

Substructure in Complex Systems and Partially Subdivided Neural Networks I: Stability of Composite Patterns

R. Sadr-Lahijany and Y. Bar-Yam

ECS Engineering and Physics Departments

44 Cummington St., Boston University, Boston MA 02215

In order to gain insight into the role of substructure in complex systems we investigate attractor neural network models with substructure. Partially subdivided neural networks enable the storage of highly correlated states, and may also be useful as a model for brain function. Two possible biological examples are the separation of color, shape and motion in visual processing, and the use of parts of speech in language. With this motivation we analyze the ability of partially subdivided attractor networks to store both imprinted neural states and composite states formed out of combinations of imprinted subnetwork states. As the strength of intersubnetwork synapses are reduced more combinations of imprinted substates are recalled by the network. We perform a mean field analysis of the stability of composite patterns. Analytic solution of the equations at zero temperature show how stability of a particular composite pattern is controlled by the number of subdivisions that represent each imprinted pattern. Numerical minimization of the free energy is used to obtain the phase transition diagrams for networks with 2, 3 or 4 subdivisions.

I. Introduction

Many complex systems have the property that they are formed from substructures that extend to essentially the same scale as the whole system. The brain is segregated into hemispheres and lobes and further divided into functional regions. The human body has physiological systems further divided into organs. Proteins are often organized out of several amino acid chains with substructure formed of α -helices and β -sheets. Life on earth, considered as a complex system, is divided among climates, ecosystems, habitats and species. Global weather patterns are formed out of large scale air and ocean currents, storms and regions of high and low pressure. In all of these systems the largest scale of subdivision comprises fewer than 100 parts, and more typically of order 10 parts of the whole system. For biological systems, particularly biological organisms, the explanation of the substructure must originate from the advantages that accrue to the system from its presence. Our ultimate objective is to understand the role of functional and structural subdivision in complex systems.

In this manuscript we begin a systematic investigation of the properties of partially subdivided neural networks. The advantages of subdivided networks may be considered in the

context of pattern recognition tasks, the storage of correlated information and as a mechanism for generalization from a small number of training examples. Subdivided networks also are a model for the functional structure of the brain. Our construction of partially subdivided networks consists of a conventional network of N neurons with Hebbian learning where the strength of synapses between q subdivisions of N/q neurons are reduced by a factor g compared to the synapses between neurons within each subnetwork. We expect systematic dilution of inter-subnetwork synapses, with g the fraction of remaining synapses, to lead to similar results.

This paper is organized as follows. Section II provides a brief review of related research on neural networks. Section III describes qualitatively the properties and possible advantages of subdivided networks. Section IV is the main body of this article where the mean field equations for the subdivided networks are derived. They are solved analytically for a particular class of patterns, called ideal composite patterns, at $T=0$. The composite patterns are generalized and followed at $T>0$ using numerical solutions of the mean field equations to obtain phase diagrams for their stability. Section V presents brief conclusions and relates this work to other complex systems and their substructure.

II. Related Research on Subdivided Neural Networks

Theoretical analysis of subdivided attractor networks has focused on hierarchical networks where the interaction between subnetworks is only through the subnetwork magnetization.¹⁻³ Higher levels of the hierarchy are networks of synapses between neuron-like subnetwork magnetizations. In effect, the inversion degeneracy of subnetwork states is used for the higher level degrees of freedom. If the magnetization is explicitly normalized to create new Ising spin variables the states of different subnetworks and different levels of the hierarchy are decoupled.³ When unnormalized, the interactions between subnetwork states are through the magnitude of the magnetization.²

In a more biologically relevant model, the possibility of training correlations between subnetwork states was considered by Sutton et al⁴ in the context of a network hierarchy with asymmetric synapses between selected neurons in different subnetworks. Our analysis is close in spirit to that of Sutton et al, however our interest is not only in describing the trained correlation between subnetwork states but also the stability of states that are formed out of other combinations of trained subnetwork states. Idiart and Theumann⁵ considered binary branching hierarchical networks where all individual spins in subnetworks interact via Hebbian synapses. They considered the retrieval of network states that differ from the imprinted states only by inversion of subnetwork states. Our study of composite states in subdivided networks generalizes their discussions.

III. Advantages of subdivision

The primary measure of neural network capability is its capacity to store imprinted patterns. The storage capacity of a neural network increases with the degree of interconnectedness. For a network where each neuron is connected to every other neuron the number of imprints that can be recalled αN is proportional to the number of neurons N with a constant of proportionality α somewhat dependent on the particular imprinting rule. When additional imprints are added an overload catastrophe causes erasure of all information. Subdivision *inherently* results in a decreased storage capacity. Biological evidence indicates that the brain has well defined functional subdivisions. We are investigating the functional advantages such an architecture might provide.

A. The left-right universe

Consider first an artificial world composed of pictures with independent (uncorrelated) left and right halves. A completely connected network is capable of recalling αN pictures. However, if we divide the network into left and right hemispheres, the subdivided network can recall $(N/2)^2$ pictures. Since the number of neurons is large this results in a huge increase in effective storage capacity. Moreover, training the network may be achieved with only $(N/2)$ imprints. Since each hemisphere acts independently all possible combinations of left and right halves are recalled. Of course, if the network were divided top-from-bottom rather than left-from-right the scheme would not work. The memory would be degraded to $(N/2)$ patterns and many spurious patterns would be introduced. The effectiveness of subdivision requires matching to the nature of information. The key property that motivates subdivision is the independence - the lack of correlation - between parts of the information. The existence of synaptic connections reflects a correlation or coupling between distinct pieces of information. The use of subdivision enables *a-priori* separation between independent aspects .

B. Artificial neural network applications

Subdivided neural networks have been used in artificial neural network applications for performing parallel or sequential subtasks. A recent illustrative example for parallel tasks makes use of a partitioned feed forward network for the recognition of Kanji.⁶ For this application the Kanji were separated into components (radicals) by a preprocessing step and independent recognition tasks were trained and performed on the separated radicals through distinct feed forward networks. The motivation for subdividing the network as a means of subtask performance is intuitively clear since, as a first approximation, the large number of Kanji may be considered to be formed by combining together comparatively few radicals. It is the independence of the distinct recognition tasks that makes this effective.

In the event that the tasks are not completely independent the introduction of weak interactions between the subnetworks should introduce correlations between the subnetwork functions .

C. Color shape and motion in vision

The human visual system does not take advantage of the two hemispheres of the brain to divide the visual information right from left because the left and right visual fields are not independent. Instead, visual processing separates three attributes of the information: color, shape and motion.⁷ The implication of this preprocessing step is that these information categories are partially independent so that, for example, visual fields with the different shapes can have the same colors. Or, vice versa the same shapes can have different colors. This independence is genetically coded into the structure of the initial information processing.

The existence of three attribute categories enables a large number of descriptive categories to be constructed out of a selection of one from each attribute. For example, by separating the color information to one subnetwork, shape information to the second, and movement information to the third, it is possible for the network to identify categories such as: RED ROUND MOVING-LEFT, and RED ROUND MOVING-UP, BLUE SQUARE MOVING-LEFT, and BLUE ROUND MOVING-UP. The network receiving color information identifies the color, and so on. In a fully connected network these categories would each require separate identification (and category correlations would severely impair operation). If the subdivided network were completely separated the total number of categories would be a product of the number of categories stored in each subnetwork.

Partial subdivision implies correlation between the different attributes is also significant. In the natural world shape, color and motion are not completely independent attributes. Local correlation in the visual field such as the coincidence of edges in color and shape maps is only part of the correlation between these attributes. At higher levels of abstraction / processing there are important correlations between the overall shape of an object its color and both its direction and likelihood of motion.

D. Parts of speech - nouns, verbs and adjectives

If subdivision provides advantages in neural networks it should be particularly relevant to man-made constructs such as language. The subdivided network provides a systematic method for information organization in terms of elements (the stable states of subnetworks) which are organized in element-categories (the stable states of a particular subnetwork) and the compatibility relationships between elements as dictated by the inter-subnetwork synapses. This is reminiscent of the structure of grammar where nouns, verbs and adjectives and other parts of speech are categories that have elements and there are compatibility relations among them. It is

tempting to speculate that different subdivisions of the brain are responsible at least for the major parts of speech and the ability to combine them in different ways results from weakening the strength of inter-subnetwork synapses compared to the intra-subnetwork synapses that store representations of each word. Unlike a dictionary where the grammatical usage of a word is identified by a label (noun, verb, adjective, etc.) the storage of a word in a particular subdivision identifies its grammatical usage.

In Fig. 1 an example illustrating a fully connected and fully subdivided network is shown. The complete network is large enough to be subdivided into three networks each of which can store three words (coded appropriately). A fully-connected network would then be able to store nine sentences with three words each since the storage capacity grows linearly with size. On the subdivided network we could imprint three sentences and twenty-seven sentences would be recognized. The central difference between the set of sentences that can be remembered by the full network and the subdivided network is summarized by the concept of 'semantic content' vs. 'grammar.' The complete network knows more full sentences but does not have knowledge of the divisibility of the sentences into parts that can be put together in different ways. The subdivided network knows the parts but has no relationship between them, thus it knows grammar but does not know any context information, like who it is that fell.

The actual process in the human brain is a combination of the two, where sentences make sense or are 'grammatically correct' if properly put together out of largely interchangeable parts, but an actual event or recalled incident is a specific combination. This can be captured in the network by having a partial interconnection between subnetworks. It is to be expected that partial subdivision leads to an intermediate situation where sentences are constructed out of largely interchangeable parts, but an actual event or recalled incident is a specific combination. This is confirmed by the analysis described below.

IV. Derivation of stability conditions for subdivided networks

We derive and solve the mean field equations for the retrieval of composite states in a partially subdivided network. Section A reformulates the problem in a separable form suitable for application of the standard mean field solution for fully connected networks. Section B obtains the mean field equations using a generalized form of self averaging. Section C obtains the analytic solution of the mean field equations for the composite patterns at $T=0$, where the details of the derivation are given in the Appendix. In Section D we discuss the conventional spurious patterns that are distinct from the composite patterns but also appear in suitably generalized form as solutions of the mean field equations. In Section E the stability of the composite patterns at $T=0$ is proven using the second derivatives of the free energy. The composite solutions are

generalized to $T > 0$ in Section F using numerical solutions that obtain phase diagrams for networks with 2,3, or 4 subdivisions.

A. Formulating the subdivided network in a separable form.

We assume a network comprised of q subnetworks (each containing $N = N / q$ neurons) that are fully internally connected but more weakly connected to each other, the ratio of connection strengths is controlled by a parameter $g \in [0,1]$. $g=0$ corresponds to a completely subdivided network and $g=1$ corresponds to a conventional Hopfield network. For arbitrary g the synaptic connection matrix is written as:

$$J_{ij} = \begin{cases} J_{ij} & \text{Integer_part} \left(\frac{i}{N} \right) = \text{Integer_part} \left(\frac{j}{N} \right) \\ gJ_{ij} & \text{otherwise} \end{cases} \quad (1)$$

The first case corresponds to i and j in the same block along the matrix diagonal, i.e. in the same subnetwork. J is the usual Hebbian matrix:

$$J_{ij} = \begin{cases} \frac{1}{N} \sum_{\mu=1}^p \mu_i^\mu \mu_j^\mu & i \neq j \\ 0 & i = j \end{cases} \quad (2)$$

The neural firing patterns $\{\mu_i^\mu = \pm 1\}$ are chosen at random, where $i \in \{1, \dots, N\}, \mu \in \{1, \dots, p\}$.

We can reformulate this in a more convenient way. First we write:

$$J = J^0 + J^1 \quad (3)$$

$$J_{ij}^0 = gJ = \begin{cases} \frac{g}{N} \sum_{\mu=1}^p \mu_i^\mu \mu_j^\mu & i \neq j \\ 0 & i = j \end{cases} \quad (4)$$

$$J_{ij}^1 = \begin{cases} \frac{1-g}{N} \sum_{s=1}^q \sum_{\mu=1}^p \mu_i^{\mu s} \mu_j^{\mu s} & i \neq j \\ 0 & i = j \end{cases} \quad (5)$$

where we have introduced a set of new "pattern"s $\mu^{\mu s}$ that are correlated with the original set.

$\mu^{\mu s}$ is identical to the pattern μ^μ in the s th subdivision and zero elsewhere:

$$\mu_i^{\mu s} = \mu_i^\mu \delta_{s, \text{Integer_part} \left(\frac{i}{N} \right)} \quad (6)$$

where

$$i \in S = \begin{cases} 1 & \text{if } i \in S \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

$$S = \{(-1)^N + 1, \dots, N\} \quad (8)$$

takes values from 1 to q throughout the following treatment.

