

## **Title**

**Symbolic functions from neural computation**

## **Author**

Paul Smolensky (smolensky@jhu.edu)

## **Short title**

Symbolic functions from neural computation

## **Keywords**

cognitive science, neural networks, connectionism, linguistics, language, Optimality Theory

## **Abstract**

Is thought computation over ideas? Turing, and many cognitive scientists since, have assumed so, and formulated computational systems in which meaningful concepts are encoded by symbols which are the objects of computation. Cognition has been carved into parts, each a function defined over such symbols. This article reports on a research program aimed at computing these symbolic functions without computing over the symbols. Symbols are encoded as patterns of numerical activation over multiple abstract neurons, each neuron simultaneously contributing to the encoding of multiple symbols. Computation is carried out over the numerical activation values of such neurons, which individually have no conceptual meaning. This is massively parallel numerical computation operating within a continuous computational medium. The article presents an axiomatic framework for such a computational account of cognition, including a number of formal results. Within the framework, a class of recursive symbolic functions can be computed; formal languages defined by symbolic rewrite rules can also be specified, the sub-symbolic computations producing symbolic outputs which simultaneously display central properties of both facets of human language: universal symbolic grammatical competence and statistical, imperfect performance.

The classical theory of computation initiated by Turing [1] has provided the foundation for mainstream theories of cognition and intelligence since the inception of cognitive science and artificial intelligence in the 1950s [2]. This article presents another type of computation which has developed more recently within cognitive science, driven primarily by two goals: improved formal characterization of human mental processing—including mental grammars—and improved reduction to neural computation. These newer methods develop a distinction, absent in earlier computational theory, between three levels of description, which I will here call *symbolic*, *vectorial*, and *neural*. These distinguish the formal characterization of the mind, using recursive functions over discrete symbols, from the characterization of the brain, using continuous, parallel, numerical computation. The vectorial level provides an interlingua. The descriptions at all three levels characterize a single computer, the mind/brain [3].

A key departure from earlier theory (including Turing’s work on network machines [4]) is that machine computations operate only at a level lower than that of meaning. A conceptually meaningful entity is realized as an activation pattern—an activation vector—over many abstract neurons, each neuron simultaneously participating in the realization of multiple meaningful entities. Knowledge is realized in connections between elements (neurons) that are individually meaningless but collectively meaningful. The objects of computational manipulation are not meaningful (the neural-level description), and the meaningful elements are not objects of computational manipulation (the symbolic-level description) [5, 6, 7]. Nonetheless, we will see that the proposed vectorial computation provides the technical bridge for computing meaningful recursive symbolic functions using subsymbolic, neural computation.

## 0. Overview

The theory of human cognition owes much of its progress to the principle (formalized in Section 1 below) that cognition within some domain (e.g., arithmetic) consists in a system of functions that take inputs, and produce outputs, which are structures built of symbols that typically refer to elements in the domain (e.g., numbers).<sup>1</sup> This raises two basic questions—and a meta-question: How should these functions be specified by cognitive scientists? How are they computed within the brain? How are these two questions connected? The answers developed here are roughly these.

*How should cognitive functions be specified?* In Section 4, we adopt the conventional assumption that cognitive functions can usefully be specified via recursive equations; but in Section 5, we conclude that—at least in one key domain, the theory of human language—an alternative—specification by constraints on (input, output) pairs—has important advantages.

*How are cognitive functions computed?* In Section 4 we challenge the assumption that cognitive functions specified via recursive equations should be computed in the conventional manner, by sequentially computing the primitive functions in terms of which those equations are ultimately defined. Instead, we show how cognitively-relevant functions in a certain class can be computed in one, massively parallel step, in which neural activations encoding the input symbol structure are transformed by a simple linear transformation to output activations that encode the output symbol structure specified by the target function. In this approach, symbolic functions are computed, but not by symbolic computation: there is no algorithm over symbols that describes the internal processes by which input becomes output. There *is* a formal specification of the process, however, using the primitives of neural computation. If this account of cognition is on the right track, the kinds of symbol-manipulating algorithms sought by traditional artificial intelligence and cognitive theories do not exist, despite the existence of symbolic descriptions of cognitive functions.

*How are these questions connected?* If symbolic cognitive functions are computed by neural, not symbolic, computation, then natural descriptions of the functions computed by neural networks provide a promising candidate framework for specifying cognitive functions. This leads, in Section 5, to the description of these functions in terms of optimization over constraints, at all three levels.

Addressing these questions in Sections 4–5 requires laying considerable groundwork, the foundation of which is a mapping from a space of symbol structures to a vector space. This embedding is presented in Section 2, and the relation between the similarity of symbol structures and the similarity of their embeddings in the vector space is then analyzed in Section 3.

---

<sup>1</sup> In accepting this principle, we differ from most research programs on neural network (or ‘connectionist’) theories of cognition, which deny the validity of our symbolic level (e.g., [8]).

The paper provides a concise synopsis of the formal core of *The Harmonic Mind* [9], in which many of the results presented here are derived. Integrated into this synopsis are new results (Sections 3, 5.3.4) and current research (Section 6). The article initiates a research program developing axiomatic theory in cognitive science, motivated by the need for more precise characterizations of the central notions of cognitive theory, which are highly abstract and therefore in need of formal support to escape crippling vagueness.

### 0.1. Summary of results

In this article we take up the following topics: the decomposition of cognition into component functions (Section 1); the representation of the inputs and outputs of these cognitive functions (Section 2); the similarity structure of these representations (Section 3); a recursive-function-theoretic approach to specifying cognitive functions, and neural computation of those functions (Section 4); a formal-language-theoretic approach to specifying cognitive functions, and neural computation of those functions (Section 5); and the challenge of producing symbolically-interpretable outputs from neural computation (Section 6).

Each of these topics (except the last) is treated at three levels of description, in separated subsections: 1, the abstract symbolic level; 3, the neural level; and, mediating between them, 2, the vectorial level.<sup>2</sup>

The main results, simplified, are approximately these (numbers identify the corresponding theorems):

- 2.2.5 *Thm.* The fundamental symbolic data structure we assume, binary trees, can be embedded in an infinite-dimensional vector space in such a way that the basic data construction and data access operations are linear transformations.
- 3.2.3 *Thm.* Dissimilarity between two symbolic representations—as measured by differences in their constituents and the structural relations between their constituents—is given by the distance in the vector space between the two vectors “embedding” those representations.
- 4.2.1 *Thm.* Each function in the class defined as the closure under composition of the basic binary tree and symbol-mapping operations can be computed by a linear transformation over the vectors embedding the function’s input and output.
- 4.3.1 *Thm.* The infinite matrices implementing these linear transformations have a particularly simple form, the tensor product of the infinite identity matrix and a finite matrix that specifies the particular recursive function.
- 5.1.2 *Thm.* A formal language—defined standardly by sequentially-applied symbol-rewriting rules—can be specified as the structures that optimize (maximize) a numerical measure of symbolic well-formedness called (symbolic) Harmony.
- 5.1.5 *Thm.* A natural language specified in a new grammatical framework called Optimality Theory can also be characterized as the optima of symbolic Harmony.
- 5.2.2 *Thm.* The symbolic Harmony of a structure can be computed as the value, for the vector that embeds that structure, of a quadratic form called (neural) markedness Harmony.
- 5.3.3 *Thm.* In evaluating errors, the dissimilarity between an output symbolic structure and an input symbolic structure considered as the desired output can be computed as the value, at the vectors embedding the input and output structures, of a bilinear form called neural faithfulness Harmony.
- 5.3.4 *Thm.* The local optima of total network Harmony—the sum of markedness and faithfulness network Harmony—can be computed by a *deterministic* neural network.
- 5.3.5 *Thm.* The global optima of total network Harmony can be computed by a *stochastic* neural network.
- 6.1.1 *Thm.* The requirement that the output of a neural computation be a vector that is the embedding of a symbol structure can be met via a network dynamics that creates an attractor at each such vector.

To obtain (or merely state) these results as theorems, a certain degree of formalization is necessary.

Through the course of the article it may prove useful to refer back to the following synopsis, in which organization of results is orthogonal to that of the paper: it gives the total picture at each level separately. It also identifies all the principal elements of the theory, providing an informal glossary of their notation.

At this point, the reader may skip directly to Section 1 without loss of continuity.

---

<sup>2</sup> A notational challenge arises because most elements of the theory need three symbols, one for each level. The choice here is to use essentially the same symbol for all three, but to distinguish them by font and/or style (italics, bold, etc.).

## 0.2. Synopsis by level

At the most abstract level, the different types of information encoded mentally, and their mutual relations, are located within a cognitive macro-architecture graph  $\mathfrak{A}$  (1.1); these mental representations are characterized as systems  $\mathcal{S}$  of structures built of symbols filling specified structural roles (2.1); dissimilarity metrics, motivated by empirical patterns of cognitive performance, are defined w.r.t. a set of abstract relations  $\{R_i\} = \mathcal{R}$  among constituents of mental representations (3.1); the symbolic functions computed over  $\mathcal{S}$  within  $\mathfrak{A}$  are characterized, including a set  $\mathcal{F}$  of recursive functions (4.1) and the set of formal languages  $\mathcal{L}_G$  generated through sequential derivation by rewrite-rule grammars  $\mathcal{G}$ ; these languages and dissimilarity metrics are cast in the form of numerical Harmony functions  $H_G$  (5.1.2) and  $H_{\mathcal{R}}$  (3.1.4), the correct mental representations being those that maximize Harmony; Harmony functions  $H_G$  can be used to characterize the grammars of natural languages, including Optimality-Theoretic grammars  $\mathbb{G}$  within a system  $\mathfrak{D}$  in terms of which the typology of possible human languages can be formally calculated (5.1).

At a level intermediate between the symbolic and the neural, linear and quasi-linear classes of dynamical systems in vector spaces are defined (1.2); a model of the symbol system of mental states  $\mathcal{S}$  is specified abstractly by a realization mapping  $\Psi$  of symbol structures into a vector space  $S$  (2.2); symbolic structural-relation dissimilarity  $H_{\mathcal{R}}$  is realized as vector distance within  $S$  (3.2); recursive functions  $g \in \mathcal{F}$  and Harmony functions  $H_G$  are reduced to linear transformations on  $S$ ,  $W_g$  (4.2) and  $W_G$  (5.2.2).

At the neural level,  $S$  is modeled with  $\mathbb{R}^n$ , the states of a neural network  $\mathcal{N}$ , by positing a distinguished neural basis (1.3) in terms of which: mental representations  $\mathbf{s}$  and the realization mapping  $\Psi$  are explicitly specified numerically (2.3); Euclidean distance is explicitly computable, and  $H_{\mathcal{R}}$  is reduced to neural faithfulness Harmony  $H_F$  (5.3.3); the linear transformations  $W_g$  and  $W_G$  are instantiated as numerical matrices  $\mathbf{W}_g$  (4.3) and  $\mathbf{W}_G$  (5.3);  $H_G$  is reduced to neural markedness Harmony  $H_M$  (5.2.1); deterministic spreading activation can be used to compute local maxima of the network Harmony  $H_{\mathcal{N}} = H_F + H_M$  (5.3.4); and stochastic spreading activation dynamics  $\mathcal{D}_{\text{opt}}$  within a diffusion network  $\mathcal{N}^T$  can be used to compute global maxima of network Harmony (5.3.5); but a (deterministic) quantization dynamics  $\mathcal{D}_{\text{quant}}$  is needed to output an activation pattern  $\mathbf{s}$  that is interpretable as a discrete symbolic state  $s$ , with  $\mathbf{s} = \Psi(s)$  (6.1.1); a total dynamics combining  $\mathcal{D}_{\text{opt}}$  and  $\mathcal{D}_{\text{quant}}$  yields the  $\lambda$ -diffusion dynamics  $\mathcal{D}_{\lambda}$  (6.1.2) of a network  $\mathfrak{N}$  which computes symbolic mental representations, enabling models simultaneously capturing central features of both the idealized computational capacity of human linguistic competence and the errorful real-time computation of human linguistic performance.

## 1. Cognitive macro-architecture

At a coarse level of description, a cognitive architecture is an interconnected collection of components, operating in parallel, each of which processes information of a certain type (e.g., for language, these types include orthographic, phonological, syntactic, and semantic information: see ① in Figure 1). Knowledge of the structure of a type of information is contained within the corresponding component, and knowledge of the relations between the types is contained in the links connecting components. A component computes a function: the input comes along inward-directed links from neighboring components, and the output is in turn processed by knowledge in the outward-directed links to generate input to other components. The flowcharts ('box and arrow' diagrams) conventionally used to depict such an architecture instantiate a kind of graph structure.

### 1.1. Symbolic graph structure

1.1.1. *Def.* A cognitive macro-architecture  $\mathfrak{A}$  is a directed graph together with:

- a. For each node ('cognitive component' or 'module')  $\mathcal{M}^{\gamma}$  in  $\mathfrak{A}$ ,
  - i. a set  $\mathcal{S}^{\gamma}$ , the 'state space' of  $\mathcal{M}^{\gamma}$ ; an  $s^{\gamma} \in \mathcal{S}^{\gamma}$  is called a 'mental representation'
  - ii. a set  $\mathcal{I}^{\gamma}$ , the 'input space' of  $\mathcal{M}^{\gamma}$
  - iii. a function  $\sigma^{\gamma}: 2^{\mathcal{I}^{\gamma}} \rightarrow \mathcal{I}^{\gamma}$  for aggregating 'partial inputs' to  $\mathcal{M}^{\gamma}$  (see b.ii below)
- b. For each edge ('connection' or 'pathway')  $\mathcal{C}^{\gamma\alpha}$ , from  $\mathcal{M}^{\alpha}$  to  $\mathcal{M}^{\gamma}$ , in  $\mathfrak{A}$ ,
  - i. a transformation  $\tau^{\gamma\alpha}: \mathcal{S}^{\alpha} \rightarrow \mathcal{I}^{\gamma}$  (the 'inter-component transformation');
  - ii.  $\iota^{\gamma\alpha} \equiv \tau^{\gamma\alpha}(s^{\alpha}) \in \mathcal{I}^{\gamma}$  is called the 'partial input' from  $\mathcal{M}^{\alpha}$  to  $\mathcal{M}^{\gamma}$  when  $\mathcal{M}^{\alpha}$  is in state  $s^{\alpha} \in \mathcal{S}^{\alpha}$
- c.  $\iota^{\gamma} \equiv \sigma^{\gamma}(\{\iota^{\gamma\alpha} \mid \text{edge } \mathcal{C}^{\gamma\alpha} \text{ is in } \mathfrak{A}\}) \in \mathcal{I}^{\gamma}$  is the 'external input' to  $\mathcal{M}^{\gamma}$

Henceforth we assume given some specific architecture  $\mathfrak{A}$  in which, for all  $\mathcal{M}^{\gamma} \in \mathfrak{A}$ , the spaces  $\mathcal{I}^{\gamma}$  and  $\mathcal{S}^{\gamma}$ , while conceptually distinct, are formally identical;  $\mathcal{I}^{\gamma}$  functions as a 'target' towards which the compo-

nents external to  $\mathcal{M}^\gamma$  drive its state  $s^\gamma$ . Knowledge within  $\mathcal{M}^\gamma$  will generally partially resist this pressure.

### 1.2. Vectorial model

The vectorial level models each  $S^\gamma = \mathcal{I}^\gamma$  as a vector space  $S^\gamma = I^\gamma$  (see Figure 1 ②). The vector space  $S^\gamma$  constitutes an abstract information-encoding medium—a space of cognitive states, of mental representations—that resides in continuous mathematics, rather than the discrete mathematics of symbolic encodings. We study two example types of vectorial models here.

In quasi-linear models, the transformations between components are linear, combining by addition; a potentially non-linear function  $\mathbf{f}$  and a linear transformation of  $S^\gamma$  to itself determines how the state of  $S^\gamma$  at one time, and the external input to  $S^\gamma$  from other components at that time, move the state forward:

- 1.2.1. *Ex.* In a *quasi-linear dynamics*  $\mathcal{D}$  [10] for the architecture  $\mathfrak{A}$ , for all cognitive components  $\mathcal{M}^\gamma$  and  $\mathcal{M}^\alpha$ :
- the operation  $\tau^{\gamma\alpha}$  on the link to  $\mathcal{M}^\gamma$  from  $\mathcal{M}^\alpha$  is a linear transformation  $T^{\gamma\alpha}: S^\alpha \rightarrow I^\gamma$ ,  $\mathbf{t}^{\gamma\alpha} = T^{\gamma\alpha} \mathbf{s}^\alpha$
  - $\sigma^\gamma$  is summation: the external input to  $\mathcal{M}^\gamma$ ,  $\mathbf{t}^\gamma$ , is the sum of the partial inputs  $\mathbf{t}^{\gamma\alpha}$  from neighboring components  $\mathcal{M}^\alpha$ :  $\mathbf{t}^\gamma = \sum_\alpha \mathbf{t}^{\gamma\alpha}$
  - a linear transformation  $W^\gamma: S^\gamma \rightarrow S^\gamma$  (the ‘intra-component input transformation’) and a function  $\mathbf{f}^\gamma: S^\gamma \rightarrow S^\gamma$  (the ‘activation function’) define  $\mathcal{D}^\gamma$ , the internal dynamics of  $\mathcal{M}^\gamma$ , by:

$$d\mathbf{s}^\gamma/dt = \mathbf{f}^\gamma[\mathbf{i}^\gamma(t)] - \mathbf{s}^\gamma(t), \quad \mathbf{i}^\gamma(t) \equiv W^\gamma \mathbf{s}^\gamma(t) + \mathbf{t}^\gamma(t)$$

$W^\gamma \mathbf{s}^\gamma$  is the internally-generated input to  $\mathcal{M}^\gamma$ ; it combines linearly with the external input  $\mathbf{t}^\gamma$ .

