, by repute, could read mirror writing as easily as ordinary writing, and it is now well established that people wearing inverting spectacles can learn quite quickly to see the world as if it was the right way up (Stratton, 1897).

For any model of human visual perception that is based on a simplistic or naive process for the matching of patterns, this aspect of visual perception would be hard to reproduce or to explain. But the SP system is different: (1) Knowledge of entities that we may recognize are always stored in a compressed form; (2) The process of recognition is a process of compressing the incoming data; (3) The overall effect is that an image of a thing to be recognized can be matched with stored knowledge of that entity, regardless of the original size of the image.

As an example, consider how the concept of an equilateral triangle (as white space bounded by three black lines all of the same length) may be stored and how an image of such a triangle may be recognized. Regarding storage, there are three main redundancies in any image of that kind of triangle: (1) The white space in the middle may be seen as repeated instances of a symbol representing a white pixel; (2) Each of the three sides of the triangle may be seen as repeated instances of a symbol representing a black pixel; and (3) There is redundancy in that the three sides of the triangle are the same.

All three sources of redundancy may be encoded recursively as suggested in Figure 5¹³, which shows a multiple alignment modeling the recognition of a one-dimensional analog of a triangle.

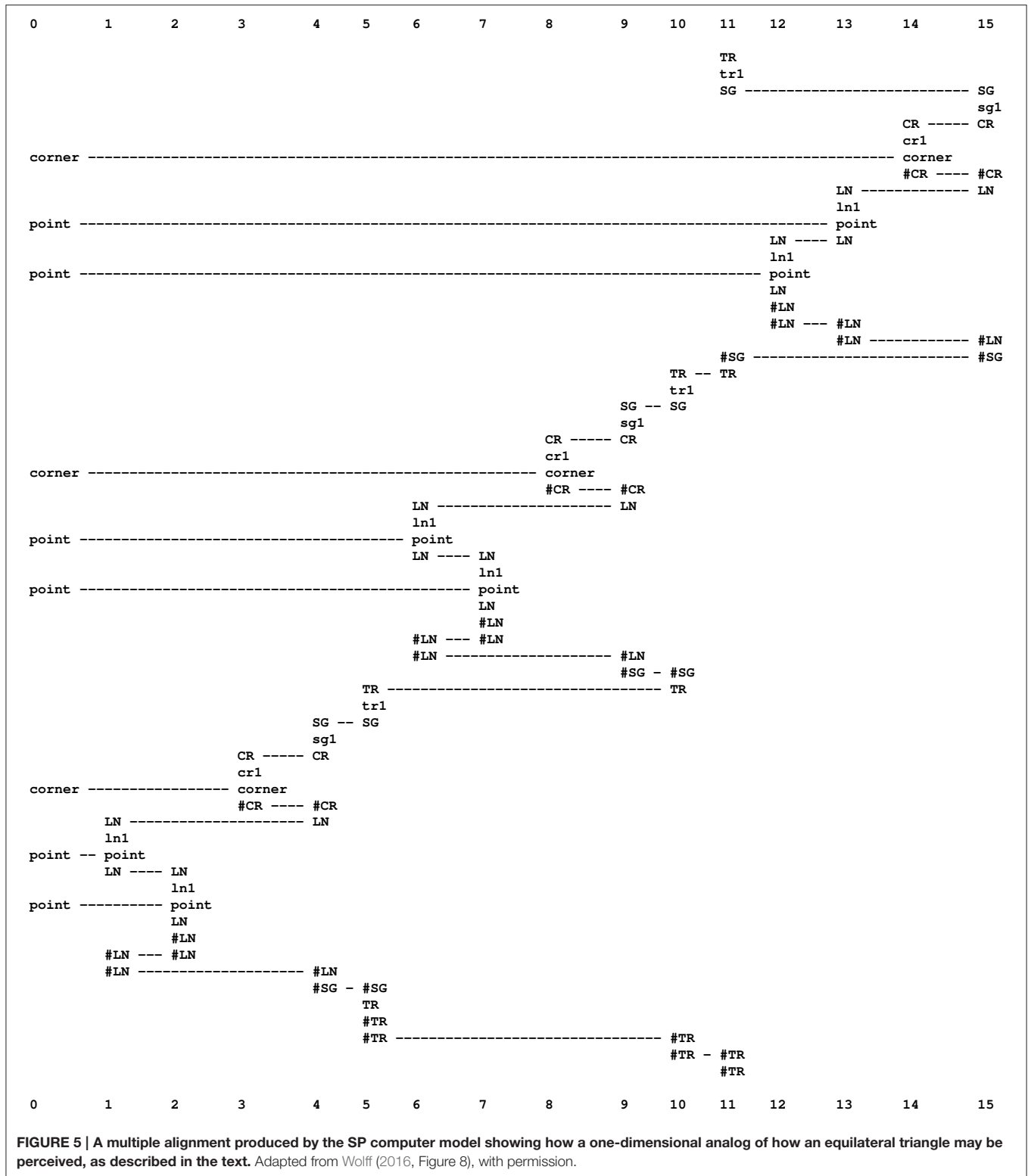
Column 0 shows information about the triangle to be recognized, comprising three "corners" and three sides of the triangle, each one represented by just two "points."

The pattern "LN ln1 point LN #LN #LN" in columns 1 and 2 is a self-referential and thus recursive definition of a line as a sequence of "points." It is self-referential because, within the body of the pattern, it contains a reference to itself via the symbols at the beginning and end of the pattern: "LN #LN." Because there is no limit to this recursion, it may represent a line containing any number of points. In a similar way, a second side is encoded via the same pattern in columns 6 and 7, and, again with the same pattern, the third line is encoded in columns 12 and 12.

In columns 4, 9 and 15 in the figure, the pattern "SG sg1 CR #CR LN #LN #SG" shows one of the three elements of a triangle as a corner ("CR #CR") followed by a line ("LN #LN"). And the recursion to encode multiple instances of that structure is in self-referential occurrences of the pattern "TR tr1 SG #SG TR #TR #TR" in columns 5, 10, and 22. Strictly speaking, the encoding is for a polygon, not a triangle, because there is nothing to stop the recursive repetition of "SG sg1 CR #CR LN #LN #SG." And, in terms of the problem, as described above, the representation is incomplete because there is nothing to show that the three sides of the triangle are the same.

These encodings account for the redundancy in the repetition of points along a line and also the redundancy in the repetition of three sides of a triangle. In a 2D version, they would also account for the redundancy in the white space within the body of the triangle, because they would allow most of the white space to be eliminated via shrinkage of the representation to the minimum needed to express the concept of a triangle.

¹³Compared with the multiple alignments shown in Figures 1, 2, this multiple alignment has been rotated by 90°. The choice between these alternative presentations of multiple alignments depends entirely on what fits best on the page.



5.6. We See Much Less than We Think We See

Most people with normal vision have a powerful sense that their eyes are a window on to a kind of cinema screen that shows what

we are looking at with great detail from left to right and from top to bottom. But research shows otherwise:

- In the phenomenon of *inattention blindness*, people may fail to notice salient things in their visual fields when they are

looking for something else, even if they are trained observers. In a recent demonstration (Drew et al., 2013), radiologists were asked to search for lung-nodules in chest x-rays but many of them (83%) failed to notice the image of a gorilla, 48 times the size of the average nodule, that was inserted into one of the radiographs.

- In the phenomenon of *change blindness*, people often fail to notice large changes to visual scenes. For example, if a conversation between two people—the investigator and the experimental subject—is interrupted by a door being carried between them, the experimental subject may fail to notice, when the door has gone by, that the person they are speaking to is different from the person they were speaking to before (Simons and Ambinder, 2005).
- Although each of our eyes has a blind spot¹⁴, we don't notice it, even when we are viewing things with one eye (so that there is no possibility that the blind spot in one eye will be filled in via vision in the other eye). Apparently, our brains interpolate what is likely to be in the blind part of our visual field.

It seems that part of the reason for this failure to see things is that photoreceptors are concentrated at the fovea (Squire et al., 2013, p. 502), and cones are only found in that region (Squire et al., 2013), so that, with two eyes, we are, to a large extent, looking at the world through a keyhole composed of two circumscribed and largely overlapping views, one from each eye.

It seems that our sense that the world is displayed to us on a wide and deep cinema screen is partly because our perception of any given scene draws heavily on our memories of similar scenes and partly because we can piece together what will normally be a partial view of what we are looking at from many short glimpses through the “keyhole” as we move our gaze around the scene.

The SP theory provides an interpretation for these things as follows:

- The theory provides an account in some detail of how New (sensory) information may be related to Old (stored) information and how an interpretation of the New information may be built up via the creation of multiple alignments. When sensory information provides an incomplete description of some entity or scene (which is normally the case), we fill in the gaps from stored knowledge.
- The theory provides an account of how we can piece together a picture of something, or indeed a 3D model of something, from many small but partially-overlapping views, in much the same way that: (1) With digital photography, it is possible to create a panoramic picture from several partially-overlapping images; (2) The views in Google's Streetview are built up from many partially-overlapping pictures; (3) A 3D digital image of an object may be created from partially-overlapping images of the object, taken from viewpoints around it. These things are discussed in Wolff (2014a, Sections 5.4, 6.1).