We augment the set of correlated patterns by including the original patterns as

$$\begin{matrix} \mu^0 \\ i \end{matrix} \quad \begin{matrix} \mu \\ i \end{matrix} \quad (9)$$

whenever the case $\mu = 0$ is included we use the index μ for μ to distinguish it from μ

Eqs. (6)-(8) can be generalized to include this case by introducing $S_0 = \{1, \dots, N\}$.

Using these definitions the synaptic matrix can be written in a separable form similar in appearance to the Hebbian form:

$$J_{ij} = \frac{1}{N} \sum_{\mu=0}^q \begin{matrix} \mu \\ i \end{matrix} \begin{matrix} \mu \\ j \end{matrix} \quad \begin{matrix} i \neq j \\ i = j \end{matrix} \quad (10)$$

where μ is proportional to μ :

$$\begin{matrix} \mu^0 \\ i \end{matrix} = \sqrt{g} \begin{matrix} \mu^0 \\ i \end{matrix} \quad \begin{matrix} \mu \\ i \end{matrix} = \sqrt{g} \begin{matrix} \mu \\ i \end{matrix} \quad (11)$$

$$\begin{matrix} \mu \\ i \end{matrix} = \sqrt{1-g} \begin{matrix} \mu \\ i \end{matrix} .$$

B. Mean field equations: Order parameters and free energy

The energy in the magnetic analog for the neural network has the standard form:

$$E\{s\} = - \sum_i h_i^{ext} s_i - \frac{1}{2} \sum_{ij} J_{ij} s_i s_j \quad (12)$$

because of the separable form of J_{ij} in Eq. (10), the conventional mean field analysis⁸ may be used without modification. This directly leads to the ensemble averaged mean field equations for the order parameters:

$$M^\mu = \frac{1}{N} \sum_{i=1}^N \begin{matrix} \mu \\ i \end{matrix} \tanh \left(\sum_{\mu=0}^q M^\mu \begin{matrix} \mu \\ i \end{matrix} + h^{ext} \right) \quad (13)$$

where the order parameters are ensemble average overlaps of the neural state s_i with the patterns μ defined in the usual way,

$$M^\mu = \frac{1}{N} \sum_{i=1}^N \begin{matrix} \mu \\ i \end{matrix} \langle s_i \rangle \quad (14)$$

Now we formulate a generalization of self-averaging for the case of our correlated patterns. It applies to any function of these correlated patterns and reads:

$$\prod_{i=1}^N f \left\{ \left. \begin{matrix} \mu \\ i \end{matrix} \right\} \right\}_{\mu} = N^q \left\langle \left\langle f \left\{ \left. \begin{matrix} \mu^0 = \sqrt{g} \mu, \mu = \sqrt{1-g} \mu \end{matrix} \right\} \right\}_{\mu} \right\rangle \right\rangle \quad (15)$$

where the double brackets indicate averaging over the quenched random variables. Using this self averaging and rescaling M^{μ}

$$M^{\mu^0} = \frac{\sqrt{g}}{q} m^{\mu^0} \quad (16)$$

$$M^{\mu} = \frac{\sqrt{1-g}}{q} m^{\mu}$$

in Eq. (13) leads to the mean field equations for m^{μ}

$$m^{\mu^0} = \frac{q}{N} \left\langle \left\langle \mu \tanh \frac{1}{q} \left(g m^{\mu^0} + (1-g) m^{\mu} \right) + h^{ext} \right\rangle \right\rangle \quad (17)$$

$$m^{\mu} = \left\langle \left\langle \mu \tanh \frac{1}{q} \left(g m^{\mu^0} + (1-g) m^{\mu} \right) + h^{ext} \right\rangle \right\rangle \quad (18)$$

note the sum rule:

$$m^{\mu^0} = \sum_{\mu=1}^q m^{\mu} \quad (19)$$

Given a solution m^{μ} of the equations, symmetries guarantee the following transformations yield new solutions:

$$m^{\mu} \rightarrow -m^{\mu} \quad (20)$$

for any fixed μ but for all μ simultaneously;

$$m^{\mu} \rightarrow m \quad (21)$$

for fixed μ and m but for all μ simultaneously; and

$$m^{\mu} \rightarrow m^{\mu} \quad (22)$$

for all μ simultaneously and any fixed m . The first two are conventional symmetries and the third is a renaming of subdivisions resulting in the switching of subdivision order parameters.

The free energy may also be obtained in terms of the order parameters. Substituting Eq. (10) into Eq. (12) and setting $h^{ext} = 0$ the energy takes the form:

$$E = -\frac{1}{2N} \sum_{ij, i \neq j} \sum_{\mu} s_i^{\mu} s_j^{\mu} \quad (23)$$

then using conventional manipulations⁸ that define the free energy, introduce the auxiliary fields M^μ (or the rescaled m^μ), perform saddle point integrations and self-averaging, this leads to the free energy

$$\begin{aligned}
F = & \frac{g}{2q^2} \sum_{\mu=1}^p (m^{\mu 0})^2 + \frac{1-g}{2q^2} \sum_{\mu=1}^q (m^\mu)^2 \\
& - \frac{1}{q} \sum_{\mu=1}^q \left\langle \left\langle \ln 2 \cosh \frac{1}{q} \sum_{\mu} \left[(g(m^{\mu 0} + h^{\mu 0}) + (1-g)(m^\mu + h^\mu)) \right] \right\right\rangle \right\rangle \\
& + \frac{1}{2N} (gp + (1-g)pq) - \frac{p(q+1)}{2N} \ln(N)
\end{aligned} \tag{24}$$

the last two terms above are negligible when pq/N is small. For $g=1, q=1$, F is the conventional free energy. One must remember that this equation makes sense only when used at the minimum of F with respect to variations in m^μ :

$$\frac{\partial F}{\partial m^\mu} = 0 \tag{25}$$

Eqs. (17) and (18) result from Eq. (25) when we set h^μ (conjugate variables to m^μ) to zero.

C. Composite pattern mean field solutions at T=0

We intend to investigate the retrieval of composite patterns. At $T=0$ the composite patterns assume an idealized form that will be generalized in Section F for $T>0$. In each subdivision the ideal composite pattern is the same as one of the imprinted patterns. However in different subdivisions the pattern may belong to a different imprinted pattern. We can identify a particular composite pattern by a set of indices $\{\mu\}$, μ represents the pattern contained in subdivision . A negative value of μ represents the presence of the inverted pattern $-\mu$. The set of subdivisions where μ occurs are called A_μ and that of $-\mu$, B_μ , that is:

$$A_\mu = \left\{ \left| \mu = \mu \right. \right\} \tag{26}$$

$$B_\mu = \left\{ \left| \mu = -\mu \right. \right\}$$

By symmetry (Eq. (22)), all rearrangements of these indices lead to the same retrieval problem. It is more convenient to characterize the composite patterns by a_μ , the number of subdivisions that contain a particular imprinted pattern, and b_μ , the number of subdivisions that contain its inverse. These numbers can be zero for some μ . There are some immediate identities:

$$\# A_\mu = a_\mu, \# B_\mu = b_\mu \tag{27}$$

$$(a + b) = q \quad \bigcup_{\mu} (A_{\mu} \quad B_{\mu}) = \{1, \dots, q\}, \quad (28)$$

$$A_{\mu} \quad B = A_{\mu} \quad A = B_{\mu} \quad B = \text{if } \mu \quad (29)$$

where # is the set number operator. These identities are useful for evaluating the expressions that follow. We note that if only one a_{μ} is non-zero and all b_{μ} are zero then the composite pattern corresponds to a particular imprinted pattern.

Using the form of ideal composite patterns described above in the definition of the order parameters Eq. (14) we derive the following form for the order parameters

$$m^{\mu 0} = m(a_{\mu} - b_{\mu}) \quad (30)$$

$$m^{\mu} = \begin{cases} m & \text{if } A_{\mu} \\ -m & \text{if } B_{\mu} \\ 0 & \text{otherwise} \end{cases} \quad (31)$$

which we will call the ideal composite form. The composite form described above would imply $m=1$. We have included the variable m to provide an extra degree of freedom in the equations. Our analysis of the mean field equations at $T=0$ will verify that the form Eqs. (30) and (31) is only possible for $m=1$.

We insert the ideal composite form of m^{μ} (Eqs. (30) and (31)) into the mean field equations (Eqs. (17) and (18)) changing the tanh function to sign function for the $T=0$ limit. We consider only the case where $h^{ext} = 0$. From Eq. (17), for $m^{\mu 0}$ we obtain:

$$m(a_{\mu} - b_{\mu}) = \left\langle \left\langle \mu \text{sign} \frac{1}{q} (gm^0 + (1-g)m) \right\rangle \right\rangle \quad (32)$$

Eq. (18) for m^{μ} separates into two cases depending on whether the order parameter m^{μ} is zero or non zero. If the order parameter is non-zero then we have A_{μ} or B_{μ} . Eq. (18) then becomes (+ signs for A_{μ} , - signs for B_{μ}):

$$\pm m = \left\langle \left\langle \mu \text{sign} \frac{m}{q} (g(a - b)) \pm (1-g) \mu \right\rangle \right\rangle \quad (33)$$

If the order parameter is zero, then we have $A_{\mu} \quad B_{\mu}$, we can write A or B for some μ . Eq. (18) becomes (+ signs for A , - signs for B):

$$0 = \left\langle \left\langle \mu \text{sign} \frac{m}{q} (g(a - b)) \pm (1-g) \right\rangle \right\rangle \quad (34)$$

We can summarize Eqs. (33) and (34) as:

$$\left\langle\left\langle \mu \operatorname{sign} \frac{m}{q} (g(a-b)) \pm (1-g) \right\rangle\right\rangle = \pm m^\mu \quad (35)$$

In order for the trial solution to be valid this equation must be satisfied, for each μ and all m for which $A = (a = 0)$ for the + sign, and $B = (b = 0)$ for the - sign. As before, Eq. (32) is the sum over Eq. (35) and therefore need not to be considered separately in obtaining the solutions.

Our objective is to solve Eq. (35) for all values of g . We note that Eq. (35) becomes easy to solve for the cases $g=0$ and $g=1$. For $g=0$ (a totally subdivided network) this equation is always valid, with $m = \pm 1$. Thus, for the totally subdivided network all possible composite patterns are solutions. For $g=1$ (a totally connected network) the equations for μ and $-\mu$ have the same left hand sides and therefore can be made consistent only when we have just one non zero a_μ or b_μ which then has to be q (in this case μ should not be considered at all). Thus, for the totally connected network only the complete imprinted patterns are solutions.

For convenience, we assume in what follows that $g \in (0,1)$ and divide the argument of the sign function by g :

$$\left\langle\left\langle \mu \operatorname{sign} \left((a-b) \pm \left(\frac{1}{g} - 1 \right) \right) \right\rangle\right\rangle = \pm |m|^\mu \quad (36)$$

First we consider the case when all of the b_μ s are zero

$$\left\langle\left\langle \mu \operatorname{sign} \left(a \pm \left(\frac{1}{g} - 1 \right) \right) \right\rangle\right\rangle = |m|^\mu \quad (37)$$

the goal is to find the range of g for which this is valid for a special choice of $\{a_\mu\}$. When only one a_μ is non zero (an imprinted pattern) all g values satisfy Eq. (37). For two or more non zero a_μ the result, after some analysis, is

$$g < g_{\max} = \frac{1}{1 + q - 2a_{\min}} \quad \& \quad |m|=1 \quad (38)$$

Details of the derivation of Eq. (38) are described in the Section A of the Appendix. An alternate derivation using signal-to noise analysis will be published separately.

For the case when not all b_μ are zero, in Section B of the Appendix we show that the condition on g becomes

$$g < g_{\max} = \frac{1}{1 + |a-b|} \quad \& \quad |m|=1 \quad (39)$$

This condition applies unless for each μ $a_\mu b_\mu = 0$. When all $a_\mu b_\mu = 0$ every imprinted pattern appears in the composite pattern with only one sign, and by reversing the sign of the patterns that have $b_\mu = 0$ we return to the previous case where all $b_\mu = 0$. This is verified by the direct analysis in Section B of the Appendix. Note that Eq. (39) does not reduce to the Eq. (38) simply by taking $b_\mu = 0$.