Note that at equilibrium, we must have  $\mathbf{s}^\gamma(t) = \mathbf{f}^\gamma[\mathbf{i}^\gamma(t)]$ ;  $\mathbf{f}^\gamma$  gives the equilibrium relation between input  $\mathbf{i}^\gamma$  and activation  $\mathbf{s}^\gamma$

A simpler type of model we will use (ultimately, as exemplified in Figure 2 below) is:

- 1.2.2. *Ex.* A *linear associator* [11] is a simple sub-architecture consisting of two components,  $\mathcal{M}^\alpha$  and  $\mathcal{M}^\gamma$ , with a single (uni-directional) pathway, from  $\mathcal{M}^\alpha$  to  $\mathcal{M}^\gamma$ , in which
- the operation  $\tau^{\gamma\alpha}$  is a linear transformation  $T^{\gamma\alpha}: S^\alpha \rightarrow I^\gamma$ , producing external input  $\mathbf{t}^\gamma = T^{\gamma\alpha} \mathbf{s}^\alpha$
  - the state  $\mathbf{s}^\gamma$  equals the external input  $\mathbf{t}^\gamma$ :  $\mathbf{s}^\gamma = T^{\gamma\alpha} \mathbf{s}^\alpha$

Both types of model are spelled out more explicitly and intuitively at the neural level.

### 1.3. Neural computation: Cognitive micro-architecture

The neural level (Figure 1 ③) invokes a coordinate system (with axes defined by a ‘neural basis’  $\{\hat{\mathbf{e}}_{\mu}^{\gamma}\}_{\mu=1}^n$  [12]) for  $S^\gamma$  in which the list of coordinates for any vector  $\mathbf{s}^\gamma \in S^\gamma$ ,  $(s^{\gamma_1}, s^{\gamma_2}, \dots, s^{\gamma_n})$ , is the list of activation values of an enumerated set of  $n$  ‘abstract neurons’ in a network  $\mathcal{N}^\gamma$ ; in this model of the vectorial theory, the mental state  $s^\gamma$  is realized by an ‘activation pattern’ described by a vector in  $\mathbb{R}^n$ . We assume given a neural basis for every  $S^\gamma$ .

Generally, neural computation is a type of analog computation: a dynamical system in  $\mathbb{R}^n$ , fundamentally continuous in time, defined by differential equations  $\mathcal{D}^\gamma$  that describe how the  $n$  neurons function as parallel processors (how they ‘spread activation’; Figure 1 ④) [13]. A particular computation is specified by values for the parameters in  $\mathcal{D}^\gamma$ ; these parameter values are interpreted in parallel by the machine that evolves in accordance with  $\mathcal{D}^\gamma$ . In simpler cases, like the linear associator (1.2.2), there may be no dynamics, just static relations between component states.

- 1.3.1. *Def.* Given a quasi-linear dynamics  $\mathcal{D}^\gamma$  for  $S^\gamma$ , we define the following, relative to the neural basis:
- The inter-component linear transformation  $T^{\gamma\alpha}$  is realized as the matrix  $\mathbf{T}^{\gamma\alpha}$ , the element  $[\mathbf{T}^{\gamma\alpha}]_{\mu\nu}$  being the ‘connection strength’, or ‘weight’, of a ‘connection’ from neuron  $\nu$  in  $\mathcal{N}^\alpha$  to neuron  $\mu$  in  $\mathcal{N}^\gamma$ . (This applies to a linear associator as well.)
  - The linear intra-component input transformation  $W^\gamma$  is realized as the weight matrix  $\mathbf{W}^\gamma$ , the element  $[\mathbf{W}^\gamma]_{\mu\nu}$  being the weight of a connection from neuron  $\nu$  in  $\mathcal{N}^\gamma$  to neuron  $\mu$  in  $\mathcal{N}^\gamma$ .
  - The *total input* to the neurons in  $\mathcal{N}^\gamma$  is  $\mathbf{i}^\gamma \equiv \mathbf{W}^\gamma \mathbf{s}^\gamma + \mathbf{t}^\gamma$ , where  $\mathbf{t}^\gamma$  is the external input from other components, while  $\mathbf{W}^\gamma \mathbf{s}^\gamma$  (the matrix-vector product of  $\mathbf{W}^\gamma$  and  $\mathbf{s}^\gamma$ ) is the input generated internally within  $\mathcal{N}^\gamma$ . At neuron  $\mu$  of  $\mathcal{N}^\gamma$ , the total input is:  $[\mathbf{i}^\gamma]_\mu = \sum_\nu [\mathbf{W}^\gamma]_{\mu\nu} [s^\gamma]_\nu + [\mathbf{t}^\gamma]_\mu$ , i.e., the weighted sum of the activation values of  $\mu$ ’s neighbors  $\nu$  within  $\mathcal{N}^\gamma$  plus the external input to  $\mu$ .
- 1.3.2. *Def.* To the requirements defining a quasi-linear dynamics at the vectorial level (1.2.1) we now add:
- the *locality* requirement  $\mathbf{f}^\gamma: \mathbf{i}^\gamma = (i_1, i_2, \dots, i_n) \mapsto (f^\gamma(i_1), f^\gamma(i_2), \dots, f^\gamma(i_n))$  for some  $f^\gamma: \mathbb{R} \rightarrow \mathbb{R}$  called the ‘neuron activation function’; and

- b. the *monotonicity* requirement that the activation function  $f^\gamma$  be non-decreasing.  
If, in addition,  $f^\gamma$  is the identity function,  $f^\gamma(i) = i$ , then the dynamics  $\mathcal{D}^\gamma$  is *linear*.

Thanks to locality, the differential equation for neuron  $\mu$  depends only the activation levels of units connected to  $\mu$ , the weights on connections to neuron  $\mu$ , and the external input to unit  $\mu$ :

$$1.3.3. \quad ds_{\mu}^{\gamma}/dt = f^{\gamma}([\mathbf{i}^{\gamma}(t)]_{\mu}) - s_{\mu}^{\gamma}(t); [\mathbf{i}^{\gamma}(t)]_{\mu} = \sum_{\nu} [\mathbf{W}^{\gamma}]_{\mu\nu} s_{\nu}^{\gamma}(t) + [\mathbf{I}^{\gamma}(t)]_{\mu}$$

Monotonicity implies that, at equilibrium, the higher the input  $i_{\mu}^{\gamma}$  to neuron  $\mu$ , the greater its activation value  $s_{\mu}^{\gamma} = f(i_{\mu}^{\gamma})$ .

Henceforth, we focus primarily on a single component  $\mathcal{M}^{\gamma}$ , and generally drop  $\gamma$  from the notation.

## 2. Symbolic functions

A component  $\mathcal{M}$  of an architecture  $\mathcal{A}$  computes a function, producing a mental representation as output.

### 2.1. Mental representations as symbol structures

At the most abstract level, mental representations in  $\mathcal{S}$  are symbol structures [14] (Figure 1 ⑤); e.g., in an orthographic component, a string of mental letters like CABS; in a phonological component, a string of mental phonemes like kæbz (or less simplistically, the binary tree [k [æ [b z]]]); in a syntactic component, a phrase-structure tree like [S [NP Frodo] [VP [V lives]]] (simplifying greatly). The key idea now is to regard a symbol structure as a set of internal roles, each filled by a constituent symbol structure (see ⑤).

2.1.1. *Def.* A *filler-role decomposition*  $\mathfrak{S} = (\mathcal{S}, \mathcal{F}, \mathcal{R}, \beta)$  consists of three sets, the ‘structures’  $\mathcal{S}$ , the ‘fillers’  $\mathcal{F}$ , and the ‘roles’  $\mathcal{R}$ , and a one-to-one function  $\beta: \mathcal{S} \rightarrow 2^{\mathcal{F} \times \mathcal{R}}$ ; the  $N$  pairs constituting  $\beta(s)$ , written  $\{f_k/r_k\}_{k=1}^N$ , are the ‘filler/role bindings’, or ‘constituents’, of  $s \in \mathcal{S}$  [15].

A symbolic data structure that is widely applicable for mental representations is the binary tree; the artificial-intelligence language LISP deploys it exclusively [16]. Two filler-role decompositions are [17]:

2.1.2. *Ex.* The following defines the *canonical filler-role decomposition of binary trees*,  $\mathfrak{T}_t$ . Let the structure-set  $\mathcal{S}_t$  be the set of binary trees labeled with ‘atomic symbols’ in the set  $\mathcal{A} \equiv \{\mathfrak{a}, \mathfrak{b}, \mathfrak{k}\}$ , e.g.,  $s_{cab} \equiv [\mathfrak{k} [\mathfrak{a} \mathfrak{b}]] \in \mathcal{S}_t$ . Define the role-set  $\mathcal{R}_t \equiv \{r_x \mid x \in \{0,1\}^*\}$ ; we think of, e.g.,  $r_{01}$  as the binary-tree position ‘left child [0] of the right child [1] of the root’, and analogously for any string  $x$  of 0s and 1s. In  $s_{cab}$ , the position  $r_{01}$  is labeled with the symbol  $\mathfrak{a}$ . Let the filler-set be  $\mathcal{F}_t \equiv \mathcal{A}$ . Then  $\mathfrak{a}/r_{01}$  is a filler-role binding of  $s_{cab}$ ; the bindings-function  $\beta_t$  pairs roles with their labels:  $\beta_t(s_{cab}) = \{\mathfrak{a}/r_{01}, \mathfrak{k}/r_0, \mathfrak{b}/r_{11}\}$ . The root position is  $r_{\varepsilon}$ , where  $\varepsilon$  is the empty string; for any atomic symbol  $\mathfrak{x} \in \mathcal{A}$ ,  $\beta_t(\mathfrak{x}) = \{\mathfrak{x}/r_{\varepsilon}\}$ .

2.1.3. *Ex.* Continuing the example of  $\mathcal{S}_t$ , binary trees labeled with  $\mathcal{A}$ , we can also define the *recursive* filler-role decomposition  $\mathfrak{T}_r = (\mathcal{S}_r, \mathcal{F}_r, \mathcal{R}_r, \beta_r)$ . The roles are simply  $\mathcal{R}_r \equiv \{r_{\varepsilon}, r_0, r_1\}$ , while the fillers are  $\mathcal{F}_r \equiv \mathcal{S}_t$ . Now, for  $s_{cab}$  above,  $\beta_r(s_{cab}) = \{\mathfrak{k}/r_0, [\mathfrak{a} \mathfrak{b}]/r_1\}$ . Here the filler of  $r_1$ ,  $f_1 \equiv [\mathfrak{a} \mathfrak{b}]$ , is a non-atomic element of  $\mathcal{F}_r = \mathcal{S}_t$ , and itself has bindings  $\beta_r(f_1) = \{\mathfrak{a}/r_0, \mathfrak{b}/r_1\}$ ; the atomic filler  $\mathfrak{k}$  has binding-set  $\beta_r(\mathfrak{k}) = \{\mathfrak{k}/r_{\varepsilon}\}$ , as in  $\beta_t$ . As in LISP, we denote by *cons* the binary-tree constructor function:  $s_{cab} = \text{cons}(\mathfrak{k}, [\mathfrak{a} \mathfrak{b}]) = \text{cons}(\mathfrak{k}, \text{cons}(\mathfrak{a}, \mathfrak{b}))$ ; thus  $\beta_r(\text{cons}(s_0, s_1)) = \{s_0/r_0, s_1/r_1\}$  for all  $s_0, s_1 \in \mathcal{S}_t$ . The access function that extracts the left (right) child of a tree will be denoted  $\text{ex}_0$  ( $\text{ex}_1$ ):  $\mathcal{S}_t \rightarrow \mathcal{F}_r \cup \{\emptyset\}$  (both return a special ‘null filler’  $\emptyset$  when applied to an atom); thus, for non-atomic  $s \in \mathcal{S}_t$ ,  $\beta_r(s) = \{\text{ex}_0(s)/r_0, \text{ex}_1(s)/r_1\}$ .  $\text{ex}_{\varepsilon}(s)$  is the symbol bound to the root of the tree  $s$  (or  $\emptyset$  if there is no such symbol).

Henceforth, unless stated otherwise, we assume that  $\mathcal{S}$  is a given set of binary trees under the canonical filler-role decomposition  $\mathfrak{T}_t$ . We also assume throughout that, in general, the sets  $\mathcal{F}$  and  $\mathcal{R}$  are denumerable (not necessarily finite); henceforth, let  $\{\hat{f}_j\}$  and  $\{\hat{r}_k\}$  be given enumerations of them.

### 2.2. Symbol structures as vectors

The vectorial level of description springs from the following fundamental definition [15]: it asserts that the vector realizing (or modeling, or instantiating, or embedding) a symbol structure is the sum of vectors that realize each constituent filler-role binding of the structure, and that the vector realizing a filler-role binding is the tensor (generalized outer) product of vectors that realize the filler and the role.

2.2.1. *Def.* A *tensor-product realization*  $(\mathfrak{S}, F, R, \psi_F, \psi_R)$  consists of a filler-role decomposition  $\mathfrak{S} = (\mathcal{S}, \mathcal{F}, \mathcal{R}, \beta)$ , two real vector spaces  $F$  and  $R$ —the ‘filler vector space’ and the ‘role vector space’, with dimensions  $\dim(F)$  and  $\dim(R)$  respectively—and two ‘realization’ functions  $\psi_F: \mathcal{F} \rightarrow F$ ,  $\psi_R: \mathcal{R} \rightarrow R$ . Here, we require that each of the ranges of  $\psi_F$  and  $\psi_R$  be linearly independent sets.

The associated *realization mapping*  $\Psi: \mathcal{S} \rightarrow S \equiv F \otimes R$  is defined by

$$\Psi(s) \equiv \sum_{k=1}^N \mathbf{f}_k \otimes \mathbf{r}_k \quad \text{where } \beta(s) = \{f_k/r_k\}_{k=1}^N \text{ and } \mathbf{f}_k \equiv \Psi_F(f_k), \mathbf{r}_k \equiv \Psi_R(r_k)$$

The vector space  $S$ , containing the range of  $\Psi$ , is the tensor product of spaces  $F$  and  $R$ :  $S \equiv F \otimes R$ . [Given respective bases  $\{\hat{\mathbf{f}}_j\}$  and  $\{\hat{\mathbf{r}}_k\}$  for  $F$  and  $R$ ,  $\{\hat{\mathbf{f}}_j \otimes \hat{\mathbf{r}}_k\}$  is a basis for  $S$ , and the mapping  $(\mathbf{f}, \mathbf{r}) \mapsto \mathbf{f} \otimes \mathbf{r}$  from  $F \times R$  to  $S$  is bilinear: linear in each of  $\mathbf{f}$  and  $\mathbf{r}$  independently;  $\dim(S) = \dim(F) \dim(R)$ .]

Corresponding to our given enumerations  $\{\hat{f}_j\}$  and  $\{\hat{r}_k\}$  of  $\mathcal{F}$  and  $\mathcal{R}$  we have their vectorial realizations  $\{\hat{\mathbf{f}}_j \equiv \Psi_F(\hat{f}_j)\} \subset F$  and  $\{\hat{\mathbf{r}}_k \equiv \Psi_R(\hat{r}_k)\} \subset R$ .

Given a vector  $\mathbf{v} \in S$  realizing some symbol structure in  $\mathcal{S}$ , we can determine that structure from  $\mathbf{v}$ .

2.2.2. *Prop.* The linear independence of the ranges of  $\Psi_F$  and  $\Psi_R$  entails that  $\Psi_F$  and  $\Psi_R$  are invertible. Furthermore,  $\Psi$  is invertible: given  $\mathbf{v} = \Psi(s)$ ,  $s$  can be recovered as the unique element of  $\mathcal{S}$  with bindings  $\{f_k/\hat{r}_k\}$ , where  $f_k \equiv \Psi_F^{-1}(\mathbf{f}_k)$ , and  $\{\mathbf{f}_k\}$  is the unique sequence in  $F$  such that  $\mathbf{v} = \sum_k \mathbf{f}_k \otimes \hat{\mathbf{r}}_k$ . [15]

Within the continuous space of vectorial mental representations  $S$ , distance can be used to model cognitive dissimilarity, once  $S$  is endowed with a metric structure.

2.2.3. *Def.* A *metric vectorial realization* is a tensor-product realization in which each vector space  $V \in \{F, R\}$  has an *inner product*; the inner product of two vectors  $\mathbf{u}, \mathbf{v} \in V$  is written  $\mathbf{u} \cdot \mathbf{v}$  [ $(\mathbf{u}, \mathbf{v}) \mapsto \mathbf{u} \cdot \mathbf{v}$  is bilinear, symmetric, and positive-definite:  $\mathbf{u} \neq \mathbf{0} \Rightarrow \mathbf{u} \cdot \mathbf{u} > 0$ ]. The *length* of  $\mathbf{u}$  is  $\|\mathbf{u}\| \equiv (\mathbf{u} \cdot \mathbf{u})^{1/2}$ ;  $\mathbf{u}$  is *normalized* iff  $\|\mathbf{u}\| = 1$ . The *distance* between  $\mathbf{u}$  and  $\mathbf{v}$  is  $\|\mathbf{u} - \mathbf{v}\|$ .  $\mathbf{u}$  and  $\mathbf{v}$  are *orthogonal* iff  $\mathbf{u} \cdot \mathbf{v} = 0$ . Inner products on  $F, R$  induce an inner product on  $F \otimes R$  satisfying  $(\mathbf{f}_1 \otimes \mathbf{r}_1) \cdot (\mathbf{f}_2 \otimes \mathbf{r}_2) = (\mathbf{f}_1 \cdot \mathbf{f}_2)(\mathbf{r}_1 \cdot \mathbf{r}_2)$ .