With regard to the second point, it should perhaps be said that partial overlap between “keyhole” views is not an essential part of building up a big picture from smaller views. But if two or more views do overlap, it is useful if they can be stitched together,

¹⁴See “Blind spot (vision),” *Wikipedia*, bit.ly/1oI0vyI, retrieved 2016-04-08.

thus removing the overlap. And partial overlap may be helpful in establishing the relative positions of two or more views.

5.7. A Resolution Problem and Its Possible Resolution

As we have seen (Section 5.1), each hypercolumn in the primary visual cortex of cats occupies about 1 mm^2 at the surface of the cortex, and it seems likely that each such hypercolumn provides a means of encoding one out of an alphabet of perceptual primitives, such as a line at a particular angle.

Assuming that this interpretation is correct, and if we view the primary visual cortex as if it was film in an old-style camera or the image sensor in a digital camera, it may seem that the encoding of perceptual primitives, with 1 mm^2 for each one, is remarkably crude. How could such a system—with the area of the primary visual cortex corresponding to the area of our field of view—create that powerful sense that, through our eyes, we see a detailed “cinema screen” view of the world (Section 5.6).

Part of the answer is probably that we see much less than we think we see (Section 5.6). But it seems likely that another part of the answer is to reject the assumption that the whole of the primary visual cortex corresponds to the area of our field of view. In the light of the remarks in Section 5.6, it seems more likely that, normally, in each of the previously-mentioned glimpses of a scene, all of the primary visual cortex or most of it is applied in the assimilation and processing of information capture by the fovea and, perhaps, parts of the retina that are near to the fovea.

In support of this idea: “*Cortical magnification* describes how many neurons in an area of the visual cortex are ‘responsible’ for processing a stimulus of a given size, as a function of visual field location. In the center of the visual field, corresponding to the fovea of the retina, a very large number of neurons process information from a small region of the visual field. If the same stimulus is seen in the periphery of the visual field (i.e., away from the center), it would be processed by a much smaller number of neurons. The reduction of the number of neurons per visual field area from foveal to peripheral representations is achieved in several steps along the visual pathway, starting already in the retina (Barghout-Stein, 1999)”¹⁵.

With this view of visual processing, what appears superficially to be a rather course-grained recording and analysis of visual data, may actually be very much more detailed. As described in Section 5.6, it seems likely that our view of any scene is built up partly from memories and partly from many small snapshots or glimpses of the scene. And it seems like that each such snapshot or glimpse is processed using a relatively large neural resource.

5.8. Grandmother Cells, Localist and Distributed Representations

In terms of concepts that have been debated about how knowledge may be represented in the brain, the ID-neural-symbols for any pattern assembly are very much like the concept of a *grandmother cell*—a cell or small cluster of cells in one's brain that represents one's grandmother so that, if the cell or

¹⁵See “Cortical magnification,” *Wikipedia*, bit.ly/1qJsQX1, emphasis in the original, retrieved 2016-04-14.

cells were to be lost, one would lose the ability to recognize one's grandmother¹⁶.

It seems that the weight of observational and experimental evidence favors the belief that such cells do exist (Gross, 2002; Roy, 2013). This is consistent with the observation that people who have suffered a stroke or are suffering from dementia may lose the ability to recognize members of their close family.

Since SP-neural, like Hebb's (1949) theory of cell assemblies, proposes that concepts are represented by coherent groups of neurons in the brain, it is very much a "localist" type of theory. As such, it is quite distinct from "distributed" types of theory that propose that concepts are encoded in widely-distributed configurations of neurons, without any identifiable location or center.

However, just to confuse matters, SP-neural does *not* propose that all one's knowledge about one's grandmother would reside in a pattern assembly for that lady. Probably, any such pattern assembly would, in the manner of object-oriented design as discussed in Section 6 and illustrated in **Figure 6**, be connected to and inherit features from a pattern assembly representing grandmothers in general, and from more general pattern assemblies such as pattern assemblies for such concepts as "person" and "woman." And again, a pattern assembly for "person" would not be the sole repository of all one's knowledge about people. That pattern assembly would, in effect, contain "references" to pattern assemblies describing the parts of a person, their physiology, their social and political life, and so on.

Thus, while SP-neural is unambiguously localist, it proposes that knowledge of any entity or concept is likely to be encoded not merely in one pattern assembly for that entity or concept but also in many other pattern assemblies in many parts of the cortex, and perhaps elsewhere.

5.9. Positional Invariance

With something simple like a touch on the skin, or a pin prick, it is not too difficult to see how the sensation may be transmitted to the brain via any one of many relevant receptors located in many different areas of the skin. But with something more complex, like an image on the retina of a table, a house, or a tree, and so on, it is less straightforward to understand how we might recognize such a thing in any part of our visual field.

For each entity to be recognized, it seems necessary at first sight to provide connections, directly or indirectly, from every part of the receptor array to the relevant pattern assembly. In terms of the schematic representation shown in **Figure 3**, it would mean repeating the connections for "t h e" and "b r a v e" in each of many parts of the receptor array. Bearing in mind the very large number of different things we may recognize, the number of necessary connections would become very large, perhaps prohibitively so.

However, things may be considerably simplified via either or both of two provisions:

1. For reasons outlined in Section 5.6, it seems likely that, with vision, we build up our perception of a scene, partly from memories of similar scenes and partly via many relatively

narrow "keyhole" views of what is in front of us. If that is correct, and if, as suggested in Section 5.7, most of the primary visual cortex is devoted to analysing information received via the fovea and, perhaps, via parts of the retina that are very close to the fovea, then the need to provide for any given pattern in many parts of the receptor array may be greatly reduced. Since, by moving our eyes, we may view any part of a scene, it is possible that any given entity would need only one or two sets of connections between the receptor array and the pattern assembly for that entity.

2. As noted in Section 4.3, it seems likely that, with regard to **Figure 3**, there would, in a more realistic example, be several levels of structure between neural symbols in the receptor array and relatively complex structures like words. At the first level above the receptor array there would be pattern assemblies for relatively small recurrent structures, and the variety of such structures would be relatively small. This should ease any possible problems in connecting the receptor array to pattern assemblies.

If it turns out that the number of necessary connections is indeed too large to be practical, or if there is empirical evidence against such numbers, then a possible alternative to what has been sketched in this paper is some kind of dynamic system for the making and breaking of connections between the receptor array and pattern assemblies. It seems likely that permanent or semi-permanent connections would be very much more efficient and the balance of probabilities seems to favor such a scheme.

In connection with positional invariance, it is relevant to note that "... lack of localization is quite common in higher-level neurons: receptive fields become larger as the features they represent become increasingly complex. Thus, for instance, neurons that respond to faces typically have receptive fields that cover most of the visual space. For these cells, large receptive fields have a distinct advantage: the preferred stimulus can be identified no matter where it is located on the retina." (Squire et al., 2013, p. 579). A tentative and partial explanation of this observation is that repetition of neurons that are sensitive to each of several categories of low-level feature—in the receptor array and as ID-neural-symbols for "low-level" pattern assemblies—is what allows positional invariance to develop at higher levels.

6. NON-SYNTACTIC KNOWLEDGE IN SP-NEURAL

As was emphasized in Section 3, the SP system (SP-abstract) has strengths and potential in the representation and processing of several different kinds of knowledge, not just the syntax of natural language. That versatility has been achieved using the mechanisms in SP-abstract that were outlined in that section. If those mechanisms can be modeled in SP-neural, it seems likely that the several kinds of knowledge that may be represented and processed in SP-abstract may also be represented and processed in SP-neural.

As an illustration, **Figure 6** shows a simple example of how, via multiple alignment, the SP computer model may recognize an unknown creature at several different levels of abstraction,

¹⁶See "Grandmother cell," *Wikipedia*, bit.ly/1UDulyV, retrieved 2016-08-26.

and **Figure 7** suggests how part of the multiple alignment, with associated patterns, may be realized in terms of pattern assemblies and their inter-connections.

Figure 6 shows the best multiple alignment found by the SP computer model with four symbols representing attributes of an unknown creature (shown in column 0) and a collection of Old patterns representing different creatures and classes of creature, some of which are shown in columns 1–4, one pattern per column. In a more detailed and realistic example, symbols like “eats,” “retractile-claws,” and “breathes,” would be represented as patterns, each with its own structure.

From this multiple alignment, we can see that the unknown creature has been identified as an animal (column 4), as a mammal (column 3), as a cat (column 2) and as a specific cat, “Tibs” (column 1). It is just an accident of how the SP computer model has worked in this case that the order of the patterns across columns 1–4 of the multiple alignment corresponds with the level of abstraction of the classifications. In general, the order of patterns in columns above 0 is entirely arbitrary, with no significance.

Figure 7 shows how part of the multiple alignment from **Figure 6** may be realized in SP-neural. The figure contains pattern assemblies for “animal” and “mammal,” corresponding to

patterns from columns 4 and 3 of the multiple alignment. Notice that the left-right order of the pattern assemblies is different from the order of the patterns in the multiple alignment, in accordance with the remarks, above, about the workings of the SP computer model, and also because there is no reason to believe that pattern assemblies are represented in any particular order.