We have shown that at $T=0$ the ideal composite patterns are solutions of the mean field equations in the restricted range $g < g_{\max}$ given by Eqs. (38) and (39). This does not guarantee that they are minima of the free energy (i.e. stable). In Section E below we will show by direct computation of the second order derivatives of F (the stability matrix) that they are stable near $T=0$ throughout the range of $g < g_{\max}$. Furthermore, in Section F we solve numerically for the phase diagrams of minima of F - the range of both g and temperature for which composite patterns are stable.

It is interesting to consider how the two ranges of g values in Eqs. (38) and (39) compare. To make this comparison, with no loss of generality, we assume that $a_\mu = b_\mu$ for all μ . Then Eq. (39) becomes :

$$g < \frac{1}{1 + q - 2} \frac{1}{b} \quad \& \quad |m|=1 \quad (40)$$

When all regions of the network that represent a single imprinted state have the same sign, the value of g is limited by the minimum size of the network representing a particular pattern (Eq. (38)). However, when patterns are present with inverted portions, the value of g is limited by the sum over all inverted parts (Eq. (40)). This occurs because the instability of the smallest portion of the composite pattern arises from the random fields generated by all other imprinted patterns. A destabilizing field arising from a portion of the network representing another pattern is coherent because their field tries to reconstruct the remainder of that pattern. However, the coherent destabilizing field is canceled to the extent that its inverse appears in the composite state. One way to look at Eq. (39) is that the smallest portion of the composite pattern (minimum b) is strengthened by the other b . This result persists to the finite temperature case in Section F.

D. Spurious solutions.

In addition to the composite solutions there are additional solutions of the mean field equations that correspond directly to the spurious states of fully connected networks. We can include some of these solutions by generalizing the forms Eqs. (30) and (31) to include a single variable instead of all previously vanishing parts:

$$\begin{aligned}
m^\mu &= 1 \quad \text{if } A_\mu \\
m^\mu &= m \times -1 \quad \text{if } B_\mu \\
&= e \quad \text{otherwise}
\end{aligned} \tag{41}$$

$$m^{\mu 0} = m \left((a_\mu - b_\mu) + e \left(q - (a_\mu + b_\mu) \right) \right)$$

inserting these in Eq (18) results in

$$\begin{aligned}
m &= \left\langle \left\langle \tanh \frac{m}{q} + \frac{[g(a_\mu - b_\mu) + ge(q - (a_\mu + b_\mu)) \pm (1 - g)]}{[g(a_\mu - b_\mu) + ge(q - (a_\mu + b_\mu)) + e(1 - g)]} \right\rangle \right\rangle \\
me &= \left\langle \left\langle \tanh \frac{m}{q} + \frac{[g(a_\mu - b_\mu) + ge(q - (a_\mu + b_\mu)) \pm e(1 - g)] + (1 - g)(1 - e)}{[g(a_\mu - b_\mu) + ge(q - (a_\mu + b_\mu)) + e(1 - g)]} \right\rangle \right\rangle
\end{aligned} \tag{42}$$

The case $e=0$ corresponds to the composite solutions that we have already examined. The symmetric spurious solutions correspond to the case $e=1$ and $B_\mu = 1$ for all μ . In this case the two equations are identical and simplify to

$$m = \left\langle \left\langle \tanh \frac{m(qg + 1 - g)}{q} \right\rangle \right\rangle \tag{43}$$

or

$$m = \left\langle \left\langle \mu \tanh m \bar{g} \right\rangle \right\rangle \tag{44}$$

where

$$\bar{g} = \frac{1}{q} (1 + g(q - 1)) \tag{45}$$

is the average synaptic strength. This case corresponds to having all patterns with the same order parameter in each subdivision. This corresponds to the symmetric spurious patterns that have been studied for a fully connected network.⁸ For the partially subdivided network the solutions are the same as for the fully connected network but with a rescaled temperature $\beta \rightarrow \bar{g}$. Because the order parameters are the same in each subdivision only the average synaptic strength plays a

role. This is the same rescaling as would occur for stability of one of the imprinted patterns, as discussed in Section F (see Eq. (55)).

In addition to the symmetric spurious patterns there are non-symmetric spurious solutions that can be identified by taking the limit of Eqs. (17) and (18) at $T=0$ and $g=0$. Eq. (17) was derived by canceling a factor of g on both sides of Eq. (13), in this limit it will be identically zero and always valid. Eq. (18) reduces to

$$m^\mu = \left\langle \left\langle \begin{matrix} \mu \text{ sign} \\ m \end{matrix} \right\rangle \right\rangle \quad (46)$$

Since $g=0$ the subdivisions are independent and different s are decoupled. This is the general equation for mean field solutions for a single network at $T=0$. The number of subdivisions q plays no role in determining the form of solutions. The symmetry $m^\mu = -m^\mu$ for any fixed μ but for all s simultaneously, reduces to $m^\mu = -m^\mu$ for any μ or s independently. In addition, if for some μ , m^μ is zero, then for this μ Eq. (46) is automatically satisfied and it does not enter the equations for other μ . To find all possible solutions we find the possible sets of positive $\{m^\mu, \mu = 1, \dots, p\}$ for fixed s that are solutions to Eq. (46). By inversion symmetry we need only consider either the plus or minus sign. Any combination of these for each s with either sign and extra zero elements is also a solution.

The simplest nontrivial solution occurs for $p=1$ and is $m=\pm 1$. Adding zeroes and combining them for different s 's will generate all of the composite patterns discussed in the previous section.

Other solutions arise when one considers special combinations of m^μ in which the sign function in Eq. (46) vanishes for some of $\{s\}$. Using the relevant convention $\text{sign}(0)=0$ the averaging results in $|m^\mu| < 1$. For $p=2$ the only extra solution of the mean field equations are m^μ equal to $\{\pm 1/2, \pm 1/2\}$, for $p=3$ the additional solution is $\{\pm 1/2, \pm 1/2, \pm 1/2\}$, and for $p=4$ there are three other solutions: $\{\pm 5/8, \pm 3/8, \pm 3/8, \pm 1/8\}$, $\{\pm 1/2, \pm 1/2, \pm 1/4, \pm 1/4\}$, $\{\pm 3/8, \pm 3/8, \pm 3/8, \pm 3/8\}$.

These solutions correspond to some of the known spurious solutions for a single network. While they are solutions of Eq. (46), not all are stable. All the patterns for even p are unstable. Our simulations confirm that they are unstable for every $g>0$ even at $T=0$ and evolve to other solutions. The $p=3$ solutions are stable. When the $p=3$ spurious solutions are present in all subdivisions (all s) these are the same symmetric spurious patterns discussed earlier. Combining spurious solutions with composite solutions in different subdivisions also leads to stable states.

E. Stability matrix

To investigate the stability of the composite patterns we study the eigenvalues of the second derivative of free energy (Eq. (24)) written in terms of m^μ . The first derivatives of F are:

$$\frac{F}{m^{\mu 0}} = \frac{g}{q^2} m^{\mu 0} - \left\langle \left\langle \mu \tanh \frac{q}{m^{\mu 0} + (1-g)m} \right\rangle \right\rangle \quad (47)$$

$$\frac{F}{m^\mu} = \frac{1-g}{q^2} m^\mu - \left\langle \left\langle \mu \tanh \frac{q}{m^{\mu 0} + (1-g)m} \right\rangle \right\rangle \quad (48)$$

Defining the second derivative matrix

$$\frac{\partial^2 F}{\partial m^\mu \partial m^{\mu'}} = A^{\mu, \mu'} \quad (49)$$

we derive the different elements as

$$A^{\mu 0, \mu' 0} = \frac{g}{q^2} \left(\mu \mu' (1-g) + \frac{g}{q} Q^{\mu \mu'} \right) \quad (50)$$

$$A^{\mu 0, \mu'} = -\frac{g(1-g)}{q^3} \left(\mu \mu' - Q^{\mu \mu'} \right)$$

$$A^{\mu, \mu'} = \frac{1-g}{q^2} \left(\mu \mu' - \frac{(1-g)}{q} \left(\mu \mu' - Q^{\mu \mu'} \right) \right)$$

where we have defined

$$Q^{\mu \mu'} = \left\langle \left\langle \mu \mu' \tanh^2 \frac{q}{m^{\mu 0} + (1-g)m} \right\rangle \right\rangle \quad (51)$$

In the limit $T \rightarrow 0$ we see that $Q^{\mu \mu'} \rightarrow \mu \mu'$ when the argument of tanh doesn't vanish for any $\{\mu\}$. This condition is guaranteed for composite solutions when $g < g_{\max}$ (with no equal sign) as shown in Section C of the Appendix. Inserting this limit of $Q^{\mu \mu'}$ into Eq. (50) implies that for $T \rightarrow 0$ and $g < g_{\max}$ the second derivative matrix is diagonal with positive eigenvalues. Thus the composite solutions are stable near $T=0$.

For the more general case of arbitrary temperature one must find the eigenvalues of A. We construct $B = A - I$

$$B^{\mu 0, \mu' 0} = \mu \mu' (1-g) - \frac{g}{q} \left(\mu \mu' - Q^{\mu \mu'} \right) \quad (52)$$

$$B^{\mu 0, \mu'} = -\frac{\sqrt{1-g}}{q} \left(\mu \mu' - Q^{\mu \mu'} \right)$$

$$B^{\mu, \mu'} = \mu^{\mu'} (1 - g) - \frac{(1 - g)}{q} (\mu^{\mu'} - Q^{\mu\mu'})$$

The determinant of B must be equated to zero to solve for g . B has a simpler form ($B = B_1$) after performing some elementary row and column operations (that do not change the determinant):

$$B_1^{\mu 0, \mu' 0} = \mu^{\mu'} (1 - g) + \frac{gq}{1 - g} \quad (53)$$

$$B_1^{\mu 0, \mu'} = B_1^{\mu, \mu' 0} = -\mu^{\mu'} (1 - g) \sqrt{\frac{g}{1 - g}}$$

$$B_1^{\mu, \mu'} = \mu^{\mu'} (1 - g) - \frac{(1 - g)}{q} (\mu^{\mu'} - Q^{\mu\mu'})$$

In order to make use of this expression it is necessary to know m^{μ} for $T = 0$. In the following section we describe numerical solution of the mean field equations that explicitly determine the phase boundaries of the stability of various composite states.

F. Phase diagrams of composite states

For $T > 0$ the ideal composite form of the order parameters, Eqs. (30) and (31), must be modified because it includes only one unknown m and there are two or more equations to be satisfied (Eqs. (17) and (18) for different μ and μ'). These equations are degenerate only at $T=0$. For $T > 0$ we generalize the composite forms by assuming the most general form of m^{μ} . For each ideal composite pattern the value of the order parameters $m^{\mu}(T, g)$ are then obtained by continuation from the ideal composite solutions at $T=0, g=0$. There is a range of g and T over which the continuation is stable. The boundary of this region is analogous to a phase transition. The ideal composite patterns at $T=0$ correspond to the retrieval of one and only one imprinted pattern μ in each subdivision. The order parameter of a particular pattern μ will decrease (from the value 1 at $T=0$) at a rate that depends on the number of subdivisions that contain it as specified by (a_{μ}, b_{μ}) . Moreover, the existence of a non-zero order parameter for a pattern in one subdivision will cause a non-zero order parameter for the same pattern in all other subdivisions.

To study the phase transitions of the composite patterns we started with the ideal form of each composite pattern for $g=0$ at $T=0$. Gradually increasing g and T we performed iterative minimization of the free energy. Conjugate gradient minimization was used to find the closest local minimum of the free energy with respect to the order parameters m^{μ} . We then located the temperature for each g at which the composite pattern would no longer be stable (at this point it will typically evolve discontinuously to other solutions that are stable, if there are any). We varied each of the m^{μ} as an independent variable. Care must be taken to break symmetry at

every step using small random perturbations to m^μ to avoid the conjugate gradient (or steepest descent) remaining upon a saddle surface. The phase diagrams and order parameters that resulted are illustrated in Figs. 2-7.

Numerical solutions of the mean field equation were performed for the cases $q=2,3$ and 4. In each case we took the number of non zero values of a_μ to range from 0 to q , the maximum allowed in a composite pattern. We specify the form of each composite pattern by the value of $(a_\mu - b_\mu)$ of each of the patterns (dropping $b_\mu=0$ for conciseness). For example, in the $q=4$ case [4000] represents $a_1=4$ and $a_2=a_3=a_4=0$. This is one imprinted pattern retrieved in the whole network. [1111] is the composite pattern with four different patterns retrieved in the four different subdivisions $a_1=a_2=a_3=a_4=1$, [2(1-1)00] is the pattern with $a_1=2, a_2=b_2=1$ and $a_3=a_4=0$. Note that Eq. (28) is satisfied.