Returning to the case of binary trees, we can now relate the two decompositions 2.1.2 and 2.1.3:

2.2.4. *Def.* A *recursive realization of binary trees* is a metric vectorial realization of  $\mathfrak{T}_t$ , built from a vector space  $R_{(1)}$ , in which

- $\Psi_R(r_0) \equiv \mathbf{r}_0 \in R_{(1)}$ ,  $\Psi_R(r_1) \equiv \mathbf{r}_1 \in R_{(1)}$
- $\mathbf{r}_{x0} = \mathbf{r}_x \otimes \mathbf{r}_0$  and  $\mathbf{r}_{x1} = \mathbf{r}_x \otimes \mathbf{r}_1$  for all  $x \in \{0, 1\}^*$  (where  $x0$  is the concatenation of string  $x$  and 0)
- The roles for tree positions at depth  $d$  lie in the vector space  $R_{(d)} \equiv R_{(1)} \otimes R_{(1)} \otimes \dots \otimes R_{(1)}$  ( $d$  factors); letting  $R_{(0)} \equiv \mathbb{R}$ , the total role space  $R$  is the direct sum of all the vector spaces  $R_{(d)}$ :
- $R \equiv R_{(0)} \oplus R_{(1)} \oplus R_{(2)} \oplus R_{(3)} \oplus \dots$  [an  $\mathbf{r} \in R$  can be represented  $\mathbf{r} = (\mathbf{r}_{(0)}, \mathbf{r}_{(1)}, \dots)$ , each  $\mathbf{r}_{(d)} \in R_{(d)}$ ].<sup>3</sup>

2.2.5. *Thm.* There are four linear transformations on  $S = F \otimes R - W_{\text{cons}0}, W_{\text{cons}1}, W_{\text{ex}0}, W_{\text{ex}1}$ —such that if  $s = \text{cons}(p, q)$ , and we define  $\mathbf{s} \equiv \Psi(s)$ ,  $\mathbf{p} \equiv \Psi(p)$ ,  $\mathbf{q} \equiv \Psi(q)$ , then [19]

$$\mathbf{s} = W_{\text{cons}0} \mathbf{p} + W_{\text{cons}1} \mathbf{q}; \quad \mathbf{p} = W_{\text{ex}0} \mathbf{s}; \quad \mathbf{q} = W_{\text{ex}1} \mathbf{s}; \quad \mathbf{s} = \mathbf{p} \otimes \mathbf{r}_0 + \mathbf{q} \otimes \mathbf{r}_1$$

This recursive realization of the canonical filler-role decomposition 2.1.2 makes it possible to write  $\mathbf{s}$  either in the form corresponding to the canonical decomposition,  $\mathbf{s} = \sum_x \mathbf{f}_x \otimes \mathbf{r}_x$ , where  $\mathbf{f}_x$  is the realization of the atomic symbol at tree position  $x$ , or as  $\mathbf{s} = \mathbf{p} \otimes \mathbf{r}_0 + \mathbf{q} \otimes \mathbf{r}_1$ , just as it would be expressed using the recursive decomposition  $\mathfrak{T}_t$ ,  $\mathbf{r}_0$  and  $\mathbf{r}_1$  realizing its two roles (left/right-child), and  $\mathbf{p}$  ( $\mathbf{q}$ ) realizing the entire left (right) sub-tree. In this sense, 2.2.4.b renders the canonical representation recursive.

Below, the linear transformations of 2.2.5 will enable computation of an interesting class of functions.

Henceforth we assume given metric vectorial realizations  $\Psi^S$  and  $\Psi^I$  of  $\mathcal{S}$  and  $\mathcal{I}$  which are defined over the same filler and role vector spaces  $F$  and  $R$ , and which obey the following *input scaling condition*:

2.2.6. For all  $f \in \mathcal{F}$ ,  $\Psi^I_F(f) = \Psi^S_F(f)$ ; there exists a function  $\tilde{\rho}: \mathcal{R} \rightarrow \mathbb{R}$  such that for all  $r \in \mathcal{R}$ ,  $\Psi^I_R(r) = \rho(r) \Psi^S_R(r)$ . Section 3.2 uses  $\tilde{\rho}$  to control the length of each input  $\mathbf{t}$  in computing its distance to a state  $\mathbf{s}$ .

### 2.3. Symbolic structures as neural activation patterns

Each vector space  $S^y$  has a given neural basis; the coordinates of  $\mathbf{s} = \Psi_S(s)$  w.r.t. this basis are the neural activation values realizing state  $s \in \mathcal{S}$  (1.3). If the individual vectors  $\{\hat{\mathbf{f}}_j\}$  and  $\{\hat{\mathbf{r}}_k\}$  realizing the individual symbolic fillers and roles each lie along a neural basis vector, then the presence of a given symbolic constituent corresponds to the activation of a single neuron: this is ‘local representation’. We will assume the general case, *distributed representation*, in which an individual constituent is realized by a vector that is a linear combination of multiple neural basis vectors: presence of that constituent is encoded by an activation pattern distributed over multiple neurons, and each neuron participates in the encoding of multiple constituents. (Distributed representations: allow many more representations per  $n$  neurons; allow en-

<sup>3</sup> Physicists will recognize this construction of an infinite-dimensional Hilbert space as isomorphic to Fock space, where  $R_{(d)}$  corresponds to the subspace of  $d$  particles, and  $\mathbf{f} \otimes \mathbf{r}$  corresponds to the tensor product binding, e.g., of the spin of an electron ( $\mathbf{f}$ ) to its orbital ( $\mathbf{r}$ ) in an atom [18].

coding of similarity among constituents; are what results from neural network learning; and are ubiquitous in the brain [20, 21, 22].) The  $[\varphi \times \varrho]^{\text{th}}$  activation value in the tensor product realization of a symbol structure is the product  $[\varphi^{\text{th}}$  activation in the realization of constituent  $k$ 's filler]  $\times$   $[\varrho^{\text{th}}$  activation in the realization of constituent  $k$ 's role], summed over all the constituents  $k$ .

2.3.1. *Def.* A ‘neural basis’  $\{\underline{\mathbf{f}}_\varphi\}, \{\underline{\mathbf{r}}_\varrho\}$  for the spaces  $F, R$  of a metric vectorial realization is a distinguished orthonormal basis. The neural coordinates, or activation pattern, realizing a symbol structure  $s \in \mathcal{S}$  is the point  $\mathbf{s} \in \mathbb{R}^{nm}$  with elements  $s_{\varphi\varrho}$  such that  $\mathbf{s} = \sum_{\varphi\varrho} s_{\varphi\varrho} \underline{\mathbf{f}}_\varphi \otimes \underline{\mathbf{r}}_\varrho$ .

2.3.2. *Prop.* For each binding of  $s$ ,  $\beta(s) = \{f_k/r_k\}$ , let  $[\mathbf{f}_k]_\varphi$  be the  $\varphi^{\text{th}}$  neural coordinate of  $\mathbf{f}_k \equiv \psi_F(f_k)$  (i.e.,  $\mathbf{f}_k = \sum_\varphi [\mathbf{f}_k]_\varphi \underline{\mathbf{f}}_\varphi$ ) and similarly let  $[\mathbf{r}_k]_\varrho$  be the  $\varrho^{\text{th}}$  neural coordinate of  $\mathbf{r}_k \equiv \psi_R(r_k)$ . Then  $s_{\varphi\varrho} = \sum_k [\mathbf{f}_k]_\varphi [\mathbf{r}_k]_\varrho$  [15].

### 3. Dissimilarity as representational distance: relational faithfulness

The total input  $\iota \in \mathcal{I}$  to a component  $\mathcal{M}$  is the target towards which the state of  $\mathcal{M}$  is driven under the combined influence of those components that send partial input to  $\mathcal{M}$ . As we see in Section 5, *all else equal*, the greater the distance  $d(s, \iota)$  of a state  $s \in \mathcal{S}$  to the total input  $\iota$ —i.e., the more dissimilar or ‘unfaithful’  $s$  is to  $\iota$ —the less likely  $\mathcal{M}$  is to be in state  $s$ . Experimental data provide evidence for uncovering the (component-specific) dissimilarity metric  $d$  in  $\mathcal{S}$ ; cognitive scientists formulate hypotheses concerning  $d$  in terms of the format of representations in  $\mathcal{S}$ : the format determines the dimensions of unfaithfulness relevant for  $d$ . We propose to formalize the problematic notion of representational format by making explicit those relations between constituents that are understood to be defined in that format (Figure 1 ⑦).

#### 3.1. Representational format as relational faithfulness

3.1.1. *Def.* A *representational format*  $\mathfrak{F}$  is a filler-role decomposition in which the roles are a set of *relations*  $\mathcal{R} = \{R_i\}$ . Each  $R_i$  has an ‘arity’  $n_i$ , and a collection of ‘argument domains’  $A_{i1}, \dots, A_{in_i}$ . Then for any particular symbol structure  $s \in \mathcal{S}$ ,  $R_i(s) \subset A_{i1} \times \dots \times A_{in_i}$  is the set of elements which in  $s$  stand in the relation  $R_i$ . The filler/role bindings are defined as  $\beta(s) = \cup_i \{(a_1, \dots, a_{n_i})/R_i \mid (a_1, \dots, a_{n_i}) \in R_i(s)\}$ .

3.1.2. *Ex.* Let  $\mathcal{S}_s \subset \mathcal{A}^*$  be the set of strings of symbols from the alphabet  $\mathcal{A}$  such that no symbol appears more than once in any string. For any  $s \in \mathcal{S}_s$ , define the arity-1 relation  $R_A(s) \equiv \{X \in s\} \subset \mathcal{A} \equiv \mathbf{A}_{A1}$ , i.e.,  $R_A(s)$  is the set of symbols that occur (once) in the string  $s$ . Define the arity-2 relation  $R_L(s) \equiv \{(l, X) \mid X \text{ is the symbol of } s \text{ in position } l \text{ (relative to the left edge)}\} \subset \mathbb{N} \times \mathcal{A} \equiv \mathbf{A}_{L1} \times \mathbf{A}_{L2}$ ; e.g.,  $R_L(\text{CAB}) = \{(2, A), (1, C), (3, B)\}$ .  $\mathcal{R}_s \equiv \{R_A, R_L\}$  is a representational format for  $\mathcal{S}_s$ , yielding a bindings function  $\beta_s$ ; e.g.,  $\beta_s(\text{BA}) = \{A/R_A, B/R_A, (1, B)/R_L, (2, A)/R_L\}$ .

The distance  $d(\iota, s)$  is measured in terms of the degree of unfaithfulness of  $s$  to  $\iota$  w.r.t.  $\mathcal{R}$  (Figure 1 ⑧):

3.1.3. *Def.* Given a representational format  $\mathfrak{F}$  including some relation  $R$ , the *relational faithfulness constraints* for  $R$  are the following functions  $\mathcal{I} \times \mathcal{S} \rightarrow \mathbb{N}$  ( $|X|$  denotes the number of members of set  $X$ ):

- a.  $\mathbb{C}_R^I(\iota, s) \equiv |R(\iota) \setminus R(s)| \equiv |\{a \in R(\iota) \text{ s.t. } a \notin R(s)\}|$
- b.  $\mathbb{C}_R^O(\iota, s) \equiv |R(s) \setminus R(\iota)| \equiv |\{a \in R(s) \text{ s.t. } a \notin R(\iota)\}|$

$\mathbb{C}_R^I(\iota, s)$  is the number of elements for which  $R$  is true of the input  $\iota \in \mathcal{I}$  but not of the output  $s \in \mathcal{S}$ ;  $\mathbb{C}_R^O(\iota, s)$  is the reverse. So for the string example  $\mathcal{S}_s$  above, for  $R = R_A$ ,  $\mathbb{C}_R^I(\iota, s)$  is the number of symbols that have been deleted, and  $\mathbb{C}_R^O(\iota, s)$  is the number of symbols that have been inserted (irrespective of position in the string, which  $R_A$  ignores). So for example  $\mathbb{C}_{R_A}^I(\text{CAB}, \text{AB}) = 1$ . For  $R = R_L$ ,  $\mathbb{C}_R^I(\iota, s)$  is the number of symbols in the input that are not in the same position (relative to the left edge) in the output; e.g.,  $\mathbb{C}_{R_L}^I(\text{CAB}, \text{AB}) = 3$ , although  $\mathbb{C}_{R_L}^I(\text{CAB}, \text{CA}) = 1$ .

It is relations like these that are implicit in arguments like [23] that the mental representation for letter strings uses a format in which position is reckoned from the right as well as the left edge (i.e., that in addition to  $R_L$ , the mental format involves an analogous relation  $R_R$  in which positions are counted from the right edge). Under neurological damage, a patient spelling a list of words may erroneously intrude letters from previous words into a given word, more likely errors being those more faithful to the previous-word target. Since these errors tend to preserve the position of letters relative to right as well as left word-edges, the relevant faithfulness constraints must penalize discrepancies w.r.t.  $R_R$  as well as  $R_L$ .

An overall measure of the discrepancy between an input (target)  $\iota$  and an output  $s$ —which plays the role here of a distance metric  $d(\iota, s)$ —is given by a weighted sum of the relational faithfulness constraints for  $R$  (3.1.3). Faithfulness is one facet of well-formedness assessed by a measure called ‘Harmony’.

3.1.4. *Def.* A pair of positive *weights*  $(w_{\mathcal{R}}^I, w_{\mathcal{R}}^O)$  defines an *R-faithfulness Harmony function*  $H_{\mathcal{R}}$ :

$$H_{\mathcal{R}}(\iota, s) = -w_{\mathcal{R}}^I \mathbb{C}_{\mathcal{R}}^I(\iota, s) - w_{\mathcal{R}}^O \mathbb{C}_{\mathcal{R}}^O(\iota, s)$$

The total relational-faithfulness Harmony of format  $\mathfrak{F}$  with relation-set  $\mathcal{R}$  is  $H_{\mathcal{R}}(\iota, s) \equiv \sum_{R \in \mathcal{R}} H_{\mathcal{R}}(\iota, s)$

The faithfulness Harmony  $H_{\mathcal{R}}$  is increasingly negative as  $\iota$  and  $s$  become more disparate along the dimensions encoded in the relations  $\{R_i\} = \mathcal{R}$ ; the maximum faithfulness-Harmony value is 0, attained when  $\forall R \in \mathcal{R}, R(s) = R(\iota)$  (e.g., when  $s = \iota$ ).

### 3.2. Relational unfaithfulness as vector distance

Given a representational format  $\mathfrak{F}$  with relation-set  $\mathcal{R}$ , the dissimilarity or relational-faithfulness Harmony between an input  $\iota$  and a state  $s$ ,  $H_{\mathcal{R}}(\iota, s)$ , is directly related to the distance, in the metric vector space  $S$ , between their vectorial realizations  $\mathbf{t}$  and  $\mathbf{s}$ :  $H_{\mathcal{R}}$  is the negative of the distance squared (see ③). This result assumes the following type of tensor product realization of the representations in  $\mathcal{S}$ .

The tensor product realization of a representational format  $\mathfrak{F}$  including an arity- $n$  relation  $R$  involves a mapping  $\psi_F$  from a filler set  $\mathcal{F} = \mathbf{A}_1 \times \dots \times \mathbf{A}_n$  to a vector space  $F$ . Treating such a filler—an element such as  $(\mathbf{a}, \mathbf{b}, \mathbf{c}) \in \mathcal{A}_1 \times \mathcal{A}_2 \times \mathcal{A}_3$ —as a structure in its own right, and then recursively applying tensor product realization to  $\mathcal{F}$  (using contextual roles [15]) leads to the following type of mapping  $\psi_F: (\mathbf{a}, \mathbf{b}, \mathbf{c}) \mapsto \mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}$ .

3.2.1. *Def.* Given a representational format  $\mathfrak{F}$  (assuming all the notation of 3.1.1), a role realization mapping  $\psi_R: \{R_i\} \rightarrow R$ , and a set of filler realization mappings  $\psi_{F_{im}}: \mathcal{A}_{im} \rightarrow F_{im}$ , the *induced* tensor product realization is defined by:

$$\Psi(s) \equiv \sum_i \sum \{ \psi_{F_{i1}}(\mathbf{a}_1) \otimes \dots \otimes \psi_{F_{in_i}}(\mathbf{a}_{n_i}) \otimes \psi_R(R_i) \mid (\mathbf{a}_1, \dots, \mathbf{a}_{n_i}) \in R_i(s) \}$$

3.2.2. *Def.* A tensor product realization is called *orthogonal* iff the ranges of  $\psi_F$  and  $\psi_R$  are each orthogonal sets:  $\forall f, f' \in \mathcal{F}, f \neq f' \Rightarrow \psi_F(f) \cdot \psi_F(f') = 0$ , and similarly for  $\psi_R$ . The realization is *orthonormal* iff, in addition,  $\forall f \in \mathcal{F}, \forall r \in R, \|\psi_F(f)\| = \|\psi_R(r)\| = 1$ .

Note that orthogonal realizations are in general distributed: they need not be local (2.3). The result is:

3.2.3. *Thm.* Let  $\Psi^S$  be a tensor-product realization of a representational format  $\mathfrak{F}$  induced by orthonormal filler and role realizations  $\psi^S_F, \psi^S_R$ . Assume given faithfulness weights  $w_{\mathcal{R}}^I, w_{\mathcal{R}}^O$  for each  $R \in \mathcal{R}$ ; these define the faithfulness Harmony function  $H_{\mathcal{R}}$  (3.1.4). For each  $R$ , replace the unit-length role vector  $\psi^S_{R}(R) \equiv \hat{\mathbf{r}}^S_R$  by the rescaled vector  $\mathbf{r}_R \equiv (2w_{\mathcal{R}}^O)^{1/2} \hat{\mathbf{r}}^S_R$ . For the scaling of  $\Psi^I$  (2.2.6):

$$\psi^I_{\mathcal{R}}(R) = \rho(R) \psi^S_{\mathcal{R}}(R) \quad \text{define } \rho(R) \equiv \frac{1}{2}(w_{\mathcal{R}}^O + w_{\mathcal{R}}^I)/w_{\mathcal{R}}^O.$$

Given any  $\iota \in \mathcal{I}, s \in \mathcal{S}$ , let  $\mathbf{t} \equiv \Psi^I(\iota)$  and  $\mathbf{s} \equiv \Psi^S(s)$ . Then, up to a constant term depending on  $\iota$ , the total relational-faithfulness Harmony of  $s$  to  $\iota$  decreases as the square of the distance between  $\mathbf{s}$  and  $\mathbf{t}$ :

$$H_{\mathcal{R}}(\iota, s) = -\frac{1}{2} \|\mathbf{t} - \mathbf{s}\|^2 + \kappa(\iota)$$

where  $\kappa(\iota) \equiv \frac{1}{4} \sum_{R \in \mathcal{R}} \kappa_R(\iota)$ ;  $\kappa_R(\iota) \equiv n_R^I (w_{\mathcal{R}}^I - w_{\mathcal{R}}^O)^2 / w_{\mathcal{R}}^O$ ;  $n_R^I \equiv |R(\iota)|$ .