Neural connections amongst the things that have been mentioned so far are very much the same as alignments between neural symbols in **Figure 6**: “eats” on the left connects with “eats” in the “animal” pattern assembly; “furry” connects with “furry” in the “mammal” pattern assembly, and the “A” and “#A” connections for those two pattern assemblies correspond with the alignments of symbols in the multiple alignment. As in **Figure 3**, some neural connections are shown with broken lines to suggest that they would be relatively inactive during the neural processing which identifies one or more “good” NAMAs. And as before, it is envisaged that there would be one or more neural connections between each neural symbol and its immediate neighbors within each pattern assembly, but these are not marked in the figure.

The inclusion of a pattern assembly for “reptile” in **Figure 7**, with some of its neural connections, is intended to suggest some of the processing involved in identifying one or more winning NAMAs. In the same way that the pattern for “mammal” is

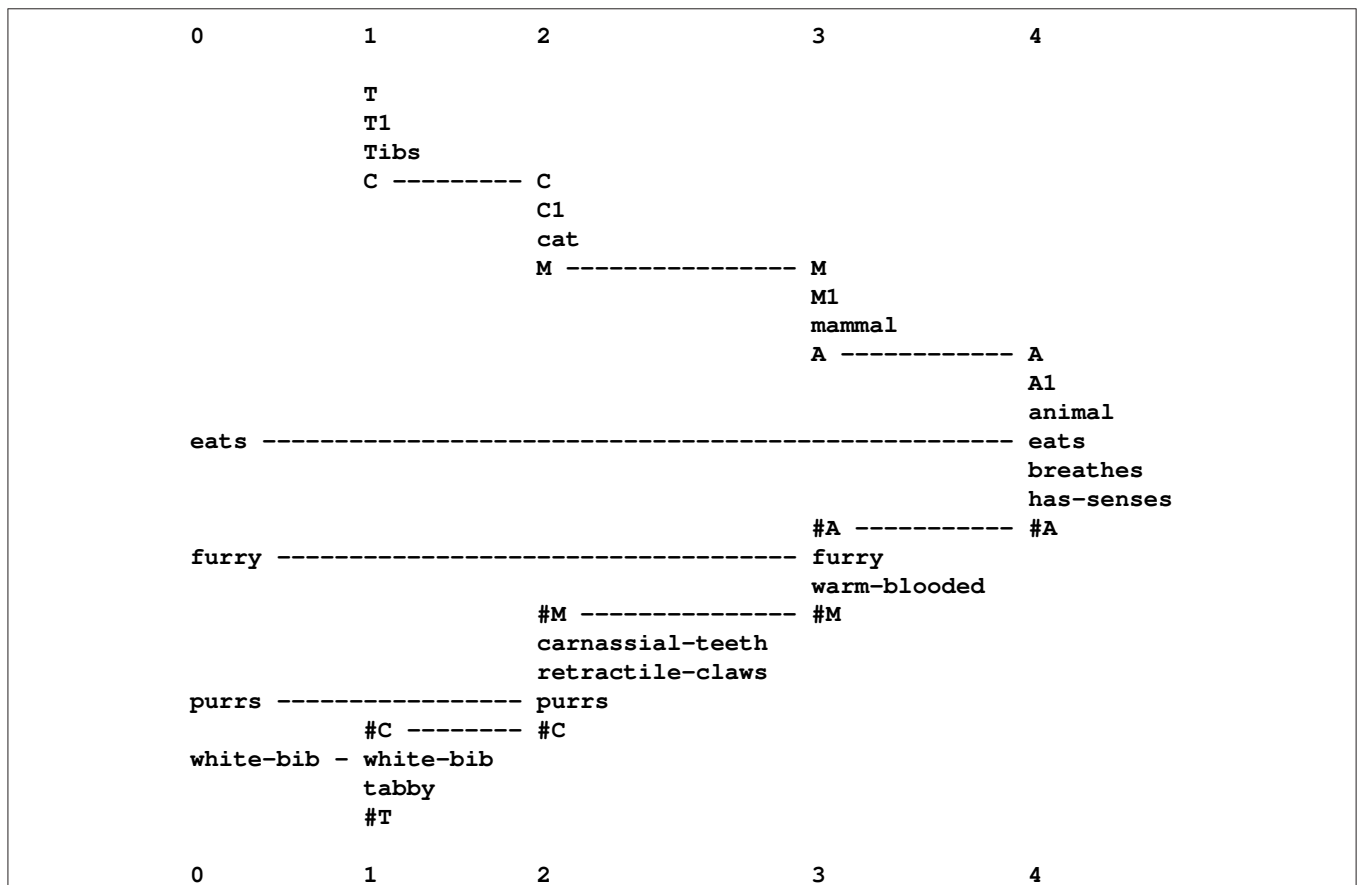


FIGURE 6 | The best multiple alignment found by the SP computer model with four one-symbol New patterns representing attributes of an unknown creature and a collection of Old patterns representing different creatures and classes of creature. Adapted from Figure 6.7 in Wolff (2006), with permission.

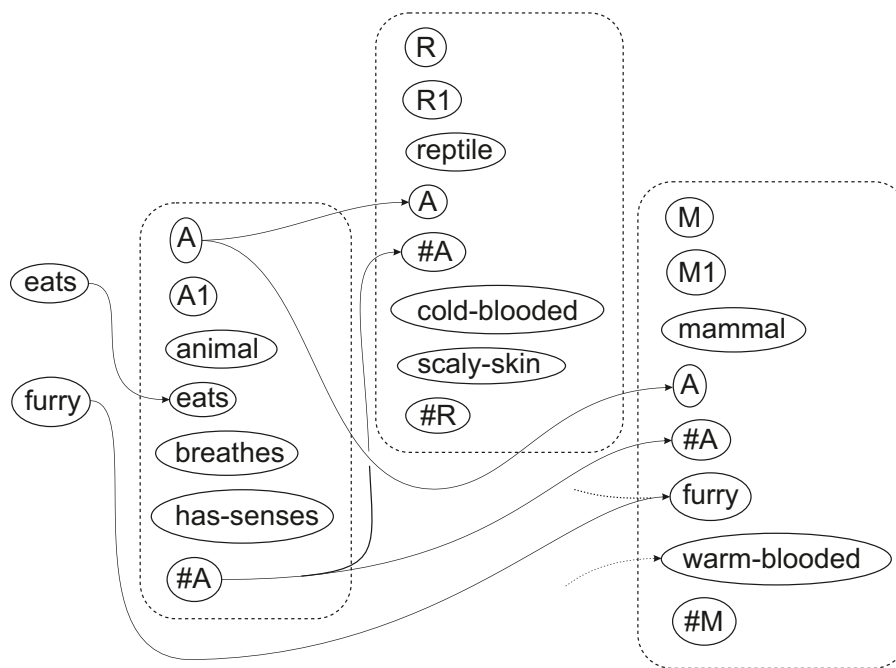


FIGURE 7 | How part of the multiple alignment shown in Figure 6 may be realized in SP-neural—showing two of the attributes from column 0 in the multiple alignment and with “animal” and “mammal” pattern assemblies corresponding to patterns from columns 4 and 3—with an associated pattern assembly for “reptile.” The conventions are the same as in Figure 3.

receiving excitatory signals from the pattern for “animal,” one would expect excitatory signals to flow to pattern assemblies for the other main groups of animals, including reptiles. Ultimately, “reptile” would fail to feature in any winning NAMA because of evidence from the neural symbols “furry,” “purrs,” and “white-bib.”

7. REPETITION AND RECURSION

Like any good database or dictionary, the repository of Old patterns in SP-abstract should only contain one copy of any given SP pattern. But in something like *Jack Sprat could eat no fat, His wife could eat no lean*, the words *could*, *eat*, and *no* each occur twice. With an example like this, it seems reasonable to suppose that there is only one stored pattern for each of the repeated words, and likewise for the many other examples of entities that are repeated within something larger, witness the many legs of a centipede.

In SP-abstract, this apparent difficulty has been overcome by saying that each SP pattern in a multiple alignment is an *appearance* of the pattern, not the pattern itself—which allows us to have multiple instances of a pattern in a multiple alignment without breaking the rule that the repository of Old patterns should contain only one copy of each pattern. But in SP-neural, it is not obvious how to create an “appearance” of a pattern assembly that is not also a physical structure of neurons and their interconnections—but the speed with which we can understand natural language seems to rule out what appears

to be the relatively slow growth of new neurons and their interconnections.

How we can create new mental structures quickly arises again in other connections, as discussed in Section 11. If we duck these questions for the time being and return to parsing, it may be argued that with something like *Jack Sprat could eat no fat, His wife could eat no lean*, the first instance of *could* is represented only for the duration of the word by the stored pattern for *could*, so that the same pattern can be used again to represent the second instance of *could*—and likewise for *eat* and *no*. But it appears that this line of reasoning does not work with a recursive structure like *the very very very fast car*.