For each q we plot phase diagrams showing the transition temperature as a function of g . Each pattern is stable below its phase transition line. These diagrams enable us to determine domains in the phase diagram where particular kinds of patterns are stable while others are not. We also show plots of the order parameters for some of the composite patterns as a function of temperature at the value $g=0.1$. We discuss below some conclusions that can be reached from these diagrams.

The phase transition diagrams show that in all three cases the imprinted patterns are the most stable. Also their transition line is straight. This can be derived directly using the form of the order parameters for retrieval of an imprinted pattern $\{ \}$:

$$\begin{aligned} m^\mu &= m^\mu \\ m^{\mu 0} &= q m^\mu \end{aligned} \tag{54}$$

in Eqs. (17) and (18). This results in the usual Ising model equation for the order parameter m with rescaled temperature

$$\begin{aligned} m &= \tanh(\beta m) \\ &= \frac{(gq + 1 - g)}{q}. \end{aligned} \tag{55}$$

Thus the pattern is stable for β greater than 1. Setting $\beta = 1$ gives the phase transition line for the imprinted pattern. As derived numerically $T(g)$ is linear and passes through the points $(g = 0, T = 1/q)$ and $(g = 1, T = 1)$.

For the composite patterns there is a hierarchy of descending ranges of stability. As a rule the composite states that have a set of $\{a_\mu\}$ with higher symmetry are stable at higher temperatures. Specifically patterns with equal a_μ or b_μ have higher stability. For example compare [111] with [210] and others in the $q=3$ case. For $q=4$ compare [2200] with [2110] or others whose transition

lines are below that of [1111]. Also for $q=4$ compare [1111] with [(1-1)110] or with [2110] or with [3100] that are less symmetric.

Comparing [2200] and [1111] shows that among the symmetric patterns those with greater a_μ are more stable. This is reasonable since the reoccurrence of one pattern in more subdivisions strengthens its retrieval. However for nonsymmetric patterns, such as [2110] as compared to the symmetric pattern [1111], increasing a_μ for one of the patterns has the effect of lowering the stability of the other patterns and thus lowering the stability of the whole composite pattern. This trend is maintained when going to the pattern [3100] which is even less stable than [2110].

For all cases $q=2,3,4$ there are distinct composite patterns that appear to have identical transitions. This usually happens when the value of an a_μ is split between a_μ and b_μ . For $q=3$ the case of [210] and [(1-1)10], and for $q=4$ the cases of [2200] and [(1-1)(1-1)00], [3100] and [(2-1)100]. Alternatively this can happen when two a_μ are combined into $(a_\mu - b_\mu)$ for one pattern. For $q=2$ the case of [11] and [(1-1)0], for $q=3$ the case of [210] and [(2-1)00] and for $q=4$ the cases of [2200] and [(2-2)000], [2(1-1)00] and [(2-1)100], [3100] and [(3-1)000]. If these splittings or recombinations change the symmetry of the pattern (as discussed in the previous paragraph) they do not overlap which indicates the priority of the symmetry (compare [111] with [(1-1)10], and [2200] with [2(1-1)00]). There are cases where the transition lines are not the same even without a change in symmetry that are as yet unexplained. For example, compare [2110] with [2(1-1)00], and [2110] with [(1-1)110]. In these cases the patterns may differ in symmetries more complicated than the one mentioned above. We have been able to show analytically the equivalence of the phase transition line for [11] and [(1-1)0]. A general derivation for all cases has not yet been found.

All transition lines have the same value at $g=0$. For the fully disconnected network the pattern in each of subdivisions is retrieved independent of the other subdivisions. Composite patterns or imprinted patterns have the same transition temperature $T=1$. This can be seen also by setting $g=0$ in Eqs. (17) and (18) which decouples different μ .

The value of g at which the transition lines reach $T=0$ agrees in all cases with the analytically derived value of g_{\max} given in Eqs. (38) and (39). When a composite pattern has higher g_{\max} it has a higher transition temperature for all g (or a higher transition g for all T). However, in many cases composite patterns with the same g_{\max} have different transition temperatures in the range $0 < g < g_{\max}$.

Plots of the order parameters, Figs. 3, 5 and 7, show that the phase transition for the imprinted patterns are second order but for composite patterns they are all of first order. The plots shown are all for the same value of $g=0.1$. The height of the discontinuity in the order parameters increases as a function of g . Starting from a value of zero at $g=0$ it always reaches the value $m = 1$ at g_{\max} .

V. Conclusions

We have analyzed the retrieval of composite patterns in a partially subdivided network. The existence of composite solutions of the mean field equations shows that imprinting patterns on a partially subdivided network results in retrieval of an expanded set of composite patterns. The kind of composite patterns that are retrieved depends on the strength of the synapses between the subdivisions.

The degree of functional localization in the brain has long been a controversial subject. We have attempted to provide a framework in the context of the theory of attractor networks in which questions about functional separation and its utility may be formulated in a more precise language. The expansion of memories from the training set to the set of combinations of trained subnetwork states is a strategy for generalization by the network that may be used to incorporate prior knowledge about correlations. Conventional attractor networks generalize because training corresponds to creating a local minimum in the vector space of network states -- the basin-of-attraction of this state becomes its generalization. Partially subdivided networks generalize by recognizing various combinations of substates. Since the network has been trained on far fewer states than it recognizes it may be said to have generalized from the training set to the set of recognized states. This is an advantage if the architecture of subdivision is in direct correspondence to the information to be represented. This is a first step to understanding various aspects of generalization and creativity in the form of combining aspects of learned information in new ways.

The use of a combinatorial expansion of substates of a particular complex system is not restricted to neural networks. Another example exists in the function of the immune system⁹ where the genetic code for the immune cell receptors that detect antigen are formed from a set of seven pieces taken from the cell genetic code. The set of imprinted states is analogous to the possible sequences of each of the DNA segments, and the state of the receptor becomes a composite state. This combination of different pieces into composites enables the cells to construct a large variety of receptors from a small set of initial components. Interactions between the different DNA components arise because of the process of expression of the gene. The final structure of a receptor relies upon all of the genetic components which therefore interact - they are not fully independent. This has a similar flavor to the consideration of partially subdivided networks.

References:

- ¹ V. S. Dotsenko J. Phys. C: Solid State Phys. 18, L1017 (1985)
- ² C. R. Wilcox, J. Phys. A: Math. Gen. 22, 4707 (1989)
- ³ V. Deshpande, and C. J. Dasgupta, Stat. Phys. 64, 755 (1991)
- ⁴ J. P. Sutton, J. S. Beis, and L. E. H. Trainor, J. Phys. A 21 4443 (1988)
- ⁵ M. A. P. Idiart and A. Theumann, J. Phys. A: Math. Gen. 24, L649 (1991); It may be helpful to point out that the "mixed" states discussed by these authors are only mixed in the hierarchical representation, but consist of pieces of only a single imprinted neural state.
- ⁶ M. W. Mao, and J. B. Kuo, Neural Networks, 5, 835 (1992)
- ⁷ for a recent discussion see S. Zeki, Scientific American, (September, 1992) p. 69
- ⁸ D. J. Amit, *Modeling Brain Function, the world of attractor neural networks*, (Cambridge University Press, Cambridge, 1989)
- ⁹ see, for example, C. A. Janeway, Jr., Scientific American, (September, 1993) p. 73

Appendix: Range of Stability of Composite Patterns at T=0

In this Appendix we describe the solution of the zero temperature mean field equations for the composite patterns (Eq. (36)). In Section A we prove Eq. (38) which applies to the composite patterns for which every stored pattern appears with only one sign. With no loss of generality this can be restated as $b = 0$ for all μ . In Section B we prove Eq. (39) which applies to the more general case where some patterns appear with both signs in the composite pattern. These sections only describe the validity of composite patterns as solutions of the mean field equations, Eqs. (17) and (18). In Section IV.E of the text we prove that the composite patterns are memories - stable states near T=0 of the neural dynamics - by showing that the eigenvalues of the second derivative matrix of the free energy are all positive. In Section C of this Appendix we prove an inequality that is needed in the analysis of Section IV.E.

We note that any imprinted patterns that do not appear in the composite state do not affect the mean field solutions at low storage. These patterns satisfy $a = b = 0$. In solving Eq. (36) we first perform the summation and averaging on all μ for which $a = b = 0$, since these have no effect. In what follows we consider only the set of μ for which a or b is non zero ($a + b \neq 0$).

A. Composite patterns with all $b = 0$

The relevant equation to be solved, after setting all b equal to zero in Eq. (36), is Eq. (37). The objective is to find the range of g for which this equation is valid for a particular choice of $\{a_\mu\}$. The result is Eq. (38), where a_{\min} represents the smallest non zero a_μ . This form is valid when two or more a_μ are non zero. The case of exactly one non zero a_μ is discussed in the text following Eq. (37): when only one a_μ is non zero all $g \in [0,1]$ satisfy Eq. (37). In the rest of this section we assume that two or more a_μ are non zero.

We first rewrite Eq. (37) for the two cases $\mu = +$, $\mu = -$ and take μ inside the sign function

$$\left\langle \left\langle \text{sign } a_\mu + \frac{a_\mu}{g} + \left(\frac{1}{g} - 1\right) \right\rangle \right\rangle = 0 \quad (\text{A1})$$

$$\left\langle \left\langle \text{sign } a_\mu + \frac{a_\mu}{g} + \left(\frac{1}{g} - 1\right) \right\rangle \right\rangle = |m| \quad (\text{A2})$$

Since

$$\frac{a_\mu}{g} + \left(\frac{1}{g} - 1\right)$$

changes sign whenever $\left\{ \frac{a_\mu}{g} + \left(\frac{1}{g} - 1\right) \right\} \in \{-1, 0, 1\}$, and a_μ is positive, Eq. (A1) is valid if and only if:

$$a_\mu - \left| \frac{a}{\mu} + \left(\frac{1}{g} - 1\right) \right| < 0 \quad (\text{A3})$$

for all $\{ \}$ and for all μ . In a similar way if we want a nontrivial solution with $m = 0$ Eq. (A2) is valid if and only if for some $\{ \}$

$$a_\mu + \left(\frac{1}{g} - 1\right) - \left| \frac{a}{\mu} \right| > 0 \quad (\text{A4})$$

The number of $\{ \}$ satisfying Eq. (A4) determines the magnitude of m . This number must be the same for all μ . This originates in the ideal composite form where we assumed that the magnitude of m^μ is independent of μ .

The inequalities (A3), (A4) may be solved to obtain the range of g in which they are valid. Inequality (A3) implies that for every μ and $\{ \}$, and for every μ either

$$a_\mu + \left(\frac{1}{g} - 1\right) + \frac{a}{\mu} < 0 \text{ or } a_\mu - \left(\frac{1}{g} - 1\right) + \frac{a}{\mu} < 0 \quad (\text{A5})$$

must be satisfied. It is sufficient to consider $= 1$ because these two inequalities transform into each other under the transformation $\{ \} \rightarrow \{- \}$. Then we derive the equivalent constraints as either

$$a_\mu < \left(\frac{1}{g} - 1\right) + \frac{a}{\mu} \text{ or } a_\mu < -\left(\frac{1}{g} - 1\right) + \frac{a}{\mu} \quad (\text{A6})$$

Since we have set $= 1$ at least one must be equal to one in the sum $\frac{a}{\mu}$. The inequalities in (A6) lead to the following compound logical statement.