Note that the complexities arise only when there is an asymmetry  $w_{\mathcal{R}}^I \neq w_{\mathcal{R}}^O$ ; otherwise,  $\rho(R) = 1$ ,  $\kappa_R(\iota) = 0$ . [For a proof, see the Supplementary Materials.]

### 3.3. Relational faithfulness as network weights

We will see in 5.3.3 how relational faithfulness Harmony  $H_{\mathcal{R}}$  is naturally realized at the network level.

## 4. Recursive functions as linear transformations

Having considered the representational states of the components of the cognitive architecture  $\mathfrak{A}$ , we turn to the functions computed over these representations.

### 4.1. PC functions

Mental processes compute recursive functions; for concreteness (without compromising generality), we take the inputs and outputs of these functions to be binary trees. We now see that an interesting class of such recursive functions, denoted here  $\mathcal{F}$ , can be computed at the neural level in an extremely simple architecture: the linear associator (1.2.2).  $\mathcal{F}$  is the closure under composition of the primitive tree-manipulating functions defined in Section 2.1 (Figure 2 gives an example). It is convenient to work directly with the filler-role bindings of binary trees, assuming the canonical filler/role decomposition  $\mathfrak{T}_t$  (2.1.2).

4.1.1. *Def.* A *pseudo-tree*  $t$  is a set of binary-tree filler/role bindings with at most one filler bound to each role;  $\emptyset$  is the null pseudo-tree with no bindings;  $\mathcal{T}$  is the set of pseudo-trees. The *unification*  $t \sqcup t'$  of two pseudo-trees is the union of their bindings, provided this yields a pseudo-tree (at most one

filler bound to each role); otherwise,  $t \sqcup t' \equiv \emptyset$ . Any  $t \in \mathcal{T}$  can be represented as  $\sqcup_k \{\mathbf{A}_k/r_k\}$  for some sequence  $\{\mathbf{A}_k\} \subset \mathcal{A}$ , the alphabet of filler symbols, and  $\{r_k\} \in \mathcal{R}$ .

A function  $g: \mathcal{T} \rightarrow \mathcal{T}$  is *first order* iff  $g(\emptyset) = \emptyset$  and  $g(\sqcup_k \{\mathbf{A}_k/r_k\}) = \sqcup_k \{g(\mathbf{A}_k/r_k)\}$ .

First-order functions process each binding separately, with no interactions among bindings.

4.1.2. *Prop.* The primitive binary-tree functions  $\mathbf{ex}_0, \mathbf{ex}_1, \mathbf{ex}_\epsilon, \mathbf{cons}_0, \mathbf{cons}_1$  (2.1.3) are first order [19: 320].

As for manipulation of fillers, the basis of the set of recursive functions we consider are those in  $\mathcal{B}$ :

4.1.3. *Def.*  $\mathcal{B}$  is the set of functions  $h: \mathcal{T} \rightarrow \mathcal{T}$  satisfying the following conditions.

- a.  $h$  is first order
- b. For all  $t \in \mathcal{T}$ ,  $h(t) = \mathbf{cons}(h(\mathbf{ex}_0(t)), h(\mathbf{ex}_1(t))) \sqcup h(\mathbf{ex}_\epsilon(t))/r_\epsilon$
- c. There is a partial function  $g_{\mathcal{A}}: \mathcal{A} \rightarrow \mathcal{A}$  such that  $\forall \mathbf{X} \in \mathcal{A}$ ,  $h(\mathbf{X}/r_\epsilon) = g_{\mathcal{A}}(\mathbf{X})/r_\epsilon$ ; if  $\mathbf{X}$  is not in the domain of the partial function  $g_{\mathcal{A}}$ , then  $h(\mathbf{X}/r_\epsilon) = \emptyset$

A function  $h \in \mathcal{B}$  is indeed very simple: it simply replaces every filler  $\mathbf{X}$  in the domain of  $g_{\mathcal{A}}$  by the filler  $g_{\mathcal{A}}(\mathbf{X})$ , and deletes all other fillers. Finally, we join the filler-manipulating functions in  $\mathcal{B}$  with role manipulation in the form of the closure, under composition, of the primitive binary-tree functions.

4.1.4. *Def.* The class of *PC functions* (primitives' closure)  $\mathcal{F}$  has the following recursive definition:

- a. Base case:  $h \in \mathcal{B} \Rightarrow h \in \mathcal{F}$
- b. Recursion
  - i. If  $h \in \mathcal{B}$  and  $g \in \mathcal{F}$ , then  $h \circ g \equiv h(g) \in \mathcal{F}$
  - ii. If  $g \in \mathcal{F}$ , then  $\mathbf{ex}_0 \circ g \in \mathcal{F}$  and  $\mathbf{ex}_1 \circ g \in \mathcal{F}$
  - iii. If  $g, g' \in \mathcal{F}$ , then  $t \mapsto \mathbf{cons}(g(t), g'(t)) \in \mathcal{F}$
- c. No other functions are in  $\mathcal{F}$

While further development of the theory to address a larger class of recursive functions is in progress, the class  $\mathcal{F}$  already includes many of the functions hypothesized to be computed in cognition, for example, recursive functions like that which maps a syntactic representation such as [[Few [world leaders]] [[are admired] [by [George Bush]]]] to a semantic representation like `admire(George Bush, few world leaders)`: see the caption of Figure 2 for the definition of this function  $g$ .

## 4.2. PC functions as linear transformations

Because the primitive binary-tree functions can be realized as linear transformations (2.2.5), and because the set of linear transformations is closed under composition and addition, we get straightforwardly:

4.2.1. *Thm.* Any PC function  $g \in \mathcal{F}$  is realized by a linear transformation  $W_g$  on  $S$ ; that is,

$$\forall t, t' \in \mathcal{T}, g: t \mapsto t' \text{ if and only if } W_g \Psi(t) = \Psi(t')$$

## 4.3. PC functions as neural networks

It follows immediately that a PC function can be computed by a simple type of network, a linear associator (1.2.2); the weight matrix turns out to have a form that enables finite specification of an infinite matrix.

4.3.1. *Thm.* For any  $g \in \mathcal{F}$ , the linear transformation  $W_g$  is realized by a linear associator with weight matrix

$$\mathbf{W}_g = \mathbf{I} \otimes \underline{\mathbf{W}}_g$$

where  $\mathbf{I}$  is the identity matrix on (infinite-dimensional)  $R$  (2.2.4.d) and  $\underline{\mathbf{W}}_g$  is a finite matrix.

The map  $g \mapsto \underline{\mathbf{W}}_g$  is compositional: it can be defined constructively [19: 324] (see Figure 2 for an example).

## 5. Grammar as optimality

Having seen that a significant class of recursive functions can be computed in one massively parallel step (fully specified at the neural level in terms of multiplication and addition operations), in order to analyze linguistic cognition, we turn from recursive function theory to another general approach for specifying functions over symbols: formal language theory. This approach conceives of the function to be computed as a kind of language, and deploys rewrite-rule grammars, interpreted sequentially. These grammars can be classified in various ways, e.g., the Chomsky Hierarchy, with its corresponding hierarchy of formal sequential automata, culminating in Turing machines. It turns out that a different approach to specifying languages, more directly reducible to neural computation, deploys Harmony. In Section 3, we encountered one facet of Harmony, faithfulness; specifying languages introduces the other facet, markedness. Markedness reduces the Harmony of symbol structures that violate grammatical constraints, and may reward structures that meet grammatical desiderata. The structures of the formal language are those with

maximal Harmony.

Because of the extended chain of inference in this section, a roadmap is provided in Figure 3.

To see the basic intuition, consider a context-free rewrite-rule grammar  $\mathcal{G}$  in Chomsky Normal Form; a legal derivation  $D$  produced by such a grammar can be represented as a binary tree  $t_D$ . For example, use of the rule  $A \rightarrow B C$  in a derivation  $D$  contributes to  $t_D$  a local tree  $[_A B C]$ : a pair of sister nodes labeled  $B$  and  $C$  immediately dominated by a mother node labeled  $A$ . We can take the language generated by  $\mathcal{G}$ ,  $\mathcal{L}_{\mathcal{G}}$ , to be the set of all binary trees representing legal derivations. To evaluate whether a given tree  $t \in \mathcal{L}_{\mathcal{G}}$ , we must check that every local tree  $[_X Y Z]$  in  $t$  is sanctioned by a rule  $X \rightarrow Y Z$  in  $\mathcal{G}$ . We can do this numerically, computing a Harmony value  $H(t)$  for  $t$ , as follows. We assess a markedness-Harmony penalty of  $-3$  (lower  $H(t)$  by 3) for every symbol in  $t$ . Then for every rule  $X \rightarrow Y Z$  in  $\mathcal{G}$ , we add  $+2$  to  $H(t)$  for every pair in  $t$  consisting of a mother node labeled  $X$  with a left (right) child node labeled  $Y$  ( $Z$ ); we conceptually split this into a reward of  $+1$  for each node in the pair. Now, consider a node  $x$  labeled by some symbol  $V$ :  $x$  will lower  $H(t)$  by 3. Then if  $x$  is dominated by a legal parent node for  $V$ ,  $x$  increases  $H(t)$  by 1; if  $x$  has two legal daughter nodes for  $V$ ,  $x$  increases  $H(t)$  by 1 for each. Thus if  $x$  is fully sanctioned by the grammar (both above and below), it will contribute 0 to  $H(t)$ ; if it is not legal, its penalty  $-3$  will not be fully offset and its net contribution to  $H(t)$  will be negative. The result is that if  $t$  is legal,  $H(t) = 0$ , otherwise,  $H(t) < 0$ . The maximal-Harmony trees are those with  $H = 0$ : exactly the trees of  $\mathcal{L}_{\mathcal{G}}$ .<sup>4</sup>

### 5.1. Rewrite-rule and Optimality-Theoretic grammars as Harmonic Grammars

5.1.1. *Def.* A *Harmonic Grammar* over  $\mathcal{S}$  [25] is a function  $H_G: \mathcal{S} \rightarrow \mathbb{R}$ ;  $H_G(s)$  is the ‘(markedness) Harmony’ of  $s \in \mathcal{S}$ . The *language specified by  $H_G$* ,  $\mathcal{L}_{H_G} \subset \mathcal{S}$ , is:

$$\mathcal{L}_{H_G} \equiv \operatorname{argmax}_{s \in \mathcal{S}} H_G(s) \equiv \{s \in \mathcal{S} \mid \nexists s' \in \mathcal{S} \text{ s.t. } H_G(s') > H_G(s)\}$$

An  $s \in \mathcal{L}_{H_G}$  has maximum Harmony, and is said to be *optimal*.

A Harmonic Grammar  $H_G$  is *second-order* iff there exists a function  $H_b: B \times B \rightarrow \mathbb{R}$  (where  $B \equiv F \times R$ ) such that if  $\beta(s) = \{f_k/r_k\} \equiv \{b_k\}$ , then  $H_G(s) = \sum_j \sum_k H_b(b_j, b_k)$  ( $H_G$  depends only on *pairs* of constituents)

5.1.2. *Thm.* Given a formal language  $\mathcal{L}_{\mathcal{G}}$  generated by a grammar  $\mathcal{G}$  in the Chomsky Hierarchy, there is a second-order Harmonic Grammar  $H_G$  that generates the same language:  $\mathcal{L}_{\mathcal{G}} = \mathcal{L}_{H_G}$ .

As sketched intuitively above,  $H_G$  can be constructed compositionally from the rewrite rules of  $\mathcal{G}$  [24: 399]. This shows that formal languages can be *specified*; we consider below how languages can be *computed*.

We turn now to the grammars relevant to cognition, those of human natural languages. It turns out that for generative linguistics—the formal theory of human language—analyzing the set of grammatical expressions as a set of optimal structures is quite fruitful: thanks primarily to Optimality Theory [26, 27].

5.1.3. *Def.* An *Optimality-Theoretic system*  $\mathfrak{D} = (\mathcal{I}, \mathcal{S}, \text{Gen}, \text{Con})$  consists of

- a. a set  $\mathcal{I}$  of symbol structures called ‘inputs’ (e.g., for syntax, a meaning to be expressed)
- b. a set  $\mathcal{S}$  of symbol structures called ‘outputs’ (e.g., for syntax, a parse tree of a sentence)
- c. a function  $\text{Gen}: \mathcal{I} \rightarrow 2^{\mathcal{S}}$ ;  $\text{Gen}(t)$  is the set of ‘candidate outputs’ for  $t \in \mathcal{I}$  (e.g., all possible parses)
- d. a finite set  $\text{Con}$  of functions  $\mathbb{C}: \mathcal{I} \times \mathcal{S} \rightarrow \mathbb{N} \cup \{0\}$  called ‘constraints’;  $\mathbb{C}(t, s)$  is the number of ‘violations of  $\mathbb{C}$  by  $(t, s)$ ’, and if  $\mathbb{C}(t, s_1) < \mathbb{C}(t, s_2)$  we say  $\mathbb{C}$  ‘prefers  $s_1$  to  $s_2$  given  $t$ ’
  - i.  $\text{Con}$  includes *faithfulness* constraints; each such constraint  $\mathbb{C}_F$  evaluates a dimension of structure, being violated by each deviation of  $s$  from  $t$  w.r.t. that dimension
  - ii.  $\text{Con}$  includes *markedness* constraints; each such constraint  $\mathbb{C}_M$  evaluates, independently of  $t$ , the inherent well-formedness of  $s$  with respect to some structural dimension.

An  $\mathfrak{D}$ -grammar  $\mathbb{G}$  is a total order  $\gg$  on  $\text{Con}$ : a ‘constraint ranking’. An input-output pair  $(t, s)$ ,  $s \in \text{Gen}(t)$ , is *optimal* w.r.t. the  $\mathfrak{D}$ -grammar  $\mathbb{G} = \gg$  iff the following holds:

$\nexists o' \in \text{Gen}(t)$  s.t.  $(t, o')$  is preferred to  $(t, o)$  by the  $\gg$ -maximal constraint that prefers one of them

The language  $\mathcal{L}_{\mathbb{G}}$  specified by an  $\mathfrak{D}$ -grammar  $\mathbb{G}$  is the set of all  $\mathbb{G}$ -optimal input-output pairs.

The *typology* specified by  $\mathfrak{D}$  is the set of all languages specified by some  $\mathfrak{D}$ -grammar.

The empirical hypothesis is that the space of possible human grammars is an  $\mathfrak{D}$ -typology for some  $\mathfrak{D}$  called ‘universal grammar’—all languages possess the same grammatical constraints or preferences: differences arise only in how conflicts among those preferences are resolved; i.e., differences are limited to

<sup>4</sup> This glosses over details of the full analysis, which requires special treatment of the tree root and the start symbol, as well as terminal nodes, and conversion of  $\mathcal{G}$  to ‘Harmonic Normal Form’, which enables a second-order Harmony function (5.1.1) [24].

how the grammar of each language ranks the universal constraints.

5.1.4. *Ex.* The Italian counterpart of English *it rains* is *piove*. An Optimality Theoretic analysis goes as follows [28]. Our meaning is  $t_0 = [\text{rain (tense=present)}]$ ; we assume here a faithfulness constraint  $\mathbb{C}_F$  that requires the verb form to be *piove* in Italian, or *rains* in English. This constraint outranks all others, ensuring that only an expression with the correct verb can be optimal. A markedness constraint  $\mathbb{C}_{M1}$  requires that every sentence have a subject; the Italian candidate output sentence *piove* (or English *rains*) violates  $\mathbb{C}_{M1}$ ; *esso piove* (or *it rains*) satisfies  $\mathbb{C}_{M1}$ . Another markedness constraint  $\mathbb{C}_{M2}$  requires that every word in an expression contribute to the meaning: *esso piove* (or *it rains*), violates  $\mathbb{C}_{M2}$  (the subject is meaningless) while *piove* (or *rains*) satisfies  $\mathbb{C}_{M2}$ . In the Italian grammar,  $\mathbb{C}_{M2} \gg \mathbb{C}_{M1}$ : avoiding meaningless words is more important than providing a subject; so  $(t_0, \textit{piove})$  is optimal, hence grammatical. In English, the ranking is reversed, so  $(t_0, \textit{it rains})$  is optimal.

Despite being violated in these optimal expressions,  $\mathbb{C}_{M1}$  is active in Italian and  $\mathbb{C}_{M2}$  is active in English: in English,  $\mathbb{C}_{M2}$  can only be violated when over-ruled by a higher-ranked constraint such as  $\mathbb{C}_{M1}$ . E.g.,  $(t_0, \textit{it rains it})$  is not grammatical—it has lower Harmony than  $(t_0, \textit{it rains})$ —because of the second violation of  $\mathbb{C}_{M2}$  incurred by the meaningless object *it*: its presence (unlike the subject *it*) is not required to satisfy any higher-ranked constraint. The difference between English and Italian is not that the former prefers and the latter disprefers meaningless items: *all* languages (according to Optimality Theory) have the same grammatical preferences—differences arise only in how to resolve conflicts among those preferences, i.e., in the ranking of constraints.