Native speakers of English know that with a phrase like *the very very very fast car*, the word *very* may in principle be repeated any number of times. This observation, coupled with the observation that recursive structures are widespread in English and other natural languages, suggests strongly that the most appropriate parsing of the phrase is something like the multiple alignment shown in Figure 8. Here, the repetition of *very* is represented via three appearances of the pattern “ri ril ri #ri i #i #ri,” a pattern which is self-referential because the inner pair of symbols “ri #ri” can be matched with the same two symbols, one at the beginning of the pattern and one at the end. Because the recursion depends on at least two instances of “ri ril ri #ri i #i #ri” being “live” at the same time, it seems necessary for SP-neural to be able to model multiple appearances of any pattern.

That conclusion, coupled with the above-mentioned arguments from the speed at which we can speak, and the speed

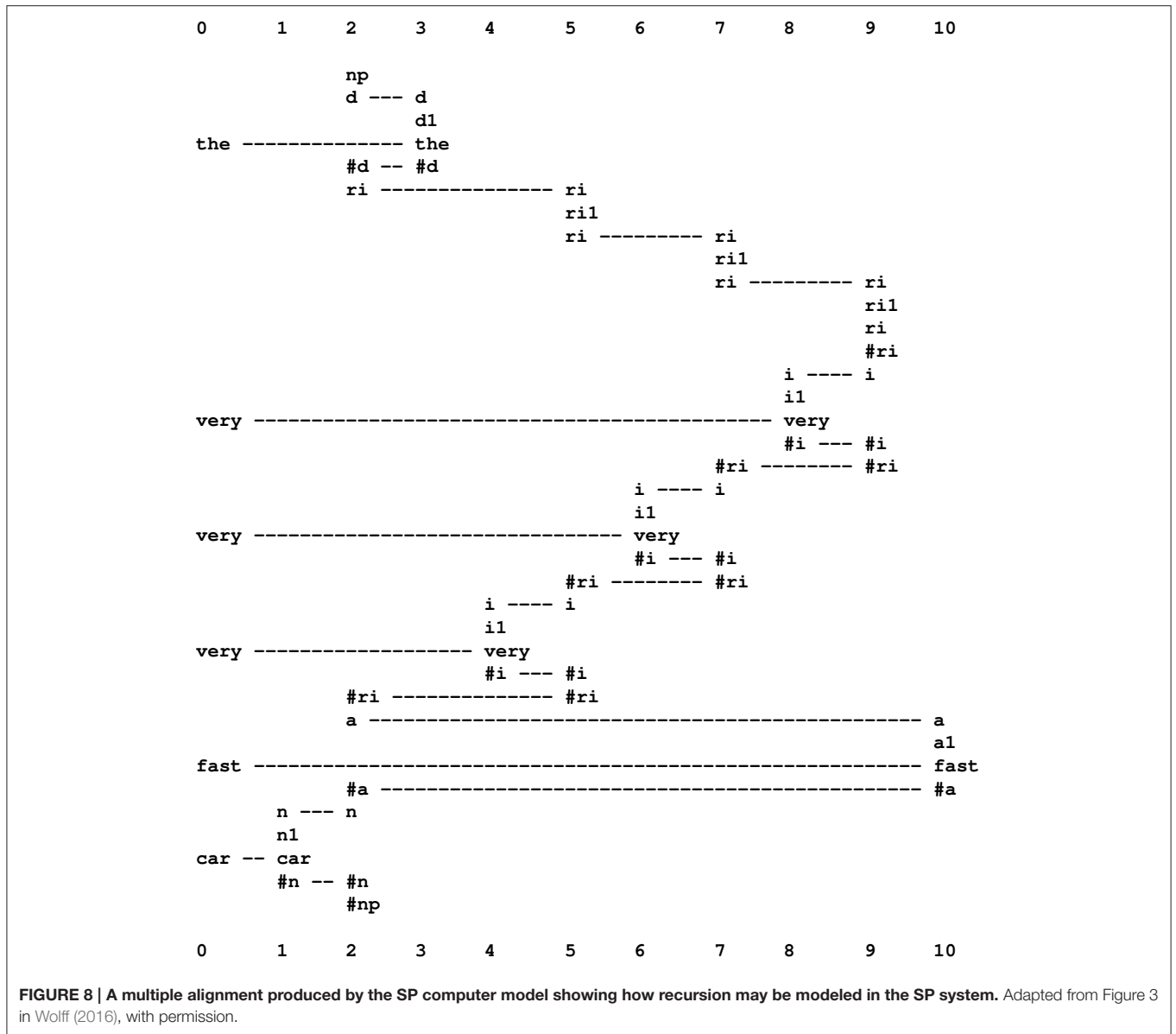


FIGURE 8 | A multiple alignment produced by the SP computer model showing how recursion may be modeled in the SP system. Adapted from Figure 3 in Wolff (2016), with permission.

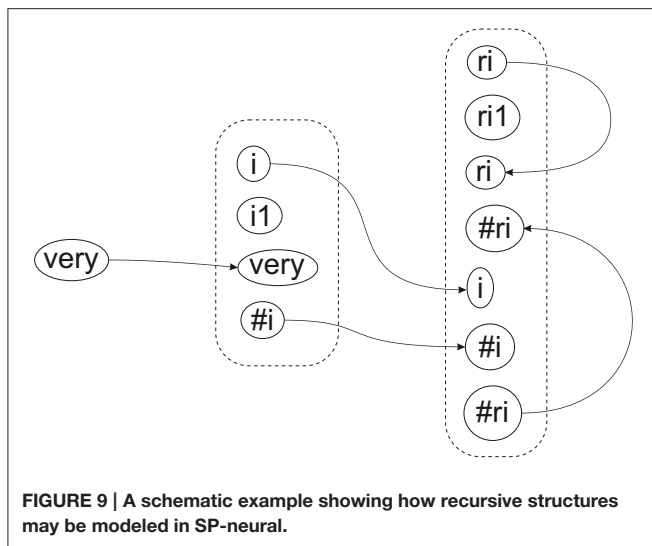
with which we can imagine new things, argues strongly that SP-neural—and any other neural theory of cognition—must have a means of creating new mental structures quickly. It seems unlikely that these things could be done via the growth of new neurons and their interconnections.

The tentative answer suggested here is that, in processes like parsing or pattern recognition, including examples with recursion like that shown in **Figure 8**, virtual copies of pattern assemblies may be created and destroyed very quickly via the switching on and switching off of synapses (Section 11). Clearly, more detail is needed for a fully satisfactory answer.

Pending that better answer, **Figure 9** shows tentatively how recursion may be modeled in SP-neural, with neural symbols and pattern assemblies corresponding to selected symbols and

patterns in **Figure 8**. On the left of that figure, we can see how the neural symbol “very” connects with a matching neural symbol in the pattern assembly “i il very #i.” Further right, we can see how the first and last neural symbols in “i il very #i” connect with matching neural symbols in the pattern assembly “ri ril ri #ri i #i #ri.”

In the figure, the self-referential nature of the pattern assembly “ri ril ri #ri i #i #ri” can be seen in the neural connection between “ri” at the beginning of that pattern assembly and the matching neural symbol in the body of the same pattern assembly, and likewise for “#ri” at the end of the pattern assembly. Although it is unclear how this recursion may achieve the effect of repeated appearances of the pattern assembly at the speed with which we understand or produce speech, the analysis appears to be more reliable than what is described



in Wolff (2006, Section 11.4.2), especially Figure 11.10 in that section.

8. SP-NEURAL: AN OUTPUT PERSPECTIVE

An inspection of **Figure 3**—showing how, in SP-neural, a small portion of natural language may be analyzed by pattern assemblies and their interconnections—may suggest that if we wish to reverse the process—to create language instead of analysing it—then the innervation would need to be reversed: we may guess that two-way neural connections would be needed to support the production of speech or writing as well as their interpretation.

But a neat feature of SP-abstract is that one set of Old patterns, together with the processes for building multiple alignments, will support both the analysis and the production of language. So it is reasonable to suppose that if SP-neural works at all, a similar duality will apply to pattern assemblies and their interconnections, without the need for two-way connections amongst pattern assemblies and neural symbols (but see Section 8.3).

Of course, speaking or writing would need peripheral motor processes that are different from the peripheral sensory processes required for listening or reading, but, more centrally, the processes for analysing language or producing it may use the same mechanisms¹⁷.

The reason that SP-abstract, as expressed in the SP computer model, can work in “reverse” so to speak, is that, from a multiple alignment like the one shown in **Figure 2**, a code pattern like “S 0 2 4 3 7 6 1 8 5 #S” may be derived, as outlined in Section 3.6. Then, if that code pattern is presented to the SP system as a New pattern, the system can recreate the original sentence, “f o r t u n e f a v o u r s t h e b r a v e,” as shown in **Figure 10**.

¹⁷Of course, things are a little more complicated with output processes because sensory feedback is normally an important part of speaking or writing.