For every μ and every $\{ \}$ with at least one (μ) equal to one

$$\text{either } \frac{a}{\mu} = a_\mu \quad (\text{P1a})$$

$$\text{or } g < \frac{1}{1 + a_\mu - \frac{a}{\mu}} \quad (\text{P1b})$$

$$\text{or } \frac{a}{\mu} < -a_\mu \text{ and } g > \frac{1}{1 - a_\mu - \frac{a}{\mu}} \quad (\text{P2})$$

and also for every μ there exist some $\{ \}$ such that

$$\text{either } \left| \frac{a}{\mu} \right| = a_\mu \quad (\text{P3a})$$

or
$$g < \frac{1}{1 - a_\mu + \frac{a}{\mu}} . \quad (\text{P3b})$$

The number of patterns that satisfy (P3a) or (P3b) determines the magnitude of m . More precisely for each μ , if the number of $\{ \}$ satisfying (P3) is K then m derived from (A2) is equal to $K/2^{\hat{p}}$ where \hat{p} is the number of non zero a_μ . This statement is composed out of several conditions. (P1a) and (P1b) arise from the first part of Eq. (A6) and (P2) arises from the second part of Eq. (A6). (P3) results from a similar analysis of (A4). In what follows we show that (P2) is never satisfied for $m = 0$. We assume this result for the moment. Then (P1a) or (P1b) must be true for all μ and $\{ \}$. The range of acceptable g values is determined by the value of μ and the pattern $\{ \}$ that sets the most restrictive limits using (P1a) or (P1b). This is obtained by considering the particular $\{ \}$ for which only one $a_\mu = 1$ where a is minimal and all other $a_\mu = -1$. Then (P1) results :

$$g < g_{\max} = \frac{1}{1 + \frac{a - a_{\min}}{a}} \quad (\text{A7})$$

which is the first part of Eq. (38). Moreover, for this range of g (P3b) will be automatically satisfied for all $\{ \}$ thus $m=1$. This is the second part of Eq. (38).

It remains to be demonstrated that condition (P2) can never be satisfied. We note that (P2) and (P3) are mutually exclusive. When m is non zero (P3) must be true for some $\{ \}$ thus (P2) cannot be true for all $\{ \}$. There remains the possibility that for some μ there exists a $\{ \}$ (the set of these is denoted as $\{ \{ \} \}^{P2}$), for which (P2) is valid, and for the rest (P1) is satisfied (denoted by $\{ \{ \} \}^{P1}$). Some of the $\{ \}$ in $\{ \{ \} \}^{P1}$ must satisfy (P3). The sets $\{ \{ \} \}^{P1}$ and $\{ \{ \} \}^{P2}$ may be μ dependent.

We begin by considering only the μ for which a_μ has its largest value. If there is more than one μ with the maximal value of a_μ anyone of them may be used. We begin by assuming a non-empty set $\{ \{ \} \}^{P2}$ and demonstrate a contradiction.

Since condition (P2) gives a lower limit on g and condition (P1) gives an upper limit on g the of a non-empty set $\{ \{ \} \}^{P2}$ implies a range for g of the form $g_{\max}^{P2} < g < g_{\min}^{P1}$. In this range m would be less than one. Here g_{\max}^{P2} is the maximum of

$$\frac{1}{1 - a_\mu - \frac{a}{\mu}} \quad (\text{A8})$$

over the set $\{ \{ \} \}^{P2}$. This expression reaches its maximum value when the sum a in the denominator is maximal over $\{ \{ \} \}^{P2}$. We call this maximum value g_{\max}^{P2} and a pattern μ for which

the maximum value is attained $\{ \}^P_{\max}$. Similarly we define $\{ \}^P_{\min}$ as the minimum of a over $\{ \}^P$ and a pattern for which the minimum value is attained $\{ \}^P_{\min}$.

There are two possibilities depending on whether $\{ \}^P_{\min}$ satisfies (P1a) or (P1b). If $\{ \}^P_{\min}$ satisfies (P1a) then all members of $\{ \}^P$ must satisfy (P1a). In this case $g^P_{\min}=1$; applying (P1a) to $\{ \}^C_{\min}$ and applying the first part of (P2) to $\{ \}^P_{\max}$ we have that $\frac{P1}{\min} - \frac{P2}{\max} > 2a_{\mu}$. The other possibility is that $\{ \}^P_{\min}$ satisfies (P1b), then g^P_{\min} is the minimum of

$$\frac{1}{1 + a_{\mu} - \frac{a}{\mu}} \quad (A9)$$

over $\{ \}^P$. Applying (P1b) to $\{ \}^P_{\min}$ and applying the second part of (P2) to $\{ \}^P_{\max}$, and using $g^P_{\min} > g^P_{\max}$ we still have that $\frac{P1}{\min} - \frac{P2}{\max} > 2a_{\mu}$, which therefore applies in all cases.

When we apply $\frac{P1}{\min} - \frac{P2}{\max} > 2a_{\mu}$ to the μ for which a_{μ} is the largest we find

$$\frac{P1}{\min} - \frac{P2}{\max} > 2a_{\max} \quad (A10)$$

This is impossible as we now demonstrate. The pattern $\{ \}^P_{\max}$ has at least one $=-1$ because it satisfies the first inequality in (P2). By changing the sign of this particular we arrive at a new pattern which must be in $\{ \}^P$. The value of a for this pattern is equal to $\frac{P2}{\max} + 2a$, so $\frac{P1}{\min}$ which is the minimum of all sums over $\{ \}^P$, must be less than or equal to $\frac{P2}{\max} + 2a$ which contradicts (A10). This proves the assertion that $\{ \}^P$ is a null set for this special μ . This shows that for the particular μ for which a_{μ} is maximal (P2) can never be valid.

For the case where a_{μ} is maximal we have shown that the condition (P1) applies to any pattern $\{ \}$. We consider the application of (P1) to a specially constructed pattern. In this pattern $\{ \}$ is +1 only for one particular and -1 for all the others. The value is chosen to be one of the for which a achieves its minimum (non zero) value. By inspection, the specially constructed pattern can not satisfy (P1a), it must therefore satisfy (P1b). (There is one special case where it satisfies (P1a), i.e. when $a_1 = a_2 = q/2$, however (P1b) gives no additional restriction for this case, so we could say it is valid). This gives the most limiting condition on the value of g . Not only is this the most limiting from all patterns for this μ , but also for all μ (it is equivalent to the condition we have previously obtained in Eq. (A7)). This implies that for all μ (P1) must be satisfied for all patterns and rules out (P2) for any μ , completing the proof.

B. Composite patterns with at least one v such that $a \neq 0, b \neq 0$.

For the general case of non zero a_μ and b_μ we start from Eq. (36) We first absorb the sign of $a - b$ into d and define $d = |a - b|$

$$\left\langle\left\langle \left[\text{sign}(a_\mu - b_\mu) \right]^\mu \text{sign} \left[\sum d \pm \text{sign}(a - b) \left(\frac{1}{g} - 1 \right) \right] \right\rangle\right\rangle = \pm |m|^\mu \quad (\text{B1})$$

or

$$\left\langle\left\langle \mu \text{sign} \left[\sum d \pm \text{sign}(a - b) \left(\frac{1}{g} - 1 \right) \right] \right\rangle\right\rangle = \pm \left[\text{sign}(a_\mu - b_\mu) \right] |m|^\mu \quad (\text{B2})$$

In these equations the choice of \pm on the left and right are coupled. As explained in the main text after Eq. (35) the equations must be satisfied for all μ with a $+$ sign when $A > 0$ ($a > 0$) and with a $-$ sign when $B > 0$ ($b > 0$). If both $a > 0$ and $b > 0$ then the equation must be satisfied with both signs. Combining this sign with the factor $\text{sign}(a - b)$ we define a variable $s = \pm \text{sign}(a - b)$ (B3)

s takes the value $+1$ when only one of (a, b) are not zero or equivalently $a b \neq 0$. s takes both possible values ± 1 when $a b = 0$ (when they are both non-zero). It has been assumed that at least one is not zero, or $a + b \neq 0$. This leads to the equations

$$\left\langle\left\langle \text{sign} \left[\sum d_\mu + \mu d + s_\mu \left(\frac{1}{g} - 1 \right) \right] \right\rangle\right\rangle = s_\mu |m|^\mu \quad (\text{B4})$$

$$\left\langle\left\langle \text{sign} \left[\sum d_\mu + \mu d + s \left(\frac{1}{g} - 1 \right) \right] \right\rangle\right\rangle = 0 \quad \mu \quad (\text{B5})$$

The only difference between Eqs. (B4), (B5) and Eqs. (A1) and (A2), besides renaming a as d , is the occurrence of the factor s on both sides and the possibility of zero d for nontrivial cases i.e. when $a + b$ is non zero. These differences will modify the results. Repeating the analysis of Section A following Eqs. (A1) and (A2) we arrive at logical conditions that are analogous to conditions (P1)-(P3). Before proceeding we note that if there is only one non-zero $a + b$ the case $\mu > 1$ doesn't apply and only the case $\mu = 1$ must be considered (see discussion after (B6)).

The set of conditions that must be satisfied by patterns $\{ \}$ are conditions (Q1)-(Q4):

$$\text{Either } \left| \frac{d}{\mu} - d_{\mu} \right| < \frac{1}{\mu} \quad (\text{Q1a})$$

$$\text{or } g < \frac{1}{1 + d_{\mu} - \frac{d}{\mu}} \quad (\text{Q1b})$$

$$\text{or } \left| \frac{d}{\mu} + d_{\mu} \right| < -d_{\mu} \text{ and } g > \frac{1}{1 - d_{\mu} - \frac{d}{\mu}}. \quad (\text{Q2})$$

This must be applied when there is at least one non zero $a + b$ with μ . There are two possible scenarios. If for all μ $a + b = 0$ (Q1)-(Q2) must be satisfied only for $\{ \}$ with at least one equal to one (μ). Otherwise for each μ the conditions (Q1)-(Q2) must be satisfied for all patterns. The former case is the same as Eq. (P1) and (P2) in part 1 with a replace by d . The latter is more restrictive and arises from a consideration of the effect of s .

Also for each μ , if $a_{\mu}b_{\mu} = 0$ then there exist some $\{ \}$ such that :

$$\text{either } \left| \frac{d}{\mu} - d_{\mu} \right| < \frac{1}{\mu} \quad (\text{Q3a})$$

$$\text{or } g < \frac{1}{1 - d_{\mu} + \left| \frac{d}{\mu} \right|} \quad (\text{Q3b})$$

If $a_{\mu}b_{\mu} \neq 0$ then consideration of the factor s_{μ} results in replacement of condition (Q3) by

$$g < \frac{1}{1 + d_{\mu} + \left| \frac{d}{\mu} \right|} \quad (\text{Q4})$$

the number of patterns that satisfy (Q3) or (Q4) determines the magnitude of m which must be the same for all μ .

We note the resemblance between (P1)-(P3) and (Q1)-(Q3). They differ in two ways: (1) that a_{μ} has been replaced by d_{μ} and (2) the extra condition after (Q2). Thus if for all μ , $a_{\mu}b_{\mu} = 0$ then (Q1)-(Q3) results in the limit analogous to (A7):

$$g < \frac{1}{1 + d - 2d_{\min}} \quad (\text{B6})$$

This is identical to Eq. (38) with a_{\min} generalized to d_{\min} which can be either an a_{μ} or a b_{μ} . This reflects the inversion symmetry Eq. (20).

For the more general case the existence of at least one pattern for which both a_μ and b_μ are non zero allows us to apply the conditions (Q1) and (Q2) to all $\{ \}$ and to obtain the limit of Eq. (39) on g

$$g < \frac{1}{1 + d} \quad (\text{B7})$$

An interpretation of this expression and comparison to Eq. (38) is given in the text.

To prove (B7) in the case when there is only one non zero $a_\mu + b_\mu$ we can only apply (Q4) which by itself will result in (B7). Otherwise if there are more than one non zero $a_\mu + b_\mu$ we first consider the conditions resulting from the largest d_μ .

If the largest d_μ is zero then all d_μ are zero. Applying (Q1)-(Q3) gives no additional restriction on $g < 1$ and also ensures that $m = 1$. This is consistent with the general result (B7).

When the maximum d_μ is non zero instead of (Q3) we must apply (Q4) and we can rule out (Q2) for this special μ the same way we treated (P2) in Section A. Thus for the maximum d_μ (Q1) must be valid. For the following analysis there are two possibilities. Either there exists a μ for which $a b = 0$ or there is no such μ . In the former case it does not matter whether $a_\mu b_\mu = 0$. In the latter case we know that $a_\mu b_\mu = 0$.

If there exists a μ not equal to μ for which $a b = 0$ then the second possibility following (Q2) holds and we must apply (Q1) to all $\{ \}$, including the $= -1$ for all μ . This directly gives (B7).

When there is no such μ for which $a b = 0$ the first possibility following (Q2) holds for μ and we only apply (Q1) to those $\{ \}$ with at least one $= +1$ for some μ . In this case the lowest limit for g derived from (Q1) is of the form (B6). However we must still consider the limits established by considering other μ for which d_μ is not maximum.