Optimality Theory arguably provides the first general formal theory of linguistic typology. This has enabled formal analysis of typologies, as well as individual languages, at all linguistic levels: phonology, syntax, semantics and pragmatics [29, 30, 31, 32] (see also the archive <http://rutgers.roa.edu>). The empirical successes of Optimality Theory are cases when what is universally shared by human languages are *preferences* that outputs of the grammar should respect, as opposed to *processes* that generate those outputs. Rewrite-rules characterize grammatical knowledge as procedures, while Optimality Theory characterizes grammatical knowledge as preferences: this is the force of the idea, ‘grammars as optimization’.

Given the potential value of Optimality Theory for understanding linguistic cognition, we turn to how neural computation might be used to compute languages specified by Optimality-Theoretic grammars, rather than languages specified by serial rewrite rules or functions specified by recursive equations.

For reduction of an  $\mathcal{D}$ -grammar  $\mathbb{G}$  to the neural level, we translate  $\mathbb{G}$  into a Harmonic Grammar  $H_{\mathbb{G}}$ :

5.1.5. *Thm.* Let  $\mathbb{G} = \gg$  be an  $\mathcal{D}$ -grammar s.t. every constraint  $\mathbb{C}_k \in \textit{Con}$  is bounded above:  $\exists M \in \mathbb{N}$  s.t.  $\forall \mathbb{C}_k \in \textit{Con}, \forall t \in \mathcal{I}, \forall s \in \textit{Gen}(t), \mathbb{C}_k(t, s) < M$ . Then there is a Harmonic Grammar  $H_{\mathbb{G}}$  over  $\mathcal{I} \times \mathcal{S}$  that specifies the same language:  $\mathcal{L}_{\mathbb{G}} = \mathcal{L}_{H_{\mathbb{G}}}$ .  $H_{\mathbb{G}} = -\sum_k w_k \mathbb{C}_k$  for a set of weights  $\{w_k\}$  with  $w_k > w_j \Leftrightarrow \mathbb{C}_k \gg \mathbb{C}_j$ . (In this context,  $(t, s^*)$  is optimal for an  $H_{\mathbb{G}}$  iff  $s^* \in \textit{argmax}_s \{H_{\mathbb{G}}(t, s) \mid s \in \textit{Gen}(t)\}$ .) [26: 10.2.2, 33, 34: 463]

The basic idea of the construction is that if  $\textit{Con} \equiv \{\mathbb{C}_k\}_{k=0}^N$  is ordered by  $\gg$  so that  $\mathbb{C}_{k+1} \gg \mathbb{C}_k$ , then the constraint weights can be  $w_k \equiv M^k$  ( $M-1$  being the largest possible number of violations of any constraint). The Harmony cost incurred by a single violation of  $\mathbb{C}_k$  ( $w_k = M^k$ ) exceeds the total maximal cost that can be incurred by violations of all lower-ranked constraints ( $\sum_{j < k} w_j [M-1] = \sum_{j=0}^{k-1} M^{j+1} - \sum_{j=0}^{k-1} M^j = M^k - 1$ ). This means that  $H_{\mathbb{G}}(t, o) > H_{\mathbb{G}}(t, o')$  if and only if the highest-ranking constraint that has a preference between  $(t, o)$  and  $(t, o')$  prefers the former: optimality under  $\mathbb{G}$  and under  $H_{\mathbb{G}}$  are equivalent.

## 5.2. Harmonic Grammars as linear input transformations

The motivation for the following definition is given shortly (5.3.4).

5.2.1. *Def.* With respect to a linear intra-component input transformation  $W: S \rightarrow S$  on the vector space  $S$  realizing  $\mathcal{S}$ , the *markedness Harmony* of a vector  $\mathbf{s} \in S$  is

$$H_{M,W}(\mathbf{s}) \equiv \frac{1}{2} \mathbf{s} \cdot W\mathbf{s}$$

A linear intra-component input transformation  $W_{\mathbb{G}}: S \rightarrow S$  realizes a Harmonic Grammar  $H_{\mathbb{G}}$  iff

$$\forall \mathbf{s} \in \mathcal{S}, H_{\mathbb{G}}(\mathbf{s}) = H_{M,W_{\mathbb{G}}}(\mathbf{s}) \text{ where } \mathbf{s} \equiv \Psi(\mathbf{s})$$

5.2.2. *Thm.* A second-order Harmonic Grammar  $H_{\mathbb{G}}$  can be realized by a linear input transformation  $W_{\mathbb{G}}$ .

Because a formal language  $\mathcal{L}_{\mathbb{G}}$  specified by a rewrite-rule grammar  $\mathcal{G}$  can also be specified by a corresponding second-order Harmonic Grammar  $H_{\mathbb{G}}$  (5.1.2), we immediately get:

5.2.3. *Cor.* Given a rewrite-rule grammar  $\mathcal{G}$ , there is a linear input transformation  $W_{\mathbb{G}}$  that realizes the

Harmonic Grammar counterpart of  $\mathcal{G}$ ,  $H_{\mathcal{G}}$  [19: 333].

Moving from rewrite-rule grammars to Optimality-Theoretic grammars realized in Harmonic Grammars requires not only markedness Harmony but also faithfulness Harmony, introduced next.

### 5.3. Harmony optimization as neural computation

Given a Harmonic Grammar  $H_G$  realized as a linear intra-component input transformation  $W_G$  on  $S$  (5.2.1), and a neural realization of  $S$  in a network  $\mathcal{N}$ , the elements of  $W_G$  w.r.t. the neural basis constitute the internal weight matrix of  $\mathcal{N}$ . The utility of this is that  $\mathcal{N}$  can perform Harmony optimization, computing maximum-Harmony representations (5.3.4): these should be the optimal states constituting the language  $\mathcal{L}_{H_G}$ . (In order to achieve this, two obstacles will need to be overcome.)

5.3.1. *Def.* Given a quasi-linear network  $\mathcal{N}$  (1.3.1) with weight matrix  $\mathbf{W}$ , external input  $\mathbf{t}$  and unit activation function  $f$ , the *network Harmony*  $H_{\mathcal{N}}: I \times S \rightarrow \mathbb{R}$  is the sum of markedness Harmony  $H_{M,W}$  (5.2.1) and *faithfulness Harmony*  $H_F$ :

$$H_{\mathcal{N}}(\mathbf{t}, \mathbf{s}) = H_{M,W}(\mathbf{s}) + H_F(\mathbf{t}, \mathbf{s}); \quad H_{M,W}(\mathbf{s}) = \frac{1}{2} \mathbf{s} \cdot \mathbf{W} \mathbf{s}, \quad H_F(\mathbf{s}, \mathbf{t}) \equiv \mathbf{s} \cdot \mathbf{t} + H_1(\mathbf{s})$$

where the *unit Harmony*  $H_1$  is:

$$H_1(\mathbf{s}) \equiv \sum_{\mu} h([\mathbf{s}]_{\mu}), \quad h(a) \equiv -\int_0^a f^{-1}(x) dx$$

5.3.2. *Ex.* Let  $\mathcal{N}$  be linear (1.3.2), i.e., have units with activation function  $f(z) = z$ . Then the unit Harmony is

$$H_1(\mathbf{s}) = \sum_{\mu} h([\mathbf{s}]_{\mu}) \quad \text{where} \quad h(a) \equiv -\int_0^a f^{-1}(x) dx = -\int_0^a x dx = -\frac{1}{2} a^2 \quad \Rightarrow \quad H_1(\mathbf{s}) = \sum_{\mu} [-\frac{1}{2}([\mathbf{s}]_{\mu})^2] = -\frac{1}{2} \|\mathbf{s}\|^2$$

What is the relation between the neural-level faithfulness Harmony  $H_F$  of the network  $\mathcal{N}$  and the symbolic-level relational faithfulness  $H_R$  that evaluates the dissimilarity of symbol structures (3.1.4)?

5.3.3. *Thm.* Let  $\Psi^S$  be a tensor-product realization in a vector space  $S$  of a representational format  $\mathfrak{F}$  satisfying the conditions of 3.2.3. Suppose given, for each  $R \in \mathcal{R}$ , a pair of weights ( $w_R^I, w_R^O$ ) that define the  $R$ -faithfulness Harmony function  $H_R$ , and hence the total relational faithfulness  $H_{\mathcal{R}}$  (3.1.4). Let  $S$  be realized in a linear network  $\mathcal{N}$ , and let  $\mathbf{t} \in \mathcal{I}$ ,  $\mathbf{s} \in \mathcal{S}$ ; write  $\mathbf{t} \equiv \Psi^I(\mathbf{t}) \in I$ ,  $\mathbf{s} \equiv \Psi^S(\mathbf{s}) \in S$ . Then:

$$H_{\mathcal{R}}(\mathbf{t}, \mathbf{s}) = H_F(\mathbf{t}, \mathbf{s}) - \kappa'(\mathbf{t})$$

where

$$\kappa'(\mathbf{t}) \equiv \sum_R w_R^I n_R^I \quad (n_R^I \equiv |\mathbf{R}(\mathbf{t})|)$$

[For a proof, see the Supplementary Materials.] Thus the relational Harmony between the symbol structures ( $\mathbf{t}, \mathbf{s}$ ) can be computed as the network faithfulness Harmony between the vectors ( $\mathbf{t}, \mathbf{s}$ ) realizing ( $\mathbf{t}, \mathbf{s}$ ), up to a constant depending on  $\mathbf{t}$ ; this constant has no effect on determining the optimal representation  $\mathbf{s}$  for a given input  $\mathbf{t}$ . By 5.2.2, the well-formedness of a symbolic state  $\mathbf{s}$  as assessed by a Harmonic Grammar  $H_G$ ,  $H_G(\mathbf{s})$ , can be computed as the network markedness Harmony. The network Harmony  $H_{\mathcal{N}}$  combines both these evaluations. The computational utility of  $H_{\mathcal{N}}$  derives from 5.3.4.

5.3.4. *Thm.* Given a network  $\mathcal{N}$  with quasi-linear dynamics (1.2.1) and external input  $\mathbf{t}$ , suppose the weight matrix is symmetric ( $\forall \mu, \nu, W_{\mu\nu} = W_{\nu\mu}$ ) and scaled so that network Harmony  $H_{\mathcal{N}}$  is bounded above.<sup>5</sup> Then as activation spreads, the Harmony of the network state  $\mathbf{s}$ ,  $H_{\mathcal{N}}(\mathbf{s})$ , is non-decreasing, and  $\mathbf{s}$  converges to a *local maximum* of  $H_{\mathcal{N}}$  (a state  $\mathbf{s}$  such that  $H_{\mathcal{N}}(\mathbf{s}) \geq H_{\mathcal{N}}(\mathbf{s} + \boldsymbol{\epsilon})$ , for all small  $\boldsymbol{\epsilon}$ ) [35, 36, 37].

We are now close to showing that neural computation can compute the grammatical expressions of a symbolic rewrite-rule grammar  $\mathcal{G}$  or an Optimality-Theoretic grammar  $\mathbb{G}$ : these expressions are the maxima of a symbolic second-order Harmonic Grammar  $H_G = H_{\mathcal{G}}$  (5.1.2) or  $H_G = H_{\mathbb{G}}$  (5.1.5), and  $H_G$  can be realized in network weights  $\mathbf{W}_G$  as the Harmony  $H_{\mathcal{N}}$  of a network  $\mathcal{N}$  (5.2.2), and spreading activation in  $\mathcal{N}$  can compute the maxima of  $H_{\mathcal{N}}$  (5.3.4). But two obstacles remain.

The first is that the network computes *local* Harmony maxima, but the Harmonic Grammar demands *global* maxima (indeed, the Harmonic Grammar recognizes no such notion as ‘local maximum’). The quasi-linear dynamics drives the network state constantly uphill in Harmony, so the network ends up at the peak of whatever Harmony mountain it happens to start on: the peak is higher than all neighboring

<sup>5</sup> In order that  $H_{\mathcal{N}}$  have a maximum, as  $\mathbf{s}$  gets large, it is necessary that the  $H_1$  term get small faster than the  $H_{M,W}$  gets large (if it does). For the linear case, of interest here,  $H_1$  goes to  $-\infty$  quadratically;  $H_{M,W}$  might go to  $+\infty$  quadratically. But if we scale  $\mathbf{W}$  appropriately, by multiplying it by a sufficiently small constant, we can ensure that  $H_1$  dominates, so that  $H_{\mathcal{N}}$  has a maximum. Rescaling  $\mathbf{W}$  does not affect the location of the maxima of  $H_{M,W}$ . We can ensure that the vectors realizing all symbol structures have the same length, so adding  $H_1$  does not affect the relative Harmonies of the symbolic states.

points, but need not be the highest peak of the mountain range, which is what the Harmonic Grammar requires us to find. To compute the global Harmony maximum, the network needs some probability of moving downhill, providing a chance to pass through valleys in order to arrive at the highest mountain. Global optimization requires some randomness in activation spreading [38, 39, 40].

5.3.5. *Thm.* The neural network  $\mathcal{N}^T$  with dynamics  $\mathcal{D}_{\text{opt}}^T$  defined by the stochastic differential equation<sup>6</sup>

$$ds_{\mu} = \partial H_{\mathcal{N}}/\partial s_{\mu} dt + (2T)^{1/2} dB \quad B \text{ a Wiener process, } H_{\mathcal{N}}:S \rightarrow \mathbb{R} \text{ a Harmony function}$$

[or the corresponding stochastic difference equation

$$\Delta s_{\mu} = \partial H_{\mathcal{N}}/\partial s_{\mu} \Delta t + (2T\Delta t)^{1/2} \underline{B} \quad \underline{B} \text{ a random variable with standard normal distribution } N(0, 1)]$$

converges to the distribution  $p_T(\mathbf{s}) \propto e^{H_{\mathcal{N}}(\mathbf{s})/T}$  [41, 42]. Thus as  $T \rightarrow 0$ ,  $p_T(\mathbf{s}) \rightarrow 0$  except for globally optimal  $\mathbf{s}$ .

In  $\mathcal{N}^T$ , the *computational temperature*  $T$  is a function of computational time  $t$ : the initial temperature  $T(0)$  is high, and  $T(t)$  decreases to 0 during the course of computation. In principle, if this is done sufficiently slowly, the network's probability of being in state  $\mathbf{s}$  at time  $t$  will remain approximately proportional to  $e^{H(\mathbf{s})/T(t)}$  [43]; then at time  $t$ , the probability of a non-globally-optimal state  $\mathbf{s}'$ , with Harmony  $H'$ , relative to the probability of the globally optimal state  $\mathbf{s}^*$ , with Harmony  $H_{\text{max}}$ , is, as  $t \rightarrow \infty$  hence  $T(t) \rightarrow 0$ ,

$$p_t(\mathbf{s}')/p_t(\mathbf{s}^*) = e^{H'/T(t)}/e^{H_{\text{max}}/T(t)} = e^{-(H_{\text{max}} - H')/T(t)} \rightarrow 0$$

since  $H_{\text{max}} - H' > 0$ . Thus the probability of any non-globally-optimal state  $\mathbf{s}'$  goes to zero as  $t \rightarrow \infty$ .

## 6. Discreteness

The stochastic neural network  $\mathcal{N}^T$  (5.3.5) computes, in the limit, a state  $\mathbf{s} \in S$  with globally-maximal Harmony. The Harmonic Grammar  $H_{\mathcal{G}}$  requires a *symbolically-interpretable* state  $\mathbf{s}$  with globally-maximal Harmony: a state  $\mathbf{s}$  realizing a symbol structure  $s$  which maximizes  $H_{\mathcal{G}}(s)$  over the set of all symbol structures  $\mathcal{S}$ . Because  $H_{\mathcal{M}}$  realizes  $H_{\mathcal{G}}$  (5.2.1) we know that  $H_{\mathcal{M}}(\mathbf{s}) = H_{\mathcal{G}}(s)$  for every  $\mathbf{s} = \Psi(s)$  that realizes a symbol structure  $s \in \mathcal{S}$ ; these  $\mathbf{s}$  comprise a discrete subset  $\mathcal{S}$  of the continuous vector space  $S$ . Our second obstacle is that the global maximum of  $H_{\mathcal{N}}$  in  $S$  is generally *not* in  $\mathcal{S}$ : conflicts between the constraints encoded in  $H$  entail that optima constitute compromises that interpolate between the alternative discrete states favored by the conflicting constraints. Harmony optimization at the neural level does not yield realizations of symbolic states. To achieve that, in addition to the optimization dynamics  $\mathcal{D}_{\text{opt}}$ , another dynamics is required: *quantization*.

6.1.1. *Thm.* There is a deterministic dynamics  $\mathcal{D}_{\text{quant}}$  on  $S$ ,  $ds/dt = \mathbf{Q}(s(t))$ , which has an attractor at every vector  $\mathbf{s} \in \mathcal{S}$ , i.e., at every vector  $\mathbf{s} = \Psi(s)$  that realizes a symbol structure  $s \in \mathcal{S}$ . [44: Sec. 2.6]

6.1.2. *Def.* Given a Harmony function  $H$ , a  $\lambda$ -diffusion network  $\mathfrak{N}$  is defined by the dynamics

$$\mathcal{D}_{\lambda}(t) \equiv \lambda(t) \mathcal{D}_{\text{opt}}^{T(t)} + (1 - \lambda(t)) \mathcal{D}_{\text{quant}} \quad \text{i.e., } ds = [\lambda \nabla H(\mathbf{s}) + (1 - \lambda) \mathbf{Q}(\mathbf{s})] dt + \lambda (2T)^{1/2} d\mathbf{B}$$

where  $\lambda$  and  $T$  are functions of computational time  $t$  such that  $\lambda(t)$  and  $T(t)$  decrease to 0 as  $t \rightarrow \infty$ .