8.1. An Answer to the Apparent Paradox of “Decompression by Compression”

That the SP system should be able to reconstruct a sentence that was originally compressed by means of the same system (Section 8) may seem paradoxical. How is it that a system that is dedicated to information compression should be able, so to speak, to drive compression in reverse?

A resolution of this apparent paradox is described in Wolff (2006, Section 3.8). In brief, the key to the conjuring trick is to ensure that, after the sentence has been compressed, there is enough residual redundancy in the code pattern to allow further compression, and to ensure that this further compression will achieve the effect of reconstructing the sentence.

8.2. Meanings in the Analysis and Production of Language

Of course, parsing a sentence (as shown in Section 3.5) or constructing a sentence from a code pattern (as shown in Section 8) are very artificial applications with natural language. Normally, when we read some text or listen to someone speaking, we aim to derive meaning from the writing or the speech. And when we write or speak, it seems, intuitively, that the patterns of words that we are creating are derived from some kind of underlying meaning that we are trying to express.

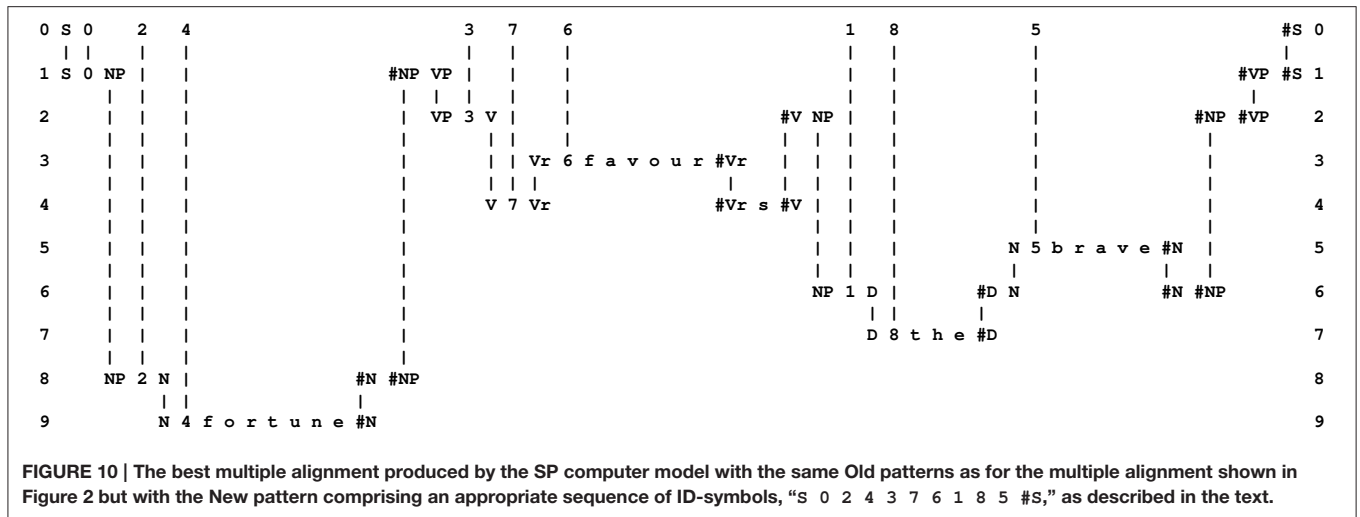
It is envisaged that, in future development of SP-abstract and the SP computer model, the ID-symbols in code patterns will provide some kind of bridge between syntactic forms and representations of meanings, thus facilitating the processes of understanding the meanings of written or spoken sentences and of creating sentences to express particular meanings.

As noted at the end of Section 3.8.2, there are preliminary examples of how, with the SP computer model, a sentence may be analyzed for its meaning (Wolff, 2006, Section 5.7, Figure 5.18), and how the same sentence may be derived from a representation of its meaning (Wolff, 2006, Figure 5.19).

8.3. But There Are Projections from the Sensory Cortex to Subcortical Nuclei

Although as we have seen earlier in Section 8, SP-neural, via principles established in SP-abstract, provides for the creation of language, and other kinds of knowledge, without the need for efferent connections from the cortex back along the path of afferent nerves, there is evidence that such connections do exist:

“Neurons of the cerebral cortex send axons to subcortical regions Subcortical projections are to those nuclei in the thalamus and brainstem that provide ascending sensory information. By far the most prominent of these is to the thalamus: the neurons of a primary sensory cortex project back to the same thalamic nucleus that provides input to the cortex. This system of descending connections is truly impressive because the number of descending corticothalamic axons greatly exceeds the number of ascending thalamocortical axons. These connections permit a particular sensory cortex to control the activity of the very neurons that relay information to it.” (Squire et al., 2013, p. 509).



But the descending nerves described in this quotation may have a function that is quite different from the creation of sentences or other patterns of activity. One possible role for such nerves may be “the focussing of activity so that relay neurons most activated by a sensory stimulus are more strongly driven and those in surrounding less well activated regions are further suppressed.” (Squire et al., 2013, p. 509).

9. THE POSSIBLE ROLES FOR INHIBITION IN SP-NEURAL

A familiar observation is that, if something like a fan is switched on near us, we notice the noise for a while and then come to ignore it. And if, later, the fan is switched off, we notice the relative quiet for a while and then cease to be aware of it. In general, it seems that we are relatively sensitive to changes in our environment and relatively insensitive to things that remain constant.

It has been accepted for some time that the way we adapt to constant stimuli is due to inhibitory neural structures and processes in our brains and nervous systems, that inhibitory structures and processes are widespread in the animal kingdom, and that they have a role in reducing the amount of information that we need to process (von Békésy, 1967).

Regarding the last point, it is clearly inefficient for anyone to be constantly registering, second-by-second, the noise of a nearby fan: “noise, noise, noise, noise, noise, . . .” and likewise for the state of relative quietness when the fan is switched off. In terms of information theory, there is *redundancy* in the second-by-second recurrence of the noise (or quietness), and we can eliminate most or all of the redundancy—and thus compress the information—by simply recording that the noise is “on” and that it is continuing (and likewise, *mutatis mutandis*, for quiet). This is the “run-length encoding” technique for compression of information,¹⁸ it is essentially what

adaptation does, and, in neural tissue, it appears to be mediated largely by “lateral” inhibition.

With lateral inhibition in sensory neurons, there are inhibitory connections between neighboring neurons so that, when they are both stimulated, they tend to inhibit each other, and thus reduce their rates of firing where there is strong uniform stimulation. But inhibition is reduced where strong stimulation gives way to weaker stimulation, leading to a local swing in the rate of firing (Ratliff et al., 1963; see also Wolff, 2006, Section 2.3.1; there is more about lateral inhibition in Squire et al., 2013, p. 505). There are similar effects in the time dimension. Again, Barlow (1982) says, in connection with neurons in the mammalian cortex that receive inputs from both eyes, “... it is now clear that input from one eye can, and frequently does, inhibit the effects of input from the other eye, ...” (p. 147).

Taking these observations together, we may abstract a general rule: *When, in neural processing, two or more signals are the same, they tend to inhibit each other, and when they are different, they don't.* The overall effect should be to detect redundancy in information and to reduce it, whilst retaining non-redundant information, in accordance with the central principle in the SP theory—that much of computing and cognition may, to a large extent, be understood as information compression.

In a similar vein: “Lateral inhibition represents the classic example of a general principle: most neurons in sensory systems are best adapted for detecting changes in the external environment. ... As a rule, it is change which has the most significance for an animal ... This principle can also be explained in terms of information processing. Given a world that is filled with constants—with uniform objects, with objects that move only rarely—it is most efficient to respond only to changes.” (Squire et al., 2013, p. 578).

In view of the widespread occurrence of inhibitory mechanisms in the brain¹⁹, and in view of their apparent

¹⁹“These [aspiny or sparsely spiny nonpyramida] interneurons constitute approximately 15–30% of the total population of cortical neurons, and they appear to be mostly GABAergic, representing the main components of inhibitory cortical circuits” (Squire et al., 2013, p. 45); “Synaptic inhibition in the mammalian brain

¹⁸See “Run-length encoding,” *Wikipedia*, bit.ly/21JlB1T, retrieved 2016-03-04.

a manner that is more like a tailor cutting up pre-woven cloth than someone knitting or crocheting each item from scratch.

In accordance with the labeled line principle (Section 5.3), the meaning of each symbol in a newly-created pattern assembly would be determined by what it is connected to, as described in Section 10.2.

Similar principles would apply when Old patterns are created from partial matches between patterns, as described in Section 3.4.

10.2. Creating Connections between Pattern Assemblies

As with the laying down of newly-created Old patterns (Section 10.1), it seems unlikely that connections between pattern assemblies, like those shown in **Figure 3**, would be created by growing new axons or dendrites. It seems much more likely that such connections would be established by switching on synapses between each of the two neurons to be connected and pre-existing axons or dendrites, somewhat like the making of connections in a telephone exchange (see Section 11).

This idea, together with the suggestions in Section 10.1 about how Old pattern assemblies may be created, is somewhat like the way in which an “uncommitted logic array” (ULA)²¹ may, via small modifications, be made to function like any one of a wide variety of “application-specific integrated circuits” (ASICs)²², or how a “field-programmable gate array” (FPGA)²³ may be programmed to function like any one of a wide variety of integrated circuits.