We consider any other non zero d which we call $d_{\mu'}$. Then either $a_{\mu'} = 0$ or $b_{\mu'} = 0$. For this μ' the second possibility following (Q2) applies and (Q1) or (Q2) must be applied to all $\{ \}$. Considering again the pattern $= -1$ for all μ , this pattern can not satisfy (Q2) because it gives a lower bound for g of the form

$$g > \frac{1}{1 + d - 2d_{\mu'}} \quad (\text{B8})$$

which is outside the range already derived (see (B6)). Thus this pattern must satisfy (Q1b) which is same as (B7).

Finally we note that the condition (B7) on g implies that for all $\{ \}$ and for all μ (Q3) or (Q4) are satisfied, so $m = 1$. This completes the proof of Eq. (39).

C. Stability of Composite Patterns at T=0:

As shown in Section E the stability matrix has elements of the form Eq. (50). To find the eigenvalues of this matrix at T=0 we find the limit of $Q^{\mu\mu'}$, defined in Eq. (51), at zero temperature.

First we prove that in this limit and for the ideal composite patterns in the range of g specified by Eqs. (38) and (39) the argument of tanh is non zero for any choice of $\{ \}$. For simplicity we call this argument $\left(\{ \} \right)$ such that Eq. (51) can be rewritten as

$$Q^{\mu\mu'} = \left\langle \left\langle \mu \mu' \tanh^2 \left(\left(\{ \} \right) \right) \right\rangle \right\rangle \quad (C1)$$

where

$$\left(\{ \} \right) = \frac{1}{q} \left(gm^0 + (1-g)m \right) \quad (C2)$$

To prove this is non zero we first insert the form of the ideal composite patterns from Eqs. (30) and (31) in Eq. (C2). To write the result in a form already used in this Appendix in Eqs. (A2) and (B4), we factor out g (is defined by the condition $(A \ B)$). With no loss of generality we absorb the sign of $a - b$ into and rename as . Thus we derive

$$\left(\{ \} \right) = \frac{gm}{q} \left(d + s \left(\frac{1}{g} - 1 \right) + d \right) \quad (C3)$$

The terms in the parenthesis in Eq. (C3) are identical to the argument of the sign function in Eq. (B4). We already proved that if g is in the limits specified Eqs. (B6) and (B7) then m=1. Using Eq. (B4) this shows that the term in the parenthesis of Eq. (C3) is positive for all $\{ \}$. Thus $\left(\{ \} \right)$ is non zero for all $\{ \}$ as claimed.

Finally, to obtain the limit of $Q^{\mu\mu'}$ as $T \rightarrow 0$, we can use the expression

$$\lim_{T \rightarrow 0} \tanh^2(x) = 1 - 4e^{-2|x|} \quad (C4)$$

for any non zero x that is independent of to prove that

$$\lim_{T \rightarrow 0} Q^{\mu\mu'} = \mu\mu' + (e^{-}) \quad (C5)$$

Using this limit in Eq. (50) Section IV.E shows that the stability matrix is diagonal with positive eigenvalues.

Figure Captions:

Fig. 1: Illustration of the use of subdivided networks in the context of language. A fully connected network with enough neurons to store exactly nine sentences shown on the left may be imprinted with and recognize these sentences. If the network is divided into three parts it may be imprinted with only three sentences (center). However, because each subnetwork functions independently, all possible twenty-seven combinations of words shown to the right are recognized. Comparing left and right columns suggests the difference between semantics and grammar in sentence construction.

Fig. 2: The transition temperature as a function of subdivision connectivity g in a network with two subdivisions ($q=2$) for an imprinted pattern, [20], and the two possible composite patterns. [11] is the composite pattern with two distinct imprinted patterns in each of the subdivisions. [(1-1)0] is the composite pattern formed out of one imprinted pattern in one subdivision and its own inverse in the other. Below each curve the corresponding pattern is stable and above it the pattern is unstable. The transition temperatures for [11] and [(1-1)0] coincide at all values of g . This can be proven analytically by analysis of the free energy, Eq. (23). Other cases of overlapping phase diagrams occur for different composite patterns for larger values of q . The phase diagrams for $q=3$ and 4 are shown in Figs. 4 and 6 respectively.

Fig. 3: Order parameters m^{μ} as a function of inverse temperature (β) for the imprinted pattern and two composite patterns. The figures show the behavior of the order parameter on the $g=0.1$ cross section of the phase diagram of Fig. 2. (a) The imprinted pattern [20]: The retrieval of the first imprinted pattern is illustrated so the order parameters for the second stored pattern ($m^{2\mu}$) are zero at all temperatures. (b) The composite pattern [11]: The composite pattern is constructed from the first imprinted pattern in the first subdivision and the second imprinted pattern in the second subdivision, so $m^{11} = m^{22} = 1$ at $T \rightarrow 0$ ($\beta \rightarrow \infty$), (see Eq.(38)) and $m^{12} = m^{21} = 0$ in the same limit. For $T > 0$ these patterns dominate in their respective subdivisions but there is some cross over. (c) the composite pattern [(1-1)0]: The composite pattern is constructed from the first imprinted pattern dominating in the first subdivision and its inverse in the second subdivision, so $m^{11} = 1, m^{12} = -1$ and $m^{2\mu} = 0$ at $T \rightarrow 0$ ($\beta \rightarrow \infty$). It is possible to map the order parameters of the composite pattern [(1-1)0] onto the order parameters of the composite pattern [11] showing the equivalence of the phase transition temperature in the two cases.

The order parameters undergo phase transitions at value of β corresponding to the value of the transition temperature at $g=0.1$ in Fig. 2. For the imprinted pattern (a) the transition is second order. For the composite patterns the transitions are first order and more detailed analysis shows that it occurs when the local energy minimum changes to a saddle point. In all cases in the limit T -

>0 (large) non-zero m^μ s become 1. Note that $m^{\mu 0}$ is the sum of m^μ s for each μ . The values of the free energy F at the minimum are also plotted. F changes only very slowly with temperature. Most of the small change occurs near the transition temperature.

Fig. 4: Similar to Fig. 2 phase transition diagram for imprinted and composite patterns in a network with three subdivisions ($q=3$). Here there are more composite patterns in addition to the imprinted pattern ([300]). [111] is the composite pattern with three distinct imprinted patterns retrieved in the three subdivisions. [210] is the composite pattern with the first imprinted pattern retrieved in two of the subdivisions, the second imprinted pattern is retrieved in the other subdivision, no other imprinted patterns are retrieved. [[(1-1)10] is similar to [210] except that in the second subdivision the first imprinted pattern is inverted. Finally [(2-1)00] is the pattern that is formed out of a single imprinted pattern but with one of the subdivisions inverted. The transition temperatures of all the composite patterns except [111] appear to coincide. The $T=0$ transition point of all of the composite patterns agrees with the analytical results Eqs.(39) and (40) which give $g_{\max}=0.5$.

Fig. 5: Order parameters plotted as a function of inverse temperature at the value $g=0.1$ for (a) the imprinted pattern and (b)-(e) several of the composite patterns of a network with $q=3$ subdivisions. See Fig. 4 for the phase transition diagram. Compare with Fig. 3. All of the composite patterns undergo first order transitions while the imprinted pattern has a second order transition.

Fig. 6: Phase transition diagram for imprinted and composite patterns in a network with three subdivisions ($q=4$) (Compare Figs. 2 and 4). The notation for patterns are similar to Fig. 4. The composite patterns fall into two groups according to the prediction of Eq.(39) and Eq.(40), those with $g_{\max}=1$ and those with $g_{\max}=1/3$ for their g intersect. The region near $g=0$, $\beta=0$ is enlarged in (b).

Fig. 7: Order parameters plotted as a function of inverse temperature at the value $g=0.1$ for (a) the imprinted pattern and (b)-(e) several of the composite patterns of a network with $q=4$ subdivisions. See Fig. 5 for the phase transition diagram. Compare with Figs. 3 and 5.

Fully connected network

Subdivided network

Imprinting and Retrieval

Imprinting

Retrieval

Big	Bob	ran.
Kind	John	ate.
Tall	Susan	fell.
Bad	Sam	sat.
Sad	Pat	went.
Small	Tom	jumped.
Happy	Nate	gave.
Mad	Dave	took.
Shy	Cathy	Slept

Big	Bob	ran.
Kind	John	ate.
Tall	Susan	fell.

Big	Bob	ran.
Big	Bob	ate.
Big	Bob	fell.
Big	John	ran.
Big	John	ate.
Big	John	fell.
Big	Susan	ran.
Big	Susan	ate.
Big	Susan	fell.
Kind	Bob	ran.
Kind	Bob	ate.
Kind	Bob	fell.
Kind	John	ran.
Kind	John	ate.
Kind	John	fell.
Kind	Susan	ran.
Kind	Susan	ate.
Kind	Susan	fell.
Tall	Bob	ran.
Tall	Bob	ate.
Tall	Bob	fell.
Tall	John	ran.
Tall	John	ate.
Tall	John	fell.
Tall	Susan	ran.
Tall	Susan	ate.
Tall	Susan	fell.