The  $\lambda$ -diffusion dynamics  $\mathcal{D}_{\lambda}$  (6.1.2) linearly combines the Harmony-optimization dynamics  $\mathcal{D}_{\text{opt}}^T$  (5.3.5), which ignores discreteness, and the quantization dynamics  $\mathcal{D}_{\text{quant}}$  (6.1.1) which ignores Harmony [45]. Although formal results have not yet been achieved, in a range of (simple) applications to modeling human language processing,  $\lambda$ -diffusion networks have proved capable of reliably computing the vectorial realizations of globally-optimal symbolic states, when the speed at which  $T(t)$  and  $\lambda(t)$  go to zero is sufficiently slow. When too fast, errors are made; the outputs are symbolic states, but not necessarily globally optimal ones. In fact, in a number of cases, the probability of outputting a symbolic state  $s$  is approximately proportional to  $e^{H_{\mathcal{M}}(s)/T}$  (for some  $T$ ): the relative Harmonies of error states  $s$  govern their relative probabilities. In these applications, the distribution of errors captures the key properties of human performance. For example, because faithfulness Harmony formalizes relational similarity of an error to a correct response (5.3.3), the probability of an error type decreases appropriately as the similarity between the error and the correct response decreases; and because markedness and faithfulness Harmony together formalize grammatical knowledge (5.2.2, 5.1.5), the probability of ungrammatical errors is appropriately small.

## 7. Conclusion

A detailed summary was provided in Section 0.2; only a few concluding remarks are given here.

<sup>6</sup> The derivative  $\partial H_{\mathcal{N}}/\partial s_{\mu}$  is simply  $[\mathbf{W}\mathbf{s} + \mathbf{t} - \mathbf{f}^{-1}(s)]_{\mu}$  which is just  $[\mathbf{W}\mathbf{s} + \mathbf{t} - \mathbf{s}]_{\mu}$  in the linear case.

This work seeks not to go beyond classical computability, but rather to explore what restrictions may be placed on cognitive functions assuming that they are computed in accord with a certain conception of neural computation. In this conception, cognition can be characterized by functions over meaningful symbols, but not as computation over such symbols. Vector spaces (of neural network states) are seen to provide a powerful representational medium for the computation of symbolic cognitive functions, even restricting to linear processing. In such vector spaces, a cognitively significant class of recursive functions can be computed in a single massively parallel step. A natural characterization of the functions computed by certain neural networks is through optimization, and as a result neural computation has led to powerful insights into one of the deepest realms of symbolic cognitive science, the theory of universal grammar. Formal languages can be specified in such neural computational terms, as well as natural languages; and while effective computability has not been proved, simulations suggest that these models simultaneously capture central features of both the idealized computational capacity of human linguistic competence and the errorful real-time computation of human linguistic performance.

## **Acknowledgments**

I warmly thank my collaborators in this work, Alan Prince, Géraldine Legendre, Matthew Goldrick, Donald Mathis, Bruce Tesar, John Hale, and Yoshiro Miyata; special thanks to Don and Matt for comments on an earlier draft of this paper (I am solely responsible for any errors, of course). For helpful discussion of the work I thank Colin Wilson, James McClelland, Robert Frank, Robert Cummins and most of all the late David E. Rumelhart. I gratefully acknowledge partial financial support for this work from the National Science Foundation, Johns Hopkins University, les Chaires internationales de recherche Blaise Pascal, la Laboratoire de Sciences Cognitives et Psycholinguistique/CNRS, and the Robert J. Glushko and Pamela Samuelson Foundation. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation or other sponsors.

## References

- [1] Turing A. M. 1936 On computable numbers, with an application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*.42:230-65.
- [2] Pylyshyn Z. W. 1984 *Computation and cognition: Toward a foundation for cognitive science*. Cambridge, MA: MIT Press.
- [3] Smolensky P. 1988 On the proper treatment of connectionism. *Behavioral and Brain Sciences*.11:1-74.
- [4] Copeland B. J., Proudfoot D. 1996 On Alan Turing's Anticipation of Connectionism. *Synthese*.108:361-77.
- [5] Hofstadter D. R. 1985 Waking up from the Boolean dream, or, subcognition as computation. In: *Metamagical themas*, pp. 631-65: Basic Books.
- [6] Cummins R., Schwarz G. 1991 Connectionism, computation, and cognition. In: *Connectionism and the Philosophy of Mind* (Horgan T. E., Tienson J., eds), pp. 60-73. Dordrecht, Holland: Kluwer.
- [7] Smolensky P. 2006 Computational levels and integrated connectionist/symbolic explanation. In: *The Harmonic Mind: From Neural Computation to Optimality-theoretic Grammar, Vol. 2*, pp. 503-92. Cambridge, MA: MIT Press.
- [8] McClelland J. L., Botvinick M. M., Noelle D. C., Plaut D. C., Rogers T. T., Seidenberg M. S., et al. 2010 Letting structure emerge: connectionist and dynamical systems approaches to cognition. *Trends in Cognitive Sciences*.14(8):348-56.
- [9] Smolensky P., Legendre G. 2006 *The Harmonic Mind: From Neural Computation to Optimality-Theoretic Grammar. Vol 1: Cognitive architecture. Vol 2: Linguistic and Philosophical Implications*. Cambridge, MA: MIT Press.
- [10] Grossberg S. 1982 *Studies of Mind and Brain: Neural Principles of Learning, Perception, Development, Cognition, and Motor Control*. Boston, MA: Reidel.
- [11] Kohonen T. 1977 *Associative memory: A system-theoretical approach*. New York: Springer.
- [12] Smolensky P. 1986 Neural and conceptual interpretations of parallel distributed processing models. In: *Parallel distributed processing: Explorations in the microstructure of cognition, Vol. 2* (McClelland J. L., Rumelhart D. E., the PDP Research Group, eds), pp. 390-431. Cambridge, MA: MIT Press.
- [13] Smolensky P. 1996 Dynamical perspectives on neural networks. In: *Mathematical Perspectives on Neural Networks* (Smolensky P., Mozer M. C., Rumelhart D. E., eds), pp. 245-70. Mahwah, NJ: Erlbaum.
- [14] Fodor J. A. 1975 *The Language of Thought*. Cambridge, MA: Harvard University Press.
- [15] Smolensky P. 1990 Tensor product variable binding and the representation of symbolic structures in connectionist networks. *Artificial Intelligence*.46:159-216.
- [16] Abelson H., Sussman G. J., Sussman J. 1985 *Structure and Interpretation of Computer Programs*. Cambridge, MA: MIT Press.
- [17] Smolensky P. 2006 Formalizing the Principles I: Representation and Processing in the Mind/Brain. In: *The Harmonic Mind: From Neural Computation to Optimality-Theoretic Grammar, Vol. 1*, pp. 147-205. Cambridge, MA: MIT Press.
- [18] Messiah A. 1961 *Quantum mechanics*. Elsevier Science/Dover.
- [19] Smolensky P. 2006 Tensor product representations: Formal foundations. In: *The Harmonic Mind: From Neural Computation to Optimality-Theoretic Grammar, Vol. 1*, pp. 271-344. Cambridge, MA: MIT Press.
- [20] Hinton G. E., McClelland J. L., Rumelhart D. E. 1986 Distributed representation. In: *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1* (Rumelhart D. E., McClelland J. L., the PDP Research Group, eds), pp. 77-109. Cambridge, MA: MIT Press.
- [21] Churchland P. S., Sejnowski T. J. 1992 *The computational brain*. Cambridge, MA: MIT Press.
- [22] Abbott L., Sejnowski T. J., eds. 1999 *Neural Codes and distributed representations*. Cambridge, MA: MIT Press.
- [23] Fischer-Baum S. Position representation: General principles or domain-specificity? [Doctoral dissertation]. Baltimore, MD: Johns Hopkins University; 2010.
- [24] Hale J., Smolensky P. 2006 Harmonic grammars and harmonic parsers for formal languages. In: *The Harmonic Mind: From Neural Computation to Optimality-Theoretic Grammar, Vol. 1*, pp. 393-415. Cambridge, MA: MIT Press.

- [25] Legendre G., Miyata Y., Smolensky P. Harmonic Grammar—A formal multi-level connectionist theory of linguistic well-formedness: Theoretical foundations. Proceedings of the Cognitive Science Society; 1990 July; Cambridge, MA: Erlbaum; 1990. p. 388–95.
- [26] Prince A., Smolensky P. 1993/2004 *Optimality Theory: Constraint Interaction in Generative Grammar*: Technical report, Rutgers University and University of Colorado at Boulder, 1993. ROA 537, 2002. Revised version published by Blackwell, 2004.
- [27] Prince A., Smolensky P. 1997 Optimality: From neural networks to universal grammar. *Science*.275:1604–10.
- [28] Grimshaw J., Samek-Lodovici V. 1998 Optimal subjects and subject universals. In: *Is the Best Good Enough? Optimality and Competition in Syntax* (Barbosa P., Fox D., Hagstrom P., McGinnis M., Pesetsky D., eds), pp. 193–219. Cambridge, MA: MIT Press.
- [29] Kager R. 1999 *Optimality Theory*. Cambridge: Cambridge University Press.
- [30] McCarthy J. J., ed. 2004 *Optimality Theory in Phonology: A Reader*. Malden, MA: Blackwell.
- [31] Legendre G., Grimshaw J., Vikner S., eds. 2001 *Optimality-Theoretic Syntax*. Cambridge, MA: MIT Press.
- [32] Blutner R., De Hoop H., Hendriks P. 2006 *Optimal communication*. Stanford, CA: CSLI Publications.
- [33] Prince A. 2002 Anything goes. In: *A New Century of Phonology and Phonological Theory* (Honma T., Okazaki M., Tabata T., Tanaka S., eds), pp. 66–90. Tokyo: Kaitakusha.
- [34] Soderstrom M., Mathis D. W., Smolensky P. 2006 Abstract genomic encoding of Universal Grammar in Optimality Theory. In: *The Harmonic Mind: From Neural Computation to Optimality-theoretic Grammar, Vol. 2*, pp. 403–71. Cambridge, MA: MIT Press.
- [35] Cohen M. A., Grossberg S. 1983 Absolute stability of global pattern formation and parallel memory storage by competitive neural networks. *IEEE Transactions on Systems, Man, and Cybernetics*.13:815–25.
- [36] Hopfield J. J. 1984 Neurons with graded response have collective computational properties like those of two-state neurons. *Proceedings of the National Academy of Sciences USA*.81:3088–92.
- [37] Golden R. M. 1986 The "Brain-State-in-a-Box" neural model is a gradient descent algorithm. *Mathematical Psychology*.30–31:73–80.
- [38] Kirkpatrick S., Gelatt C. D., Jr., Vecchi M. P. 1983 Optimization by simulated annealing. *Science*.220:671–80.
- [39] Ackley D. H., Hinton G. E., Sejnowski T. J. 1985 A learning algorithm for Boltzmann machines. *Cognitive Science*.9:147–69.
- [40] Smolensky P. 1986 Information processing in dynamical systems: Foundations of Harmony Theory. In: *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1* (Rumelhart D. E., McClelland J. L., the PDP Research Group, eds), pp. 194–281. Cambridge, MA: MIT Press.
- [41] Movellan J. R. 1998 A learning theorem for networks at detailed stochastic equilibrium. *Neural Computation*.10:1157–78.
- [42] Movellan J. R., McClelland J. L. 1993 Learning continuous probability distributions with symmetric diffusion networks. *Cognitive Science*.17:463–96.
- [43] Geman S., Geman D. 1984 Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.6:721–41.
- [44] Baird B., Eeckman F. 1993 A normal form projection algorithm for associative memory. In: *Associative neural memories* (Hassoun M. H., ed., pp. 135–66). New York, NY: Oxford University Press.
- [45] Smolensky P., Goldrick M., Mathis D. W. in press Optimization and quantization in Gradient Symbol Systems: A framework for integrating the continuous and the discrete in cognition. *Cognitive Science*.

## Figure Captions

### Figure 1 of 3

A schematic depiction of the proposed theory for neural computation of symbolic cognitive functions.

### Figure 2 of 2

Applying the technique of 4.2.1, this linear associator network computes the function

$$g(s) = \text{cons}(\text{ex}_1(\text{ex}_0(\text{ex}_1(s))), \text{cons}(\text{ex}_1(\text{ex}_1(\text{ex}_1(s))), \text{ex}_0(s)))$$

with the corresponding weight matrix

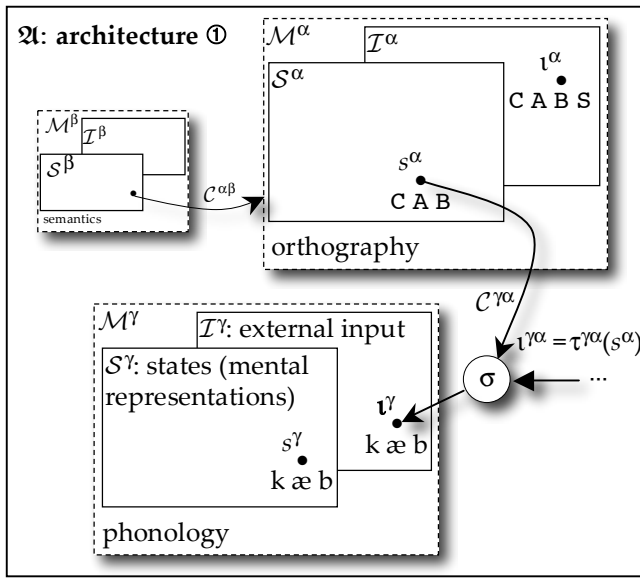
$$\mathbf{W}_g = \mathbf{W}_{\text{cons0}}[\mathbf{W}_{\text{ex1}}\mathbf{W}_{\text{ex0}}\mathbf{W}_{\text{ex1}}] + \mathbf{W}_{\text{cons1}}[\mathbf{W}_{\text{cons0}}(\mathbf{W}_{\text{ex1}}\mathbf{W}_{\text{ex1}}\mathbf{W}_{\text{ex1}}) + \mathbf{W}_{\text{cons1}}(\mathbf{W}_{\text{ex0}})].$$

Disk area displays magnitude of activations and weights; black/white denotes positive/negative. The unbounded competence of the network results from the unbounded weight matrix  $\mathbf{W}_g$ , which is finitely-specified as the tensor product of a finite matrix  $\underline{\mathbf{W}}_g$  and the unbounded identity matrix  $\mathbf{I}$  (4.3). (Reprinted with permission from [9].)

### Figure 3 of 3

A roadmap showing the chain of inference in Section 5.

Symbolic Level



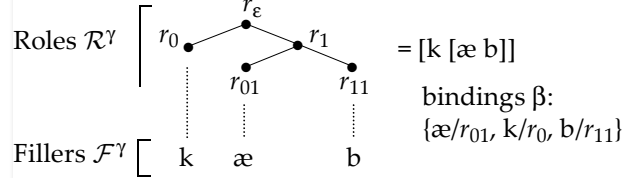
**R: Representational format ②**

$\mathcal{R}^\alpha = \{R_A, R_L\}$   
 $R_A(\mathbf{t}^\alpha) = \{\mathcal{A}, \mathcal{C}, \mathcal{S}, \mathcal{B}\}$   
 $R_A(s^\alpha) = \{\mathcal{A}, \mathcal{C}, \mathcal{B}\}$   
 $R_L(\mathbf{t}^\alpha) = \{(1, \mathcal{C}), (2, \mathcal{A}), (3, \mathcal{B}), (4, \mathcal{S})\}$   
 $R_L(s^\alpha) = \{(1, \mathcal{C}), (2, \mathcal{A}), (3, \mathcal{B})\}$

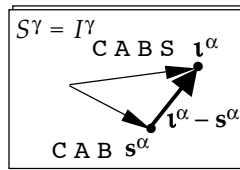
**Representational faithfulness ③**

$C_{R_A}^I(\mathbf{t}^\alpha, s^\alpha) = 1$  [S deleted]  
 $C_{R_A}^O(\mathbf{t}^\alpha, s^\alpha) = 0$  [∅ inserted]  
 $C_{R_L}^I(\mathbf{t}^\alpha, s^\alpha) = 1$   
 $C_{R_L}^O(\mathbf{t}^\alpha, s^\alpha) = 0$

**S: Filler-role decomposition of  $\mathcal{S}^\gamma$  ⑥**

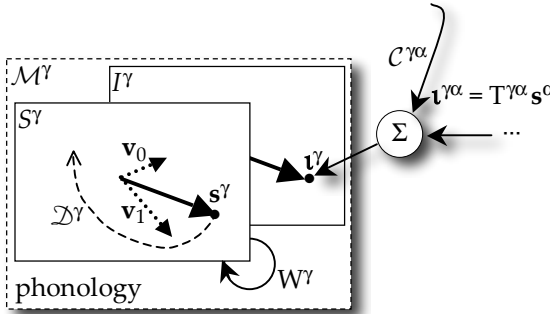


Vectorial Level



**Symbolic dissimilarity as vectorial distance ④**

representational faithfulness of  $s_\alpha$  to  $\mathbf{t}^\alpha$   
 $= -\frac{1}{2} \text{distance}(\mathbf{s}^\alpha, \mathbf{t}^\alpha)^2$   
 $= -\frac{1}{2} \text{length}(\mathbf{i}^\alpha - \mathbf{s}^\alpha)^2$



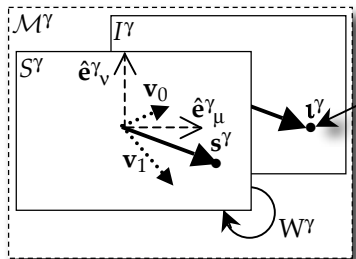
**D: Vectorial dynamics ②**

$ds^\gamma/dt = f^\gamma[\mathbf{i}^\gamma(t)] - s^\gamma(t); \quad \mathbf{i}^\gamma(t) \equiv W^\gamma s^\gamma(t) + \mathbf{t}^\gamma(t)$