10.3. Destruction of Pattern Assemblies and Their Interconnections

In the SP theory, patterns and pattern assemblies are never modified—they are either created or destroyed. The latter process occurs mainly in the process of searching for “good” grammars to describe a given set of New patterns, as outlined in Section 3.7. At each stage, when a few “good” grammars are retained in the system, the rest are discarded. This means that any pattern assembly in one or more of the “bad” grammars that is not also in one or more of the “good” grammars may be destroyed.

It seems likely that, in a process that may be seen as a reversal of the way in which pattern assemblies and their interconnections are created, the destruction of a pattern assembly does not mean the physical destruction of its neurons. It seems more likely that all neural connections from the pattern assembly are broken by switching off relevant synapses (Sections 10.3, 11) and that its constituent neurons are retained for later use in other pattern assemblies.

10.4. Searching for Good Grammars

It must be admitted that, apart from the remarks in forgoing subsections about the creation and destruction of pattern

²¹See “Gate array,” *Wikipedia*, bit.ly/1UdB46j, retrieved 2016-03-20.

²²See “Application-specific integrated circuit,” *Wikipedia*, bit.ly/1pUs2y8, retrieved 2016-03-20.

²³See “Field-programmable gate array,” *Wikipedia*, bit.ly/1Hgi9iH, retrieved 2016-03-20.

assemblies and their inter-connections, it is unclear how, in SP-neural, one may achieve anything equivalent to the process of searching the abstract space of possible grammars that has been implemented in the SP computer model.

One possibility is to simplify things as follows. Instead of evaluating whole grammars, as in the SP computer model, it may be possible to achieve roughly the same effect by evaluating pattern assemblies in terms of their effectiveness or otherwise for the economical encoding of New information and, periodically, to discard those pattern assemblies that do badly.

10.5. What about Hebbian Learning?

Readers familiar with issues in AI or neuroscience may wonder what place, if any, there may be in SP-neural for the concept of “Hebbian” learning. This idea, proposed by Hebb (1949), is that:

“When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A’s efficiency, as one of the cells firing B, is increased.” (p. 62).

Variants of this idea are widely used in versions of “deep learning” in artificial neural networks (Schmidhuber, 2015) and have contributed to success with such systems²⁴.

But in Wolff (2016, Section V-D) I have argued that:

- The gradual strengthening of neural connections which is a central feature of Hebbian learning (and deep learning) does not account for the way that people can, very effectively, learn from a single occurrence or experience (sometimes called “one-trial” learning)²⁵.
- Hebb was aware that his theory of learning with cell assemblies would not account for one-trial learning and he proposed a “reverberatory” theory for that kind of learning (Hebb, 1949, p. 62). But, as noted in Wolff (2016, Section V-D), Milner has pointed out (Milner, 1996) that it is difficult to understand how this kind of mechanism could explain our ability to assimilate a previously-unseen telephone number: for each digit in the number, its pre-established cell assembly may reverberate; but this does not explain memory for the *sequence* of digits in the number. And it is unclear how the proposed mechanism would encode a phone number in which one or more of the digits is repeated.
- One-trial learning is consistent with the SP theory because the direct intake and storage of sensory information is bedrock in how the system learns (Section 3.4).
- The SP theory can also account for the relatively slow learning of complex skills such as how to talk or how to play tennis at a

²⁴See, for example, “Don’t despair if Google’s AI beats the world’s best Go player,” *MIT Technology Review*, bit.ly/1p7Wzb7, 2016-03-08; and “Google unveils neural network with “superhuman” ability to determine the location of almost any image,” *MIT Technology Review*, bit.ly/1p5qmSe, 2016-02-24.

²⁵It may be argued that Hebbian learning may apply in such cases because a single experience may be mentally rehearsed. But that begs the question of how the one experience is remembered between when it first occurred and the first rehearsal—and likewise later on. And, while rehearsal may be helpful in some cases, it seems that there are many things we do remember after a single experience, without rehearsal.

high standard—because of the complexity of the abstract space of possible solutions that needs to be searched.

Does this mean that Hebbian learning is dead? Probably not:

- In some forms, the phenomena of “long-term potentiation” (LTP) in neural functioning seem to be linked to Hebbian types of learning (Squire et al., 2013, pp. 1022–1023).
- Gradual strengthening of neural connections may have a role to play in SP-neural because some such mechanism is needed to record, at least approximately, the frequency of occurrence of neural symbols and pattern assemblies (Sections 3.4, 4.5).

11. THE PROBLEMS OF SPEED AND EXPRESSIVENESS IN THE CREATION AND DESTRUCTION OF NEURAL STRUCTURES

A general issue for any neural theory of the representation and processing of knowledge, is how to account for the speed with which we can create neural structures, and, probably, destroy them, bearing in mind that such structures must be sufficiently versatile to accommodate the representation and processing of a wide range of different kinds of knowledge. This issue arises mainly in the following connections:

- *One-trial learning.* In keeping with the remarks above about one-trial learning (Section 10.5), it is a familiar feature of everyday life that we can see and hear something happening—a football match, a play, a conversation, and so on—and then, immediately or some time later, give a description of the event. This implies that we can lay down relevant memories at speed.
- *The learning of complex knowledge and skills.* If we accept the view of unsupervised learning which is outlined in Sections 3.4, 3.7, and 10, then it seems necessary to suppose that pattern assemblies are created and destroyed during the search for one or two grammars that provide a “good” description of the knowledge or skills that is being learned—and it seems likely that the creation and destruction of pattern assemblies would be fast.
- *The interpretation of sensory data.* In processes like the parsing of natural language or, more generally, understanding natural language, and in processes like pattern recognition, reasoning, and more, it seems necessary to create intermediate structures like those shown in **Figure 2**, and for those structures to be created at speed.
- *Speech and action.* In a similar way, it seems necessary for us to create mental structures fast in any kind of activity that requires thought, such as speaking in a way that is meaningful and comprehensible, most kinds of sport, most kinds of games, and so on.
- *Imagination.* Most people have little difficulty in imagining things they are unlikely ever to have seen—such as a cat with a coat made of grass instead of fur, or a cow with two tails. We can create such ideas fast and, if we like them well enough, we may remember them for years.

One possible solution, which is radically different from SP-neural, is to suppose that our knowledge is stored in some chemical form

such as DNA, and that the kinds of mental processes mentioned above might be mediated via the creation and modification of such chemicals. Another possibility is that learning is mediated by epigenetic mechanisms, as outlined in Baars and Gage (2010, Section 7.4). Without wishing to prejudge what the primary mechanism of learning may be, or whether perhaps there are several such mechanisms, this paper focusses on SP-neural and how it may combine speed with expressiveness, as seems to be required for the kinds of functions outlined above.

At first sight, the problem of speed in the creation of neural structures is solved via the long-established idea that we can remember things for a few seconds via a “short-term memory²⁶” that is distinct from “long-term memory²⁷,” and “working memory²⁸.” But there is some uncertainty about the extent to which these three kinds of memory may be distinguished, one from another, and there is considerable uncertainty about how they might work, and how information may be transferred from one kind of memory to another.

As a proffered contribution to discussions in this area, the suggestion here is that, in any or all of short-term memory, working memory, and long-term memory, SP-neural may achieve the necessary speed in the creation of new structures, combined with versatility in the representation and processing of diverse kinds of knowledge, by the switching on and off of synapses in pre-established neural structures and their inter-connections, as outlined in Sections 10.1, 10.2.

With regard to possible mechanisms for the switching on and off of synapses:

- It appears that, in the entorhinal cortex between the hippocampus and the neocortex, there are neurons that can be switched “on” and “off” in an all-or-nothing manner (Tahvildari et al., 2007), and we may suppose that synapses have a role to play in this behavior.
- “The efficacy of a synapse can be potentiated through at least six mechanisms” (Squire et al., 2013, Caption to Figure 47.10) and it is possible that at least one them has the necessary speed, especially since “[Long-term potentiation] is defined as a persistent increase in synaptic strength ... that can be induced rapidly by a brief burst of spike activity in the presynaptic afferents.” (emphasis added) (Squire et al., 2013, p. 1016).
- “[Long-term depression] is believed by many to be ... a process whereby [Long-term potentiation] could be reversed in the hippocampus and neocortex” (Squire et al., 2013, p. 1023).
- “... it is now evident that [Long-term potentiation], at least in the dentate gyrus, can either be ... stable, lasting months or longer.” (Abraham, 2003, Abstract), although there appears to be little or no evidence with a bearing on whether or not there might be an upper limit to the duration of long-term potentiation.
- There is evidence that the protein kinase M ζ (PKM ζ) may provide a means of turning synapses on and off, and thus perhaps storing long-term memories (Ogasawara and Kawato, 2010).

²⁶“Short-term memory,” *Wikipedia*, bit.ly/1RzAVHN, retrieved 2016-04-04.

²⁷“Long-term memory,” *Wikipedia*, bit.ly/1M9uPhh, retrieved 2016-04-04.