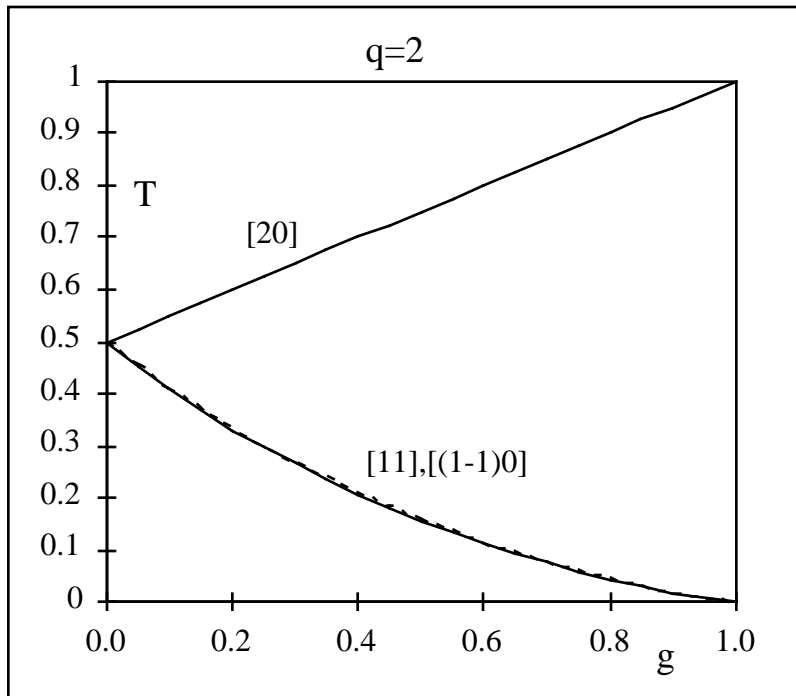


Figure 2

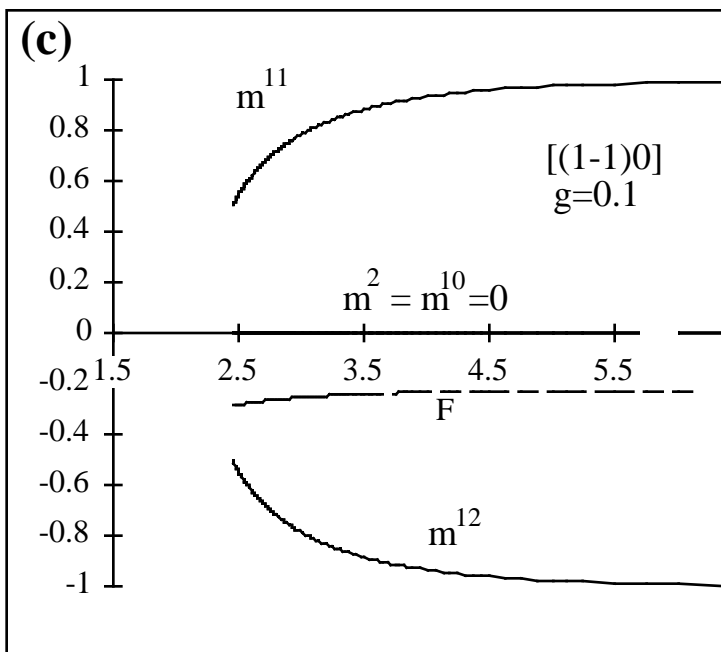
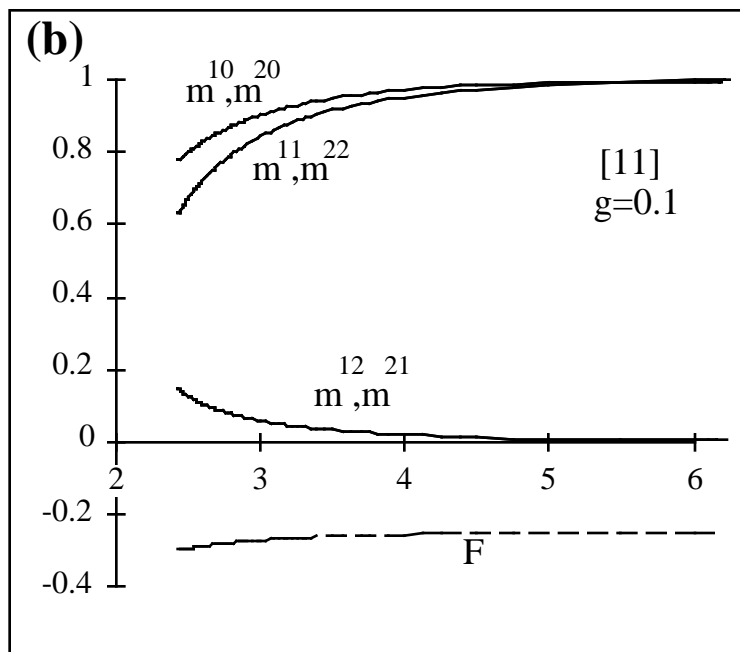
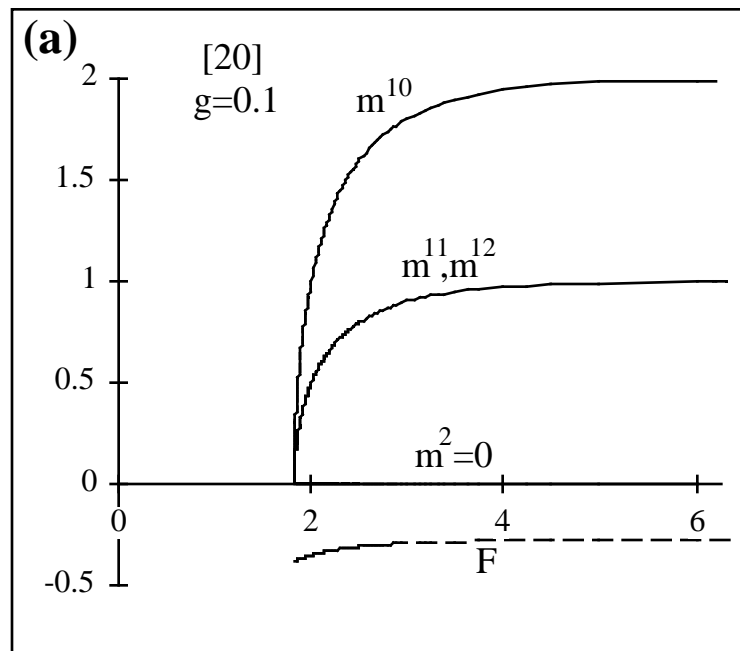


Figure 3

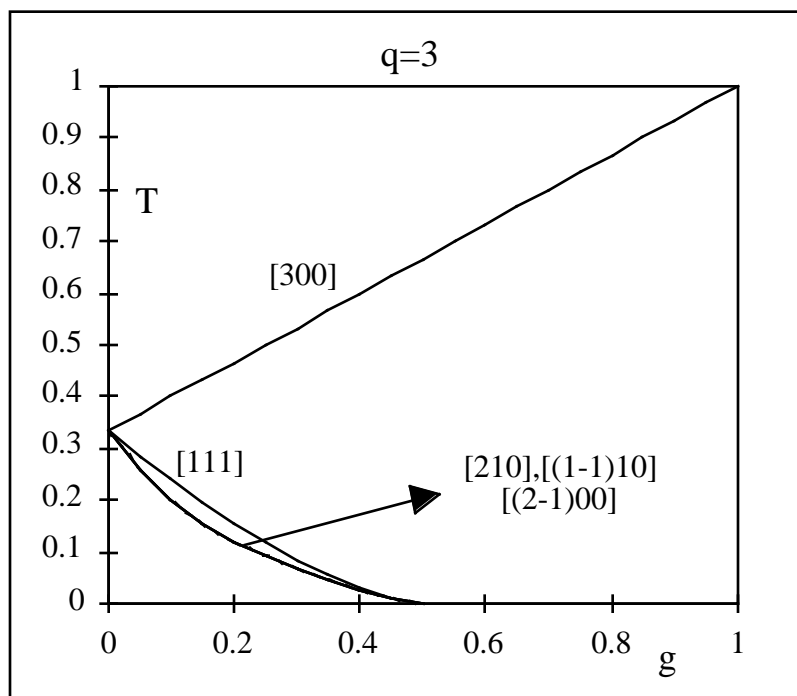


Figure 4

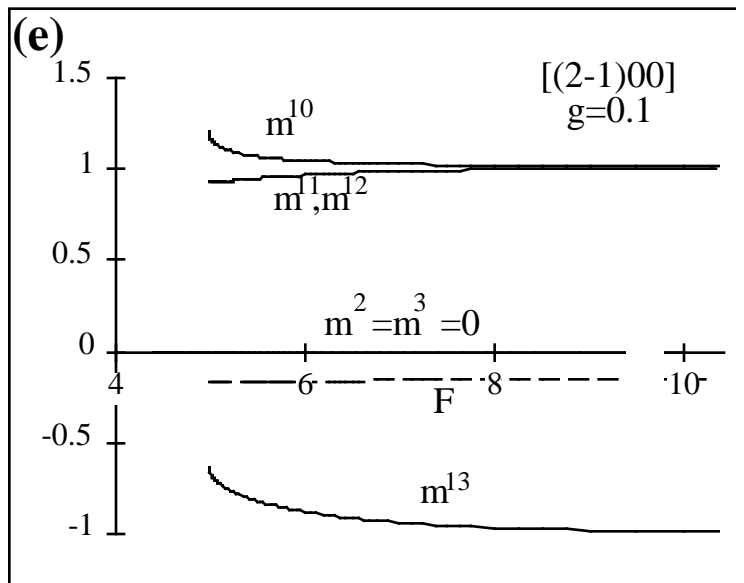
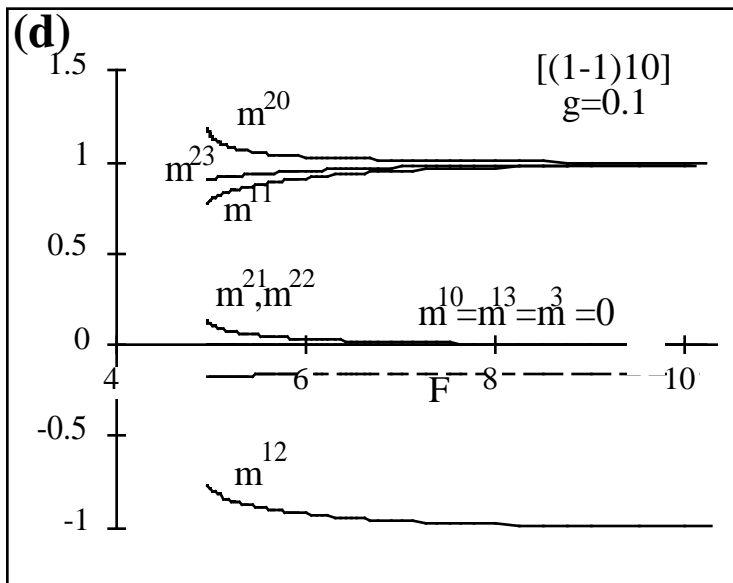
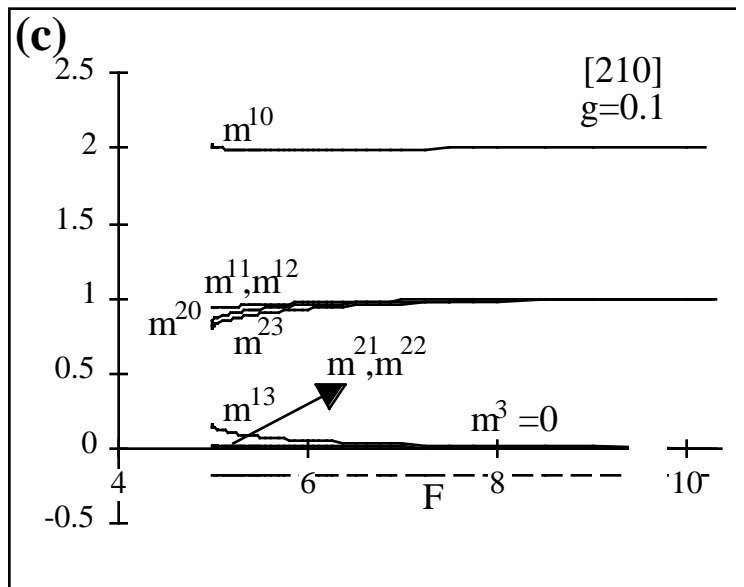
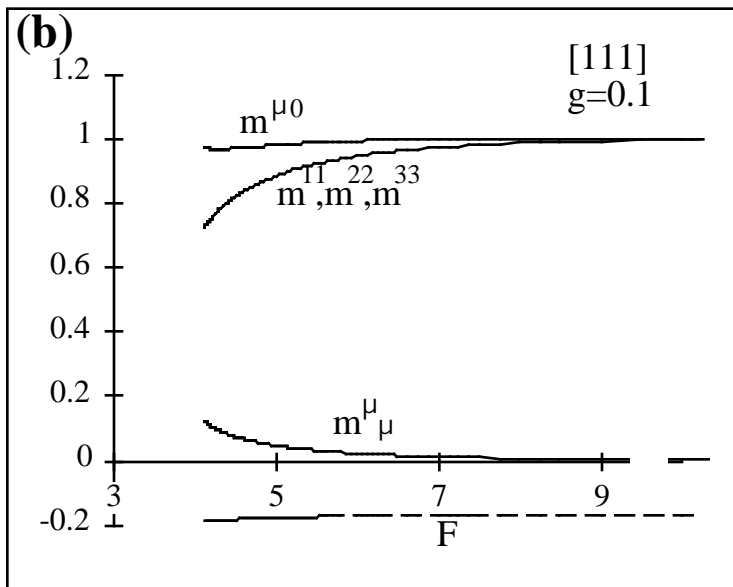
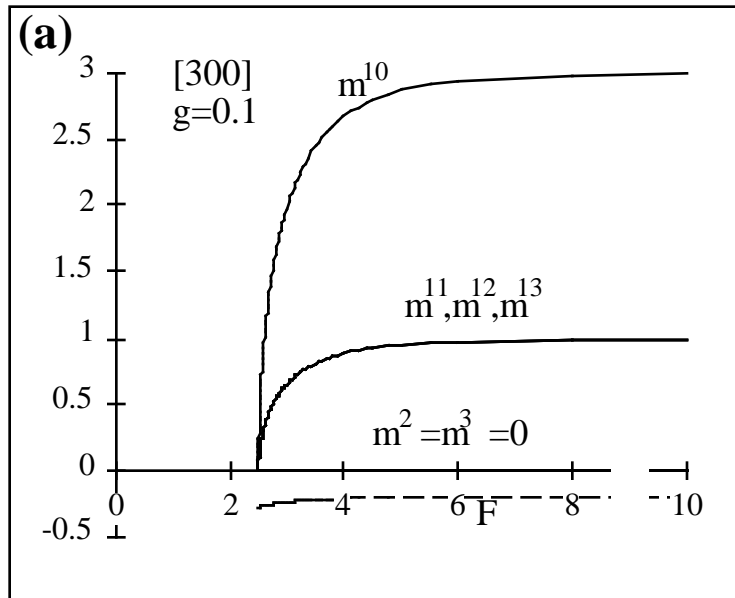