**Ψ: Tensor product (vectorial) realization of S ⑥**

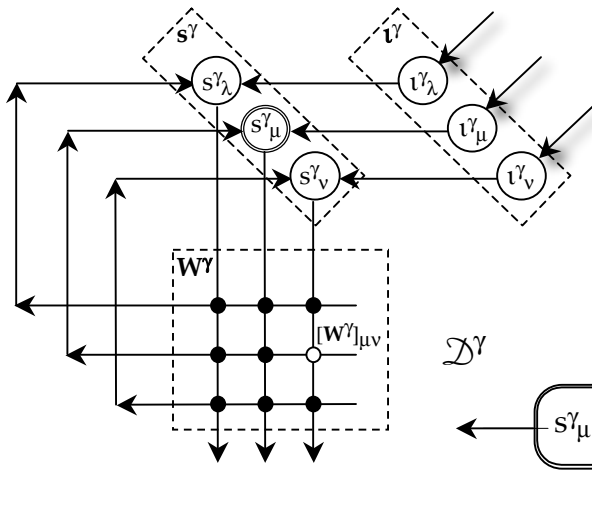
Role vectors:  $R = \{\mathbf{r}_0, \mathbf{r}_1\}$ ; Filler vectors:  $F = \{\mathbf{f}_\text{æ}, \mathbf{f}_\text{k}, \dots\}$   
 $s^\gamma = \Psi(s^\gamma \Psi([k \ \text{æ} \ b])) = \mathbf{v}_0 + \mathbf{v}_1$   
 $\mathbf{v}_0 = \mathbf{f}_\text{k} \otimes \mathbf{r}_0; \quad \mathbf{v}_1 = (\mathbf{f}_\text{æ} \otimes \mathbf{r}_0 + \mathbf{f}_\text{b} \otimes \mathbf{r}_1) \otimes \mathbf{r}_1 = \mathbf{f}_\text{æ} \otimes \mathbf{r}_{01} + \mathbf{f}_\text{b} \otimes \mathbf{r}_{11}$

Neural Level



**Representations in neural coordinates ③**

relative to the neural basis  $\{\dots, \hat{\mathbf{e}}_{\mu}^\gamma, \dots, \hat{\mathbf{e}}_{\nu}^\gamma, \dots\}$ :  
 $s^\gamma = (\dots, 1.3, \dots, -0.5, \dots)$   
 $\mathbf{v}_0 = (\dots, 0.6, \dots, 0.3, \dots) = \mathbf{f}_\text{k} \otimes \mathbf{r}_0$   
 $\mathbf{v}_1 = (\dots, 0.7, \dots, -0.8, \dots)$   
 $\mathbf{f}_\text{k} = (1, -0.3, 2, \dots, 0.6, \dots), \quad \mathbf{r}_0 = (1, -1) \Rightarrow$   
 $\mathbf{f}_\text{k} \otimes \mathbf{r}_0 = (1, -0.3, 2, \dots, \underline{0.6}, \dots; -1, \underline{0.3}, -2, \dots, -0.6, \dots)^T$



**D: Neural dynamics ④**

$ds^\gamma_\mu/dt = f^\gamma([\mathbf{i}^\gamma(t)]_\mu) - s^\gamma_\mu(t)$   
 $[\mathbf{i}^\gamma(t)]_\mu = \Sigma_{\nu} [W^\gamma]_{\mu\nu} s^\gamma_\nu(t) + [\mathbf{t}^\gamma(t)]_\mu$   
 $= [\text{internal input}]_\mu + [\text{external input}]_\mu$   
 At equilibrium:  $s^\gamma_\mu = f^\gamma([\mathbf{i}^\gamma]_\mu)$

Figure 1

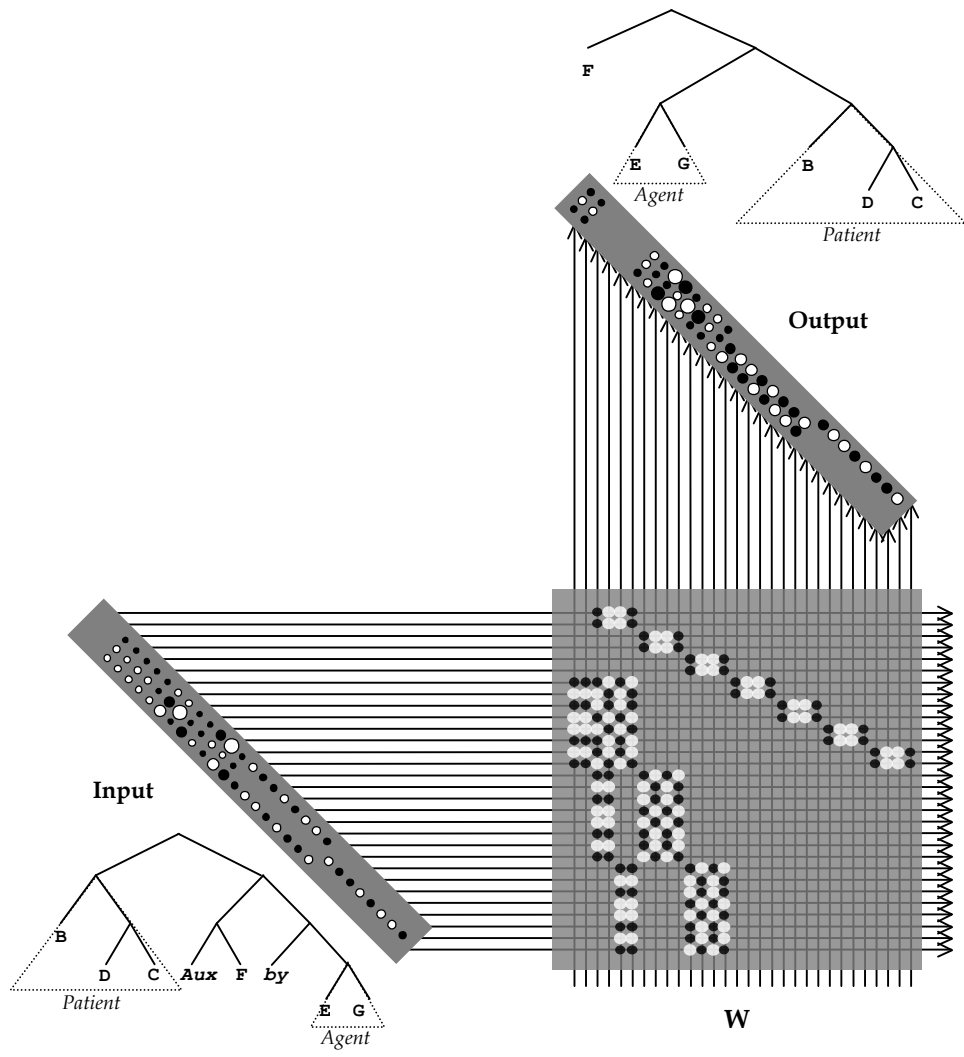


Figure 2

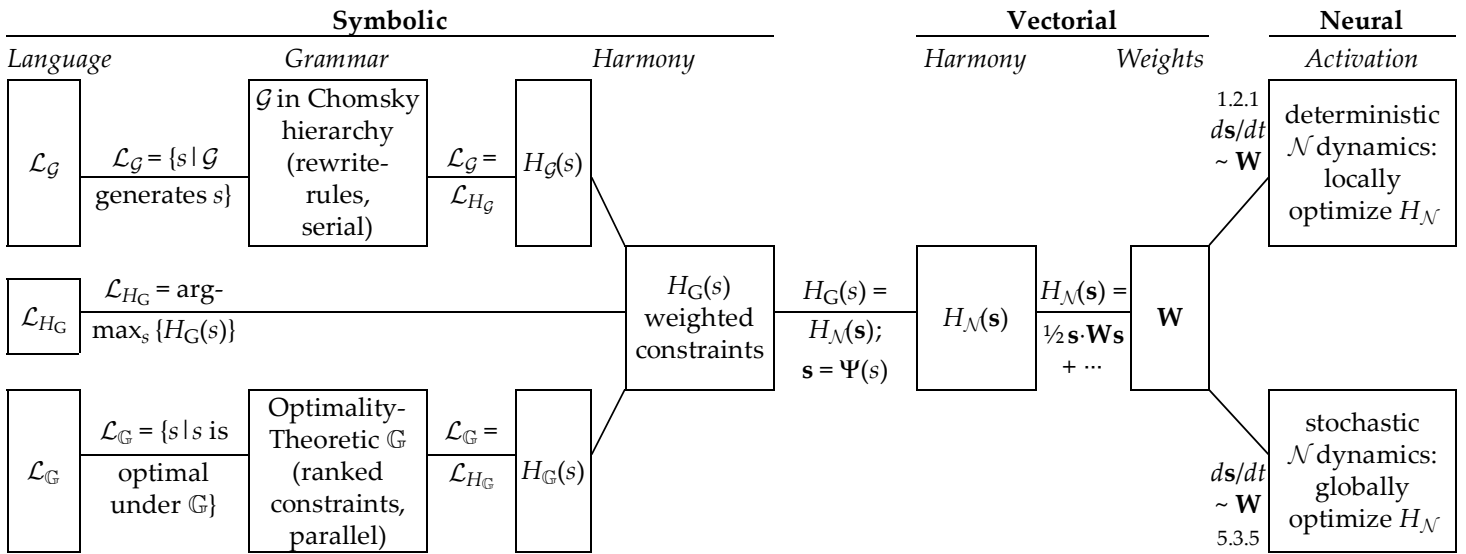


Figure 3

### Supplementary material: Proofs of Theorems 3.2.3 and 5.3.3

#### Theorem 3.2.3

From the text:

2.2.6 For all  $f \in \mathcal{F}$ ,  $\Psi^I(f) = \Psi^S_F(f)$ ; there exists a function  $\rho: \mathcal{R} \rightarrow \mathbb{R}$  such that for all  $r \in \mathcal{R}$ :  $\Psi^I_{\mathbf{R}}(r) = \rho(r)\Psi^S_{\mathbf{R}}(r)$ .

3.1.3 *Def.* Given a representational format including some relation  $\mathbf{R}$ , the *relational faithfulness constraints* for  $\mathbf{R}$  are the following functions  $\mathcal{I} \times \mathcal{S} \rightarrow \mathbb{N}$  ( $|X|$  denotes the number of members of set  $X$ ):

$$\mathbb{C}_{\mathbf{R}}^I(\mathfrak{t}, s) \equiv |\mathbf{R}(\mathfrak{t}) \setminus \mathbf{R}(s)| \equiv |\{a \in \mathbf{R}(\mathfrak{t}) \text{ s.t. } a \notin \mathbf{R}(s)\}|$$

$$\mathbb{C}_{\mathbf{R}}^O(\mathfrak{t}, s) \equiv |\mathbf{R}(s) \setminus \mathbf{R}(\mathfrak{t})| \equiv |\{a \in \mathbf{R}(s) \text{ s.t. } a \notin \mathbf{R}(\mathfrak{t})\}|$$

3.1.4 *Def.* A pair of positive *weights*  $(w_{\mathbf{R}}^I, w_{\mathbf{R}}^O)$  defines an *R-faithfulness Harmony function*  $H_{\mathbf{R}}$ :

$$H_{\mathbf{R}}(\mathfrak{t}, s) = -w_{\mathbf{R}}^I \mathbb{C}_{\mathbf{R}}^I(\mathfrak{t}, s) - w_{\mathbf{R}}^O \mathbb{C}_{\mathbf{R}}^O(\mathfrak{t}, s)$$

The total relational-faithfulness Harmony is  $H_{\mathcal{R}}(\mathfrak{t}, s) \equiv \sum_{\mathbf{R} \in \mathcal{R}} H_{\mathbf{R}}(\mathfrak{t}, s)$

3.2.3 *Thm.* Let  $\Psi^S$  be a tensor-product realization of a representational format  $\mathfrak{F}$  induced by orthonormal filler and role realizations  $\Psi^S_F, \Psi^S_R$ . Assume given faithfulness weights  $w_{\mathbf{R}}^I, w_{\mathbf{R}}^O$  for each  $\mathbf{R} \in \mathcal{R}$ ; these define the faithfulness Harmony function  $H_{\mathbf{R}}$  (3.1.4). For each  $\mathbf{R}$ , replace the unit-length role vector  $\Psi^S_{\mathbf{R}}(\mathbf{R}) \equiv \hat{\mathbf{r}}^S_{\mathbf{R}}$  by the rescaled vector  $\mathbf{r}^S_{\mathbf{R}} \equiv (2w_{\mathbf{R}}^O)^{1/2} \hat{\mathbf{r}}^S_{\mathbf{R}}$ . For the scaling of  $\Psi^I$  (2.2.6):

$$\Psi^I_{\mathbf{R}}(\mathbf{R}) = \rho(\mathbf{R})\Psi^S_{\mathbf{R}}(\mathbf{R}) \quad \text{define } \rho(\mathbf{R}) \equiv 1/2(w_{\mathbf{R}}^O + w_{\mathbf{R}}^I)/w_{\mathbf{R}}^O.$$

Given any  $\mathfrak{t} \in \mathcal{I}, s \in \mathcal{S}$ , let  $\mathfrak{t} \equiv \Psi^I(\mathfrak{t})$  and  $\mathfrak{s} \equiv \Psi^S(s)$ . Then, up to a constant term depending on  $\mathfrak{t}$ , the total relational-faithfulness Harmony of  $s$  to  $\mathfrak{t}$  decreases as the square of the distance between  $\mathfrak{s}$  and  $\mathfrak{t}$ :

$$H_{\mathcal{R}}(\mathfrak{t}, s) = -1/2 \|\mathfrak{t} - \mathfrak{s}\|^2 + \kappa(\mathfrak{t})$$

where  $\kappa(\mathfrak{t}) \equiv 1/4 \sum_{\mathbf{R} \in \mathcal{R}} \kappa_{\mathbf{R}}(\mathfrak{t})$ ;  $\kappa_{\mathbf{R}}(\mathfrak{t}) \equiv n_{\mathbf{R}}^I (w_{\mathbf{R}}^I - w_{\mathbf{R}}^O)^2 / w_{\mathbf{R}}^O$ ;  $n_{\mathbf{R}}^I \equiv |\mathbf{R}(\mathfrak{t})|$ .

*Proof of 3.2.3.*

Assume given a pair  $(\mathfrak{t}, s)$ . By 3.1.3, letting  $n_{\mathbf{R}}^I \equiv |\mathbf{R}(\mathfrak{t})|$ ,  $n_{\mathbf{R}}^S \equiv |\mathbf{R}(s)|$  we have

$$c_{\mathbf{R}}^I \equiv \mathbb{C}_{\mathbf{R}}^I(\mathfrak{t}, s) \equiv |\mathbf{R}(\mathfrak{t}) \setminus \mathbf{R}(s)|; \quad c_{\mathbf{R}}^S \equiv \mathbb{C}_{\mathbf{R}}^O(\mathfrak{t}, s) \equiv |\mathbf{R}(s) \setminus \mathbf{R}(\mathfrak{t})|; \quad g_{\mathbf{R}} \equiv |\mathbf{R}(\mathfrak{t}) \cap \mathbf{R}(s)| \Rightarrow$$

$$n_{\mathbf{R}}^I = c_{\mathbf{R}}^I + g_{\mathbf{R}}; \quad n_{\mathbf{R}}^S = c_{\mathbf{R}}^S + g_{\mathbf{R}} \Rightarrow$$

$$(1) \quad n_{\mathbf{R}}^I - c_{\mathbf{R}}^I = g_{\mathbf{R}} = n_{\mathbf{R}}^S - c_{\mathbf{R}}^S$$

Let the (now rescaled) role vectors be

$$\mathbf{r}^S_{\mathbf{R}} = (2w_{\mathbf{R}}^O)^{1/2} \hat{\mathbf{r}}^S_{\mathbf{R}}$$

$$\mathbf{r}^I_{\mathbf{R}} \equiv \rho(\mathbf{R}) \mathbf{r}^S_{\mathbf{R}} = \rho(\mathbf{R}) (2w_{\mathbf{R}}^O)^{1/2} \hat{\mathbf{r}}^S_{\mathbf{R}}.$$

Now  $\Psi^S(\mathfrak{t}) \equiv \mathfrak{t} = \sum_{\mathbf{R}} \mathfrak{t}_{\mathbf{R}}$  with

$$\mathfrak{t}_{\mathbf{R}} \equiv \sum_k n_{\mathbf{R}}^I k \mathbf{f}_k \otimes \mathbf{r}^I_{\mathbf{R}}$$

where  $n_{\mathbf{R}}^I k$  is 1 if  $f_k \in \mathbf{R}(\mathfrak{t})$  and 0 otherwise, and  $\mathbf{f}_k \equiv \Psi^I_F(f_k) = \Psi^S_F(f_k)$  by 2.2.6. The sum ranges over all possible  $f_k$ . Note that

$$(2) \quad (n_{\mathbf{R}}^I k)^2 = n_{\mathbf{R}}^I k$$

Analogously,  $\mathfrak{s} = \sum_{\mathbf{R}} \mathfrak{s}_{\mathbf{R}}$  where

$$\mathfrak{s}_{\mathbf{R}} \equiv \sum_l n_{\mathbf{R}}^S l \mathbf{f}_l \otimes \mathbf{r}^S_{\mathbf{R}}$$

where  $n_{\mathbf{R}}^S l$  is 1 if  $f_l \in \mathbf{R}(s)$  and 0 otherwise; the sum is again over all possible  $f_l$ .