²⁸“Working memory,” *Wikipedia*, bit.ly/1PQq0UA, retrieved 2016-04-04.

With all these possible mechanisms, key questions are: do they act fast enough to account for the speed of the phenomena described above; and can they provide the basis for memories that can last for 50 years or more.

12. ERRORS OF OMISSION, COMMISSION, AND SUBSTITUTION

A prominent feature of human perception is that we have a robust ability to recognize things despite disturbances of various kinds. We can, for example, recognize a car when it is partially obscured by the leaves and branches of a tree, or by falling snow or rain.

One of the strengths of SP-abstract and its realization in the SP computer model is that, in a similar way, recognition of a New pattern or patterns is not unduly disturbed by errors of omission, commission, and substitution in those data (Wolff, 2006, Chapter 6, Wolff, 2013, Section 4.2.2). This is because of the way the SP computer model searches for a global optimum in the building of multiple alignments, so that it does not depend on the presence or absence of any particular feature or combination of features in the New information that is being analyzed.

In its overall structure, SP-neural seems to lend itself to that kind of robustness in recognition in the face of errors in data. But the devil is in the detail. In further development of the theory, and in the development of a computer model of SP-neural, it will be necessary to clarify the details of how that kind of robustness may be achieved. In shaping this aspect of SP-neural, the principles that have been developed in SP-abstract are likely to prove useful and, with empirical evidence from brains and nervous systems, they may serve as a touchstone of success.

13. CONCLUSION

As was mentioned in the Introduction, SP-neural is a tentative and partial theory. That said, the close relationship between SP-neural and SP-abstract, the incorporation into SP-abstract of many insights from research on human perception and cognition, strengths of SP-abstract in terms of simplicity and power (Section 3.8.1), and advantages of SP-abstract compared with other AI-related systems (Section 3.8.3)—lend support to SP-neural as it is now as a conceptual model of the representation and processing of knowledge in the brain, and a promising basis for further research.

Naturally, we may have more confidence in some parts of the theory than others. Arguably, the parts that inspire most confidence are these:

- *Neural symbols and pattern assemblies.* All knowledge is represented in the cerebral cortex with *pattern assemblies*, the neural equivalent of patterns in SP-abstract. Each such pattern assembly is an array of *neural symbols*, each of which is a single neuron or a small cluster of neurons—the neural equivalent of a symbol in SP-abstract. Topologically, each array has one or two dimensions, perhaps parallel to the surface of the cortex.
- *Information compression via the matching and unification of patterns.* As in SP-abstract, SP-neural is governed by the overarching principle that many aspects of perception

and cognition may be understood in terms of information compression via the matching and unification of patterns.

- *Information compression via multiple alignment.* More specifically, SP-neural is governed by the overarching principle that many aspects of perception and cognition may be understood via a neural equivalent of the powerful concept of *multiple alignment*.
- *Unsupervised learning.* As in SP-abstract, unsupervised learning in SP-neural is the foundation for other kinds of learning—supervised learning, reinforcement learning, learning by imitation, learning by being told, and so on. And as in SP-abstract, unsupervised learning in SP-neural is achieved via a search through alternative grammars to find one or two that score best in terms of the compression of sensory information. As noted in Section 10.5, this is quite different from the kinds of “Hebbian” learning that are popular in artificial neural networks.
- *Problems of speed and expressiveness in the creation of pattern assemblies and their interconnections.* To account for the speed with which we can assimilate new information, and the speed of other mental processes (Section 11), it seems necessary to suppose that pattern assemblies and their interconnections may be created from pre-existing neural structures by the making and breaking of synaptic connections, somewhat like the making and breaking of connections in a telephone exchange, or the creation of a bespoke electronic system from an “uncommitted logic array” (ULA) or a “field-programmable gate array” (FPGA).

As with SP-abstract, areas of uncertainty in SP-neural may be clarified by casting the theory in the form of a computer model and testing it to see whether or not it works as anticipated. It is envisaged that this would be part of a proposed facility for the development of the SP machine (Wolff and Palade, 2016), a means for researchers everywhere to explore what can be done with the SP machine and to create new versions of it.

At all stages in its development, the theory may suggest possible investigations of the workings of brains and nervous systems. And any neurophysiological evidence may have a bearing on the perceived validity of the theory and whether or how it may need to be modified.

AUTHOR CONTRIBUTIONS

This is part of a long term research programme by JGW, developing the SP theory of intelligence. SP-neural, which is the subject of this paper, was first outlined in Chapter 11 of “Unifying Computing and Cognition” and, in this paper, has been considerably refined and developed.

FUNDING

The research is funded by CognitionResearch.org.

ACKNOWLEDGMENTS

I’m grateful to referees for constructive comments on earlier drafts of this paper.