Figure 5

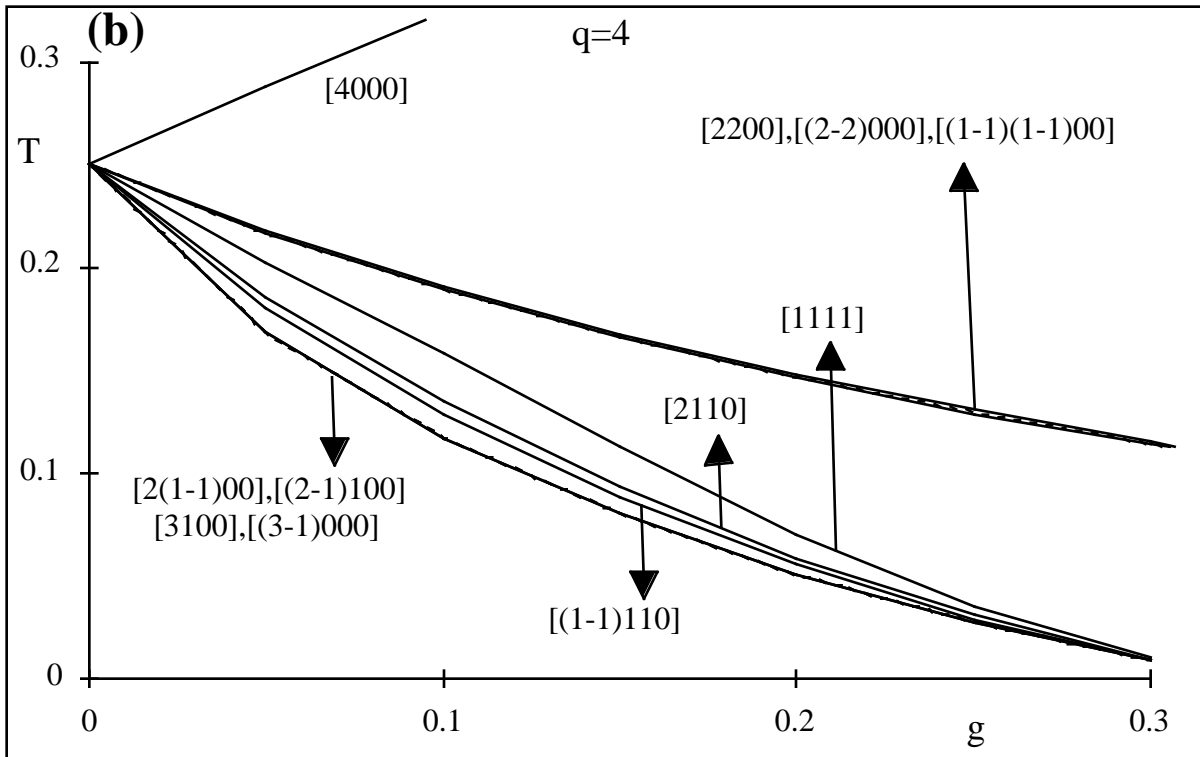
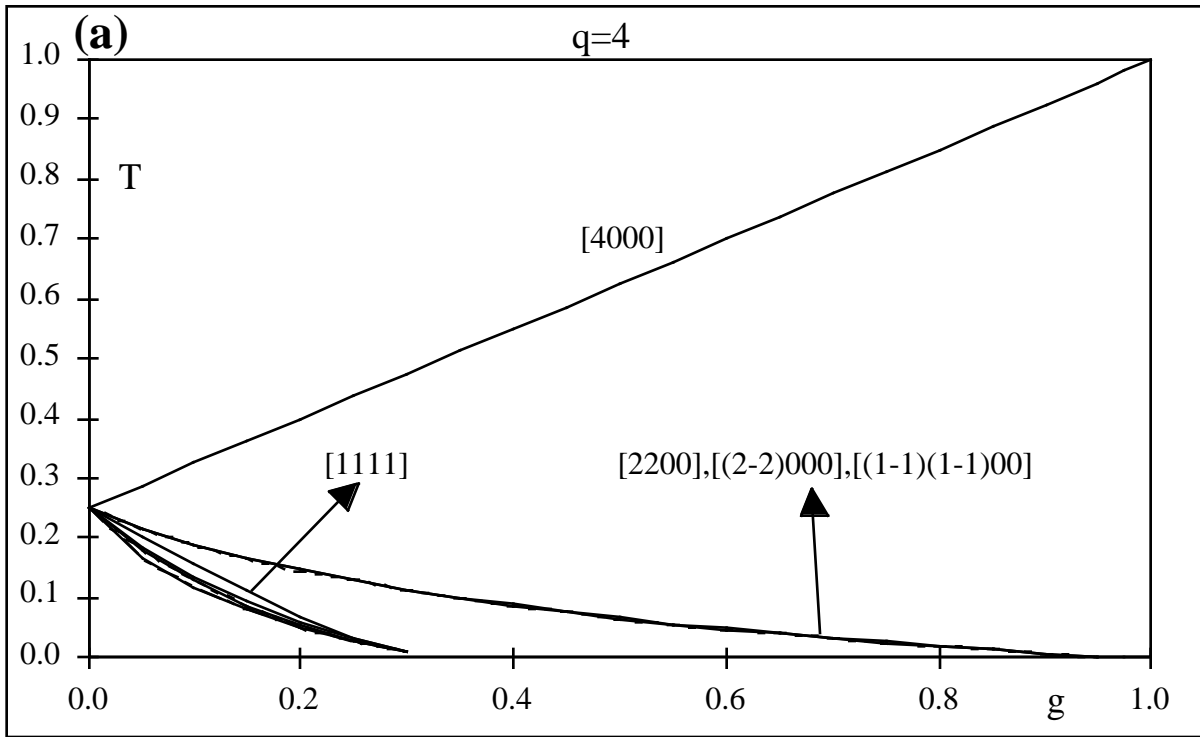


Figure 6

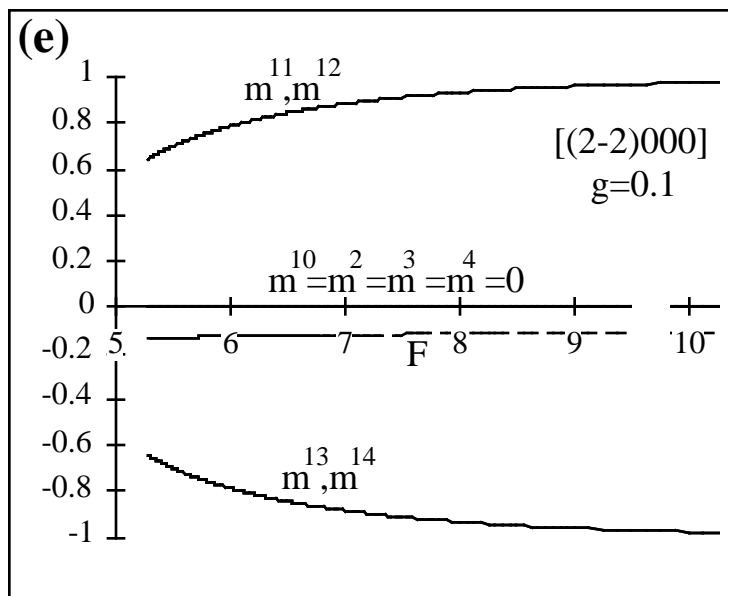
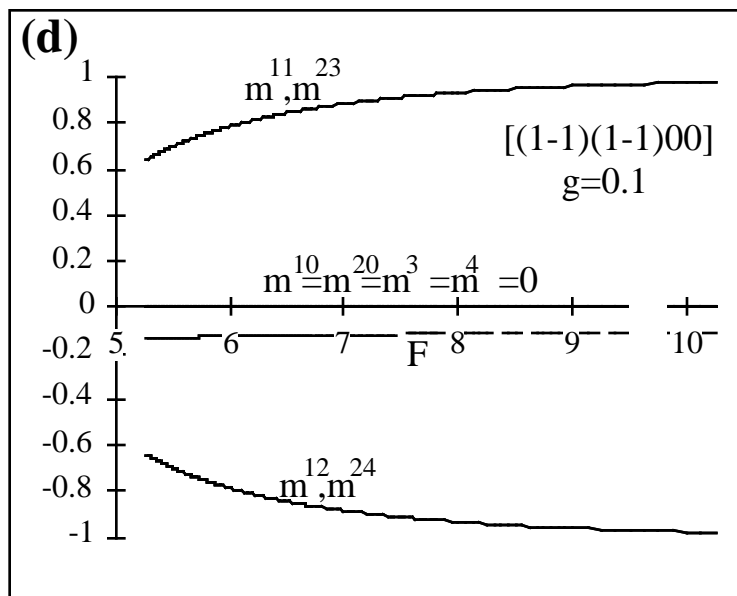
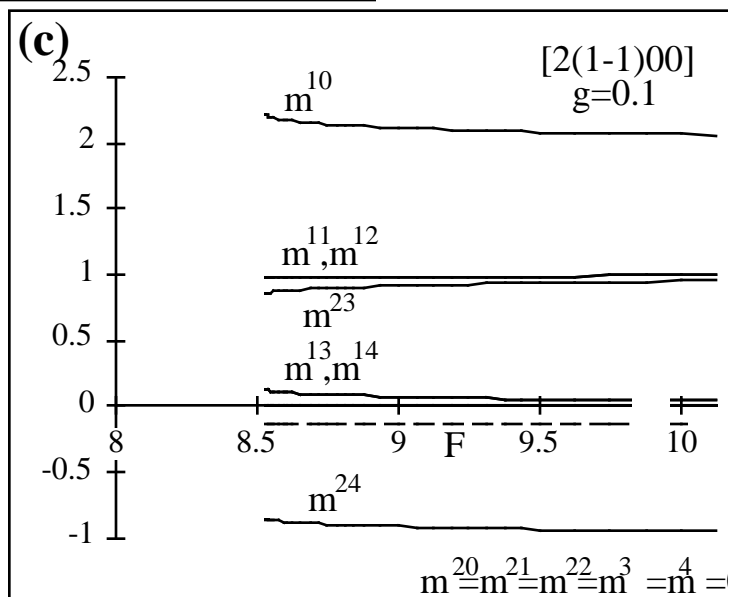
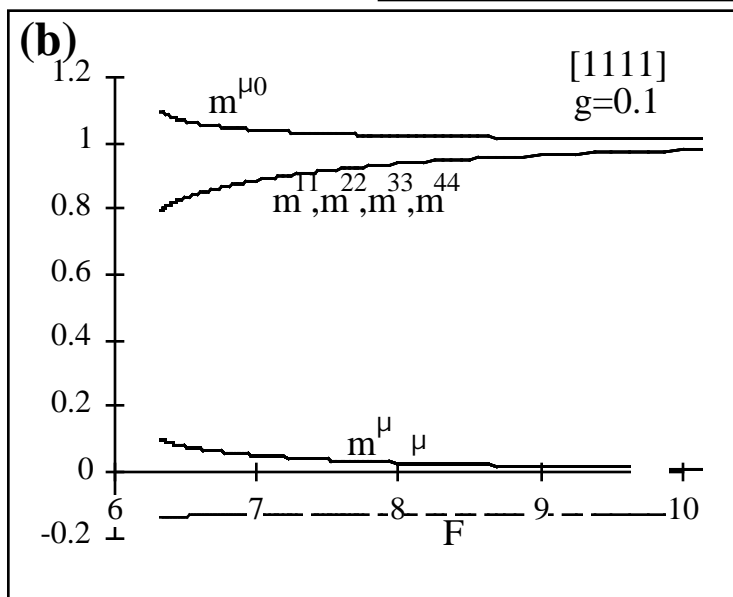
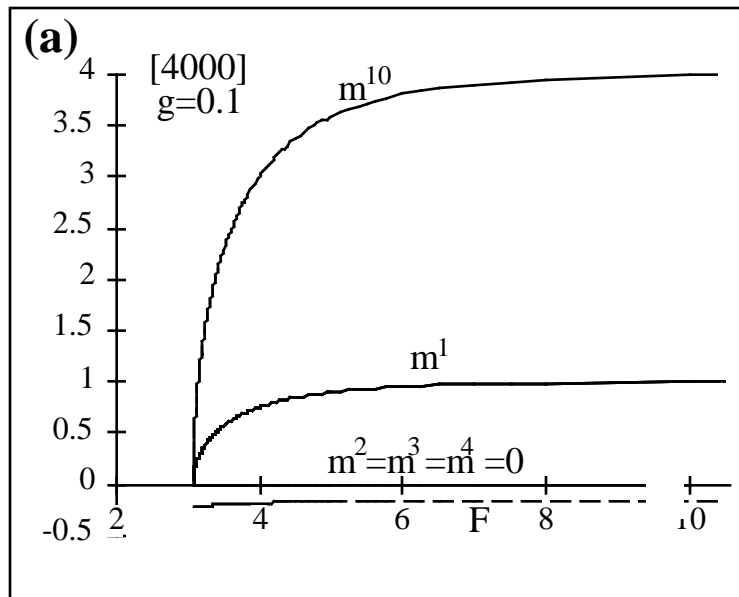


Figure 7