Now

$$\begin{aligned} \|\mathfrak{t}\|^2 &= \mathfrak{t} \cdot \mathfrak{t} = [\sum_{\mathbf{R}} \sum_k n_{\mathbf{R}}^I k \mathbf{f}_k \otimes \mathbf{r}^I_{\mathbf{R}}] \cdot [\sum_{\mathbf{R}'} \sum_l n_{\mathbf{R}'}^I l \mathbf{f}_l \otimes \mathbf{r}^I_{\mathbf{R}'}] = \sum_{\mathbf{R}} \sum_k \sum_{\mathbf{R}'} \sum_l n_{\mathbf{R}}^I k n_{\mathbf{R}'}^I l (\mathbf{f}_k \otimes \mathbf{r}^I_{\mathbf{R}} \cdot \mathbf{f}_l \otimes \mathbf{r}^I_{\mathbf{R}'}) \\ &= \sum_{\mathbf{R}} \sum_k \sum_{\mathbf{R}'} \sum_l n_{\mathbf{R}}^I k n_{\mathbf{R}'}^I l (\mathbf{f}_k \cdot \mathbf{f}_l) (\mathbf{r}^I_{\mathbf{R}} \cdot \mathbf{r}^I_{\mathbf{R}'}) \end{aligned}$$

Since the filler and role realizations are orthogonal, this collapses to

$$\|\mathfrak{t}\|^2 = \sum_{\mathbf{R}} \sum_k (n_{\mathbf{R}}^I k)^2 \|\mathbf{f}_k\|^2 \|\mathbf{r}^I_{\mathbf{R}}\|^2$$

The filler vectors have not been rescaled and so they have retained their original normalization:  $\|\mathbf{f}_k\| = 1$ .

Using (2) we have:

$$\|\mathbf{t}\|^2 = \sum_{\mathbf{R}} \sum_k n_{\mathbf{R}}^I k \|\mathbf{r}^I_{\mathbf{R}}\|^2$$

The role vectors  $\mathbf{r}^I_{\mathbf{R}}$  have been rescaled from the unit vectors  $\hat{\mathbf{f}}^S_{\mathbf{R}}$  by  $\rho(\mathbf{R}) \equiv \frac{1}{2}(w_{\mathbf{R}}^{\mathbf{O}} + w_{\mathbf{R}}^{\mathbf{I}})/w_{\mathbf{R}}^{\mathbf{O}}$  so now

$$\|\mathbf{r}^I_{\mathbf{R}}\|^2 = \|\rho(\mathbf{R})(2w_{\mathbf{R}}^{\mathbf{O}})^{1/2} \hat{\mathbf{f}}^S_{\mathbf{R}}\|^2 = \rho(\mathbf{R})^2 (2w_{\mathbf{R}}^{\mathbf{O}})$$

Since  $\sum_k n_{\mathbf{R}}^I k = |\mathbf{R}(\mathbf{t})| \equiv n_{\mathbf{R}}^{\mathbf{I}}$  we get

$$(3) \quad \|\mathbf{t}\|^2 = \sum_{\mathbf{R}} \sum_k n_{\mathbf{R}}^I k [2w_{\mathbf{R}}^{\mathbf{O}} \rho(\mathbf{R})^2] = \sum_{\mathbf{R}} [2w_{\mathbf{R}}^{\mathbf{O}} \rho(\mathbf{R})^2] \sum_k n_{\mathbf{R}}^I k = \sum_{\mathbf{R}} 2w_{\mathbf{R}}^{\mathbf{O}} \rho(\mathbf{R})^2 n_{\mathbf{R}}^{\mathbf{I}}$$

Analogously, using the rescaled role vectors  $\mathbf{r}^S_{\mathbf{R}}$  we have

$$\|\mathbf{s}\|^2 = \sum_{\mathbf{R}} \sum_k (n_{\mathbf{R}}^S k)^2 \cdot \|\mathbf{f}_k\|^2 \|\mathbf{r}^S_{\mathbf{R}}\|^2 = \sum_{\mathbf{R}} n_{\mathbf{R}}^S [(2w_{\mathbf{R}}^{\mathbf{O}})^{1/2}]^2 = \sum_{\mathbf{R}} 2w_{\mathbf{R}}^{\mathbf{O}} n_{\mathbf{R}}^S$$

Then, by the orthogonality of the filler and (rescaled) role vectors, we have

$$\begin{aligned} \mathbf{t} \cdot \mathbf{s} &= \sum_{\mathbf{R}} \mathbf{t}_{\mathbf{R}} \cdot \sum_{\mathbf{R}'} \mathbf{s}_{\mathbf{R}'} = \sum_{\mathbf{R}} \sum_k n_{\mathbf{R}}^I k \mathbf{f}_k \otimes \mathbf{r}^I_{\mathbf{R}} \cdot \sum_{\mathbf{R}'} \sum_l n_{\mathbf{R}'}^S l \mathbf{f}_l \otimes \mathbf{r}^S_{\mathbf{R}'} = \sum_{\mathbf{R}} \sum_k \sum_{\mathbf{R}'} \sum_l n_{\mathbf{R}}^I k n_{\mathbf{R}'}^S l (\mathbf{f}_k \cdot \mathbf{f}_l) (\mathbf{r}^I_{\mathbf{R}} \cdot \mathbf{r}^S_{\mathbf{R}'}) \\ &= \sum_{\mathbf{R}} \sum_k n_{\mathbf{R}}^I k n_{\mathbf{R}}^S k (\mathbf{r}^I_{\mathbf{R}} \cdot \mathbf{r}^S_{\mathbf{R}}) = \sum_{\mathbf{R}} \sum_k n_{\mathbf{R}}^I k n_{\mathbf{R}}^S k [\rho(\mathbf{R})(2w_{\mathbf{R}}^{\mathbf{O}})^{1/2} \hat{\mathbf{f}}^S_{\mathbf{R}}] \cdot [(2w_{\mathbf{R}}^{\mathbf{O}})^{1/2} \hat{\mathbf{f}}^S_{\mathbf{R}}] \\ &= \sum_{\mathbf{R}} \sum_k n_{\mathbf{R}}^I k n_{\mathbf{R}}^S k [2w_{\mathbf{R}}^{\mathbf{O}} \rho(\mathbf{R})] \end{aligned}$$

Since  $n_{\mathbf{R}}^I n_{\mathbf{R}}^S$  is 1 if  $\mathbf{f}_k \in \mathbf{R}(\mathbf{t}) \cap \mathbf{R}(\mathbf{s})$  and 0 otherwise, we have

$$\sum_k n_{\mathbf{R}}^I k n_{\mathbf{R}}^S k = |\mathbf{R}(\mathbf{t}) \cap \mathbf{R}(\mathbf{s})| \equiv g_{\mathbf{R}}$$

hence

$$\mathbf{t} \cdot \mathbf{s} = \sum_{\mathbf{R}} g_{\mathbf{R}} [2w_{\mathbf{R}}^{\mathbf{O}} \rho(\mathbf{R})]$$

Finally, then:

$$\begin{aligned} \|\mathbf{t} - \mathbf{s}\|^2 &= (\mathbf{t} - \mathbf{s}) \cdot (\mathbf{t} - \mathbf{s}) = \|\mathbf{t}\|^2 + \|\mathbf{s}\|^2 - 2\mathbf{t} \cdot \mathbf{s} \\ &= [\sum_{\mathbf{R}} 2w_{\mathbf{R}}^{\mathbf{O}} \rho(\mathbf{R})^2 n_{\mathbf{R}}^{\mathbf{I}}] + [\sum_{\mathbf{R}} 2w_{\mathbf{R}}^{\mathbf{O}} n_{\mathbf{R}}^{\mathbf{S}}] - 2[\sum_{\mathbf{R}} 2w_{\mathbf{R}}^{\mathbf{O}} \rho(\mathbf{R}) g_{\mathbf{R}}] \\ &= 2\sum_{\mathbf{R}} w_{\mathbf{R}}^{\mathbf{O}} [\rho(\mathbf{R})^2 n_{\mathbf{R}}^{\mathbf{I}} + n_{\mathbf{R}}^{\mathbf{S}} - 2\rho(\mathbf{R}) g_{\mathbf{R}}] \\ &\equiv \sum_{\mathbf{R}} D_{\mathbf{R}} \end{aligned}$$

where we now evaluate a single term  $D_{\mathbf{R}}$ , omitting the  $\mathbf{R}$ , and using (1):

$$\begin{aligned} D &\equiv 2w^{\mathbf{O}} [\rho^2 n^{\mathbf{I}} + n^{\mathbf{S}} - 2\rho g] \\ &= 2w^{\mathbf{O}} [\rho^2 n^{\mathbf{I}} + (n^{\mathbf{I}} - c^{\mathbf{I}} + c^{\mathbf{S}}) - 2\rho(n^{\mathbf{I}} - c^{\mathbf{I}})] \\ &= 2w^{\mathbf{O}} [(\rho^2 - 2\rho + 1)n^{\mathbf{I}} - c^{\mathbf{I}} + c^{\mathbf{S}} + 2\rho c^{\mathbf{I}}] \\ &= 2w^{\mathbf{O}} [(\rho - 1)^2 n^{\mathbf{I}} + (2\rho - 1)c^{\mathbf{I}} + c^{\mathbf{S}}] \end{aligned}$$

Now

$$2\rho - 1 \equiv 2[\frac{1}{2}(w^{\mathbf{O}} + w^{\mathbf{I}})/w^{\mathbf{O}}] - 1 = w^{\mathbf{I}}/w^{\mathbf{O}}$$

so

$$\frac{1}{2}D = w^{\mathbf{O}} [(\rho - 1)^2 n^{\mathbf{I}} + (w^{\mathbf{I}}/w^{\mathbf{O}})c^{\mathbf{I}} + c^{\mathbf{S}}] = w^{\mathbf{O}} (\rho - 1)^2 n^{\mathbf{I}} + w^{\mathbf{I}} c^{\mathbf{I}} + w^{\mathbf{O}} c^{\mathbf{S}}$$

Thus, using 3.1.4,

$$\begin{aligned} -\frac{1}{2}\|\mathbf{t} - \mathbf{s}\|^2 &= -\frac{1}{2}\sum_{\mathbf{R}} D_{\mathbf{R}} = -\sum_{\mathbf{R}} [w_{\mathbf{R}}^{\mathbf{I}} c_{\mathbf{R}}^{\mathbf{I}} + w_{\mathbf{R}}^{\mathbf{O}} c_{\mathbf{R}}^{\mathbf{S}}] - \sum_{\mathbf{R}} w_{\mathbf{R}}^{\mathbf{O}} (\rho(\mathbf{R}) - 1)^2 n_{\mathbf{R}}^{\mathbf{I}} \\ &= H_{\mathcal{R}}(\mathbf{t}, \mathbf{s}) - \kappa(\mathbf{t}) \end{aligned}$$

where

$$\begin{aligned} \kappa(\mathbf{t}) &\equiv \sum_{\mathbf{R}} w_{\mathbf{R}}^{\mathbf{O}} (\rho(\mathbf{R}) - 1)^2 n_{\mathbf{R}}^{\mathbf{I}} = \sum_{\mathbf{R}} w_{\mathbf{R}}^{\mathbf{O}} [\frac{1}{2}(w_{\mathbf{R}}^{\mathbf{O}} + w_{\mathbf{R}}^{\mathbf{I}})/w_{\mathbf{R}}^{\mathbf{O}} - 1]^2 n_{\mathbf{R}}^{\mathbf{I}} \\ &= \sum_{\mathbf{R}} w_{\mathbf{R}}^{\mathbf{O}} [(w_{\mathbf{R}}^{\mathbf{O}} + w_{\mathbf{R}}^{\mathbf{I}} - 2w_{\mathbf{R}}^{\mathbf{O}})/2w_{\mathbf{R}}^{\mathbf{O}}]^2 n_{\mathbf{R}}^{\mathbf{I}} \\ &= \frac{1}{4} \sum_{\mathbf{R}} n_{\mathbf{R}}^{\mathbf{I}} (w_{\mathbf{R}}^{\mathbf{I}} - w_{\mathbf{R}}^{\mathbf{O}})^2 / w_{\mathbf{R}}^{\mathbf{O}} \\ &\equiv \frac{1}{4} \sum_{\mathbf{R}} \kappa_{\mathbf{R}}(\mathbf{t}) \end{aligned}$$

□

### Theorem 5.3.3

From the text:

5.3.1 *Def.* Given a quasi-linear network  $\mathcal{N}$  with weight matrix  $\mathbf{W}$ , external input  $\mathbf{t}$  and unit activation function  $f$ , network Harmony  $H_{\mathcal{N}}: I \times S \rightarrow \mathbb{R}$  is the sum of markedness Harmony  $H_{M,W}$  (5.2.1) and faithfulness Harmony  $H_F$ :

$$H_{\mathcal{N}}(\mathbf{t}, \mathbf{s}) = H_M(\mathbf{s}) + H_F(\mathbf{t}, \mathbf{s}); \quad H_M(\mathbf{s}) = \frac{1}{2} \mathbf{s} \cdot \mathbf{W} \mathbf{s}, \quad H_F(\mathbf{s}, \mathbf{t}) \equiv \mathbf{s} \cdot \mathbf{t} + H_1(\mathbf{s})$$

where the unit Harmony  $H_1$  is:

$$H_1(\mathbf{s}) \equiv \sum_{\mu} h([\mathbf{s}]_{\mu}), \quad h(a) \equiv -\int_0^a f^{-1}(x) dx$$

5.3.2 *Ex.* Let  $\mathcal{N}$  be linear (1.3.2), i.e., have units with activation function  $f(i) = i$ . Then the unit Harmony is

$$H_1(\mathbf{s}) = \sum_{\mu} h([\mathbf{s}]_{\mu}) \quad \text{where} \quad h(a) \equiv -\int_0^a f^{-1}(x) dx = -\int_0^a x dx = -\frac{1}{2} a^2 \quad \Rightarrow \quad H_1(\mathbf{s}) = -\frac{1}{2} \|\mathbf{s}\|^2$$

5.3.3 *Thm.* Let  $\Psi^S$  be a tensor-product realization in a vector space  $S$  of a representational format  $\mathfrak{F}$  satisfying the conditions of 3.2.3. Suppose given, for each  $R \in \mathcal{R}$ , a pair of weights  $(w_R^I, w_R^O)$  that define the  $R$ -faithfulness Harmony function  $H_R$  and total relational faithfulness  $H_{\mathcal{R}}$  (3.1.4). Let  $S$  be realized in a linear network  $\mathcal{N}$ , and let  $\mathbf{t} \in \mathcal{I}$ ,  $\mathbf{s} \in \mathcal{S}$ ; write  $\mathbf{t} \equiv \Psi^I(\mathbf{t}) \in I$ ,  $\mathbf{s} \equiv \Psi^S(\mathbf{s}) \in S$ . Then:

$$H_{\mathcal{R}}(\mathbf{t}, \mathbf{s}) = H_F(\mathbf{t}, \mathbf{s}) - \kappa'(\mathbf{t})$$

where

$$\kappa'(\mathbf{t}) \equiv \sum_{R} w_R^I n_R^I \quad (n_R^I \equiv |\mathbf{R}(\mathbf{t})|)$$

*Proof of 5.3.3.*

By Thm. 3.2.3,

$$\begin{aligned} H_{\mathcal{R}}(\mathbf{t}, \mathbf{s}) &= -\frac{1}{2} \|\mathbf{t} - \mathbf{s}\|^2 + \kappa(\mathbf{t}) \\ &= -\frac{1}{2} \|\mathbf{t}\|^2 - \frac{1}{2} \|\mathbf{s}\|^2 + \mathbf{s} \cdot \mathbf{t} + \kappa(\mathbf{t}) \end{aligned}$$

Since for a linear network

$$H_F(\mathbf{s}, \mathbf{t}) \equiv \mathbf{s} \cdot \mathbf{t} + H_1(\mathbf{s}) = \mathbf{s} \cdot \mathbf{t} - \frac{1}{2} \|\mathbf{s}\|^2$$

this gives

$$\begin{aligned} H_{\mathcal{R}}(\mathbf{t}, \mathbf{s}) &= -\frac{1}{2} \|\mathbf{t}\|^2 - \frac{1}{2} \|\mathbf{s}\|^2 + [H_F(\mathbf{s}, \mathbf{t}) + \frac{1}{2} \|\mathbf{s}\|^2] + \kappa(\mathbf{t}) \\ &= -\frac{1}{2} \|\mathbf{t}\|^2 + H_F(\mathbf{s}, \mathbf{t}) + \kappa(\mathbf{t}) \\ &= H_F(\mathbf{t}, \mathbf{s}) - \kappa'(\mathbf{t}) \end{aligned}$$

where

$$\kappa'(\mathbf{t}) \equiv \frac{1}{2} \|\mathbf{t}\|^2 - \kappa(\mathbf{t})$$

By (3):

$$\begin{aligned} \|\mathbf{t}\|^2 &= \sum_{R} 2w_R^O \rho(\mathbf{R})^2 n_R^I \\ &= \sum_{R} 2w_R^O [\frac{1}{2}(w_R^O + w_R^I)/w_R^O]^2 n_R^I \end{aligned}$$

so

$$\begin{aligned} \kappa'(\mathbf{t}) &\equiv \sum_{R} w_R^O [\frac{1}{2}(w_R^O + w_R^I)/w_R^O]^2 n_R^I - \frac{1}{4} \sum_{R} n_R^I (w_R^I - w_R^O)^2 / w_R^O \\ &= \sum_{R} \kappa'_R(\mathbf{t}) n_R^I \end{aligned}$$

where

$$\begin{aligned} \kappa'_R(\mathbf{t}) &\equiv w_R^O [\frac{1}{2}(w_R^O + w_R^I)/w_R^O]^2 - \frac{1}{4} (w_R^I - w_R^O)^2 / w_R^O \\ &= \frac{1}{4} [w_R^O]^{-1} [(w_R^O + w_R^I)^2 - (w_R^I - w_R^O)^2] \\ &= \frac{1}{4} [w_R^O]^{-1} [(w_R^O)^2 + (w_R^I)^2 + 2w_R^O w_R^I - ((w_R^O)^2 + (w_R^I)^2 - 2w_R^O w_R^I)] \\ &= \frac{1}{4} [w_R^O]^{-1} [4(w_R^O w_R^I)] \\ &= w_R^I \end{aligned}$$

Hence

$$\kappa'(\mathbf{t}) = \sum_{R} w_R^I(\mathbf{t}) n_R^I$$

□