REFERENCES

- Abraham, W. C. (2003). How long will long-term potentiation last? *Philos. Trans. R. Soc. B* 358, 735–744. doi: 10.1098/rstb.2002.1222
- Attneave, F. (1954). Some informational aspects of visual perception. *Psychol. Rev.* 61, 183–193. doi: 10.1037/h0054663
- Baars, B. J., and Gage, N. M. (2010). *Cognition, Brain, and Consciousness: Introduction to Cognitive Neuroscience, 2nd Edn.* Amsterdam: Elsevier.
- Barghout-Stein, L. (1999). *On Differences between Peripheral and Foveal Pattern Masking.* Technical report, University of California, Berkeley, Master's thesis. Available online at: bit.ly/1SCoUO4
- Barlow, H. B. (1959). "Sensory mechanisms, the reduction of redundancy, and intelligence," in *The Mechanisation of Thought Processes* (London: Her Majesty's Stationery Office), 535–559.
- Barlow, H. B. (1969). "Trigger features, adaptation and economy of impulses," in *Information Processes in the Nervous System*, ed K. N. Leibovic (New York, NY: Springer), 209–230.
- Barlow, H. B. (1982). David Hubel and Torsten Wiesel: their contribution towards understanding the primary visual cortex. *Trends Neurosci.* 5, 145–152. doi: 10.1016/0166-2236(82)90087-X
- Barrow, J. D. (1992). *Pi in the Sky.* Harmondsworth: Penguin Books.
- de Penning, H. L. H., d'Avila Garcez, A. S., Lamb, L. C., and Meyer, J.-J. C. (2011). "A neural-symbolic cognitive agent for online learning and reasoning," in *Proceedings of the International Joint Conferences on Artificial Intelligence* (Palo Alto, CA), 1653–1658.
- Drew, T., Vö M. L.-H., and Wolfe, J. M. (2013). The invisible gorilla strikes again: sustained inattention blindness in expert observers. *Psychol. Sci.* 24, 1848–1853. doi: 10.1177/0956797613479386
- d'Avila Garcez, A., Besold, T. R., de Raedt, L., Földiák, P., Hitzler, P., Icard, T., et al. (2015). "Neural-symbolic learning and reasoning: contributions and challenges," in *Proceedings of the AAAI Spring Symposium on Knowledge Representation and Reasoning, 2015*, ed T. Walsh (Stanford, CA), 18–21.
- d'Avila Garcez, A. S. (2007). "Advances in neural-symbolic learning systems: modal and temporal reasoning," in *Perspectives of Neural-Symbolic Integration*, eds B. Hammer and P. Hitzler (Heidelberg), 265–282.
- d'Avila Garcez, A. S., Lamb, L. C., and Gabbay, D. M. (2009). *Neural-Symbolic Cognitive Reasoning.* Heidelberg: Springer.
- Gold, M. (1967). Language identification in the limit. *Inform. Control* 10, 447–474. doi: 10.1016/S0019-9958(67)91165-5
- Gross, C. G. (2002). Genealogy of the "Grandmother Cell". *Neuroscientist* 8, 512–518. doi: 10.1177/107385802237175
- Hebb, D. O. (1949). *The Organization of Behaviour.* New York, NY: John Wiley & Sons.
- Herculano-Houzel, S. (2012). The remarkable, yet not extraordinary, human brain as a scaled-up primate brain and its associated cost. *Proc. Natl. Acad. Sci. U.S.A.* 109(Suppl. 1), 10661–10668. doi: 10.1073/pnas.1201895109
- Huth, A. G., de Heer, W. A., Griffiths, T. L., Theunissen, F. E., and Gallant, J. L. (2016). Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature* 532, 453–458. doi: 10.1038/nature17637
- Isaacson, J. S., and Scanziani, M. (2011). How inhibition shapes cortical activity. *Neuron* 72, 231–243. doi: 10.1016/j.neuron.2011.09.027
- Komendantskaya, E., Lane, M., and Seda, A. K. (2007). "Connectionist representation of multi-valued logic programs," in *Perspectives of Neural-Symbolic Integration*, eds B. Hammer and P. Hitzler (Heidelberg), 283–313. doi: 10.1007/978-3-540-73954-8_12
- McCorduck, P. (2004). *Machines Who Think: A Personal Inquiry Into the History and Prospects of Artificial Intelligence, 2nd Edn.* Natick, MA: A. K. Peters Ltd.
- Milner, P. M. (1996). Neural representations: some old problems revisited. *J. Cogn. Neurosci.* 8, 69–77. doi: 10.1162/jocn.1996.8.1.69
- Newell, A. (1973). "You can't play 20 questions with nature and win: projective comments on the papers in this symposium," in *Visual Information Processing*, ed W. G. Chase (New York, NY: Academic Press), 283–308.
- Newell, A. (ed.). (1990). *Unified Theories of Cognition.* Cambridge, MA: Harvard University Press.
- Ogasawara, H., and Kawato, M. (2010). The protein kinase m ζ network as a bistable switch to store neuronal memory. *BMC Syst. Biol.* 4:181. doi: 10.1186/1752-0509-4-181
- Ratliff, F., Hartline, H. K., and Miller, W. H. (1963). Spatial and temporal aspects of retinal inhibitory interaction. *J. Opt. Soc. Am.* 53, 110–120. doi: 10.1364/JOSA.53.000110
- Roy, A. (2013). An extension of the localist representation theory: grandmother cells are also widely used in the brain. *Front. Psychol.* 4:300. doi: 10.3389/fpsyg.2013.00300
- Schmidhuber, J. (2015). Deep learning in neural networks: an overview. *Neural Netw.* 61, 85–117. doi: 10.1016/j.neunet.2014.09.003
- Shamma, S. A. (1985). Speech processing in the auditory system II: lateral inhibition and the central processing of speech evoked activity in the auditory nerve. *J. Acoust. Soc. Am.* 76, 1622–1632. doi: 10.1121/1.392800
- Simons, D. J., and Ambinder, M. S. (2005). Change blindness: theory and consequences. *Curr. Direct. Psychol. Sci.* 14, 44–48. doi: 10.1111/j.0963-7214.2005.00332.x
- Squire, L. R., Berg, D., Bloom, F. E., du Lac, S., Ghosh, A., and Spitzer, N. C. (eds.). (2013). *Fundamental Neuroscience, 4th Edn.* Amsterdam: Elsevier.
- Stratton, G. M. (1897). Upright vision and the retinal image. *Psychol. Rev.* 4, 182–187. doi: 10.1037/h0064110
- Tahvildari, B., Fransén, E., Alonso, A. A., and Hasselmo, M. E. (2007). Switching between "on" and "off" states of persistent activity in lateral entorhinal layer III neurons. *Hippocampus* 17, 257–263. doi: 10.1002/hipo.20270
- von Békésy, G. (1967). *Sensory Inhibition.* Princeton, NJ: Princeton University Press.
- Wolff, J. G. (1988). "Learning syntax and meanings through optimization and distributional analysis," in *Categories and Processes in Language Acquisition*, eds Y. Levy, I. M. Schlesinger, and M. D. S. Braine (Hillsdale, NJ: Lawrence Erlbaum), 179–215. Available online at: bit.ly/ZIGjyc
- Wolff, J. G. (2006). Medical diagnosis as pattern recognition in a framework of information compression by multiple alignment, unification and search. *Decis. Support Syst.* 42, 608–625. doi: 10.1016/j.dss.2005.02.005
- Wolff, J. G. (2006). *Unifying Computing and Cognition: the SP Theory and Its Applications.* Menai Bridge: CognitionResearch.org. ISBNs: 0-9550726-0-3 (ebook edition), 0-9550726-1-1 (print edition). Distributors, including Amazon.com. Available online at: bit.ly/WmB1rs
- Wolff, J. G. (2007). Towards an intelligent database system founded on the SP theory of computing and cognition. *Data Knowl. Eng.* 60, 596–624. doi: 10.1016/j.datak.2006.04.003
- Wolff, J. G. (2013). The SP theory of intelligence: an overview. *Information* 4, 283–341. doi: 10.3390/info4030283
- Wolff, J. G. (2014a). Application of the SP theory of intelligence to the understanding of natural vision and the development of computer vision. *SpringerPlus* 3, 552–570. doi: 10.1186/2193-1801-3-552
- Wolff, J. G. (2014b). Autonomous robots and the SP theory of intelligence. *IEEE Access* 2, 1629–1651. doi: 10.1109/ACCESS.2014.2382753
- Wolff, J. G. (2014c). Big data and the SP theory of intelligence. *IEEE Access* 2, 301–315. doi: 10.1109/ACCESS.2014.2315297
- Wolff, J. G. (2014d). *Information Compression, Intelligence, Computing, and Mathematics.* Technical report, CognitionResearch.org. Available online at: bit.ly/1jEoECH
- Wolff, J. G. (2014e). The SP theory of intelligence: benefits and applications. *Information* 5, 1–27. Available online at: bit.ly/1lcquWF
- Wolff, J. G. (2016). The SP theory of intelligence: its distinctive features and advantages. *IEEE Access* 4, 216–246. doi: 10.1109/ACCESS.2015.2513822
- Wolff, J. G., and Palade, V. (2016). *Short Proposal for the Development of the SP Machine.* Technical report, CognitionResearch.org. Available online at: bit.ly/1SKAjhZ

Conflict of Interest Statement: The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2016 Wolff. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

APPENDIX

Cell Assemblies and Pattern Assemblies

The main differences between Hebb's (1949) concept of a "cell assembly" and the SP-neural concept of a "pattern assembly" are:

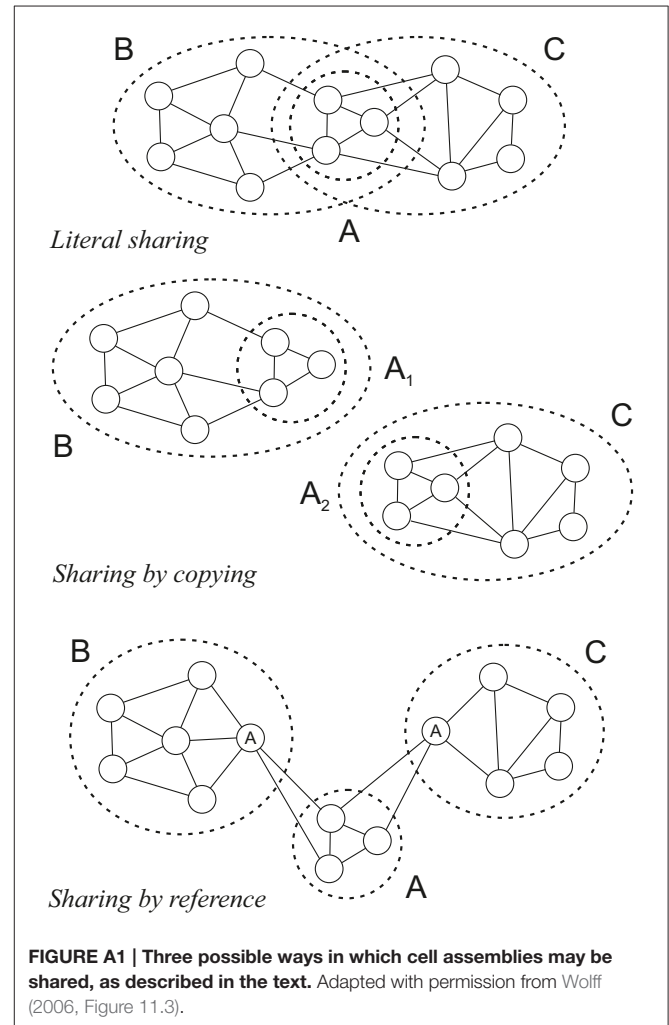
- The concept of a pattern assembly has had the benefit of computer modeling of SP-abstract—reducing vagueness in the theory and testing whether or not proposed mechanisms actually work as anticipated. These things would have been difficult or impossible for Hebb to do in 1949.
- Cell assemblies were seen largely as a vehicle for recognition, whereas, as neural realizations of SP "patterns," pattern assemblies should be able to mediate several aspects of intelligence, including recognition.
- Anatomically, pattern assemblies are seen as largely flat groupings of neurons in the cerebral cortex (Section 4.2), whereas cell assemblies are seen as structures in three dimensions.
- As described below, a fourth difference between cell assemblies and pattern assemblies is in how structures may be shared.

With regard to the last point, possible models for sharing of structures are illustrated in **Figure A1**.

In literal sharing, structures B and C in the figure both contain structure A. In sharing by copying, structures B and C each contains a copy of structure A. While in sharing by reference, structures B and C each contains a reference to structure A, in much the same way that a paper like this one contains references to other publications.

From Hebb's (1949) descriptions of the cell assembly concept, it is difficult to tell which of these three possibilities are intended.

By contrast with the concept of a pattern assembly in SP-neural, sharing is almost always achieved by means of neural "references" between structures. For example, a noun like "table" is likely to have neural connections to the many grammatical contexts in which it may occur, as suggested by the two broken-line connections from each of "N" and "#N" in the pattern assembly for "table" shown in **Figure 3**. Notice that, in this example, the putative direction of travel of nerve impulses is not relevant—it is the neural connection that counts.



In the SP system, it is intended that literal sharing should be impossible and that sharing by copying may only occur on the relatively rare occasions when the system has failed to detect the corresponding redundancy, and not always then.