st unanticipated perturbations,
' (HOT) [12].

ological aspects of robustness
bustness. Dynamic aspects of
ecently [35,47,60,63,65].

**rarchies**

tion of modules and commu-
it is assumed that many com-
sibly overlapping) *modules* or
uctures has attracted substan-
ct of complex network analy-

finition of what constitutes a
work. As a working definition,
sets of vertices that are densely
nnected to other vertices outside
fluential definition of network
van, who define a score function
].
uctures is closely related to the
ntities" into distinct categories
he context of complex network
rent manners, for example, with
total number of paths between
several other possibilities. Like-
ng algorithms, modules can be
hat is, by grouping similar ver-
a complex network into smaller
oposed novel algorithms specif-
tworks, it should be noted that a
of clustering data, for example,
tperform newer methods.
can be seen from two different
t design principle of many com-
etworks shaped by evolution, it
ages in terms of their ability to
ng new functionality within the
les if often helpful to understand
or example, many technological
not be rationalized on the basis
struction of functional modules,
consisting of many lower level
be achieved.

### 3.4.3 Subgraphs and Motifs in Networks

Closely related to the detection of community structures or functional modules is the notion of *motifs* as building blocks of complex networks [31,39,43,44,59]. Providing a bridge between local vertex-related properties and global functional properties of networks, the basic idea of network motifs is that large complex networks are essentially composed of small interlinked subgraphs. Similar to the notion of sequence motifs, we can thus look for reoccurring patterns within the network, that is, small sets of vertices that exhibit an identical local topological structure. Indeed, as shown in previous studies [44], the transcriptional interaction network of *Escherichia coli* is essentially composed of repeated appearances of three highly overrepresented network motifs. Furthermore, similar studies reveal that the distribution of motifs is characteristic for certain classes of networks, that is, networks with similar overall functionality (such as communication networks or food-webs) also exhibit a highly similar motif distribution [43]. Figure 3.12 shows all possible directed subgraphs composed of three vertices. Unfortunately, the systematic enumeration of network motifs is computationally demanding.

In general, the concept of network motifs is not restricted to subgraphs with a fixed number of vertices. Rather, it allows to look for any reoccurring patterns, including more complicated topological structures, such as multi-input motifs, regulator chains, or dense overlapping regulons (DORs) [39,59]. The most intriguing aspect of network motifs, however, results from the asserted connection of local topological structure to
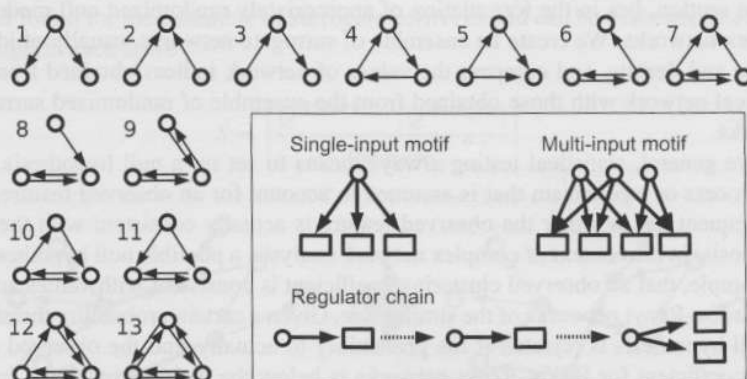


**FIGURE 3.12** All 13 possible three-vertex motifs, according to Refs. [44,59]. *Inset:* The concept of motifs can be generalized to account for specific repeated subgraphs within a complex network. Shown are possible motifs of the yeast transcription regulatory network, with circles denoting regulators and rectangles denoting gene promoters, according to Ref. [39].

dynamic function of a network motif. While such a correspondence between topology and function, that is, each motif has also a specific "hard-wired" function, was claimed in several early studies on network motifs [44,59], later results indicate that this is not the case. Though there is increasing interest in the analysis of dynamic aspects of complex networks, the precise relationship between structure and function of complex networks is still largely unclear [30,37,55,63,68]. Additional aspects in the interpretation of motif distributions, including possible pitfalls in their statistical estimation, are discussed in the next section. More details on network motifs are given in Chapter 5.

## 3.5   STATISTICAL TESTING OF NETWORK PROPERTIES

The most crucial and probably most widely underestimated aspect of complex network analysis is the statistical testing of network properties. As yet, all network indices were described as numerical quantities that can usually be straightforwardly estimated from any given network topology. However, in most applications, network indices are also associated with a biological meaning or interpretation—we seek to uncover those features of the network that are characteristic of the underlying system. Thus, for example, in the case of a metabolic network, the question is not whether the clustering coefficient takes a specific numerical value, but rather whether this value distinguishes the network from other networks of similar size, that is, whether the metabolic network can be regarded as "highly clustered." Only in the latter case, that is, if the clustering coefficient indeed deviates from what could be expected for networks of similar size, it represents a distinctive feature of the respective network. But then, how should such a deviation be detected or quantified? What values of clustering coefficient should be considered "usual" or "typical" for a network of given size?

One answer to these questions, in addition to the prototype models described in the last section, lies in the formulation of appropriately randomized null models of complex networks. We create an ensemble of surrogate networks, usually of identical size and density, and compare the values of network indices obtained from the empirical network with those obtained from the ensemble of randomized surrogate networks.

More general, statistical testing always means to set up a null hypothesis, that is, a process or mechanism that is assumed to account for an observed feature, and a subsequent test whether the observed feature is actually consistent with the null hypothesis. In the context of complex network analysis, a possible null hypothesis is, for example, that an observed clustering coefficient is consistent with values arising from Erdös–Rényi networks of the similar size. Given a certain probability threshold, the null hypothesis is rejected if the probability to actually find the observed clustering coefficient for Erdös–Rényi networks is below the defined threshold. In this case, the deviation of the clustering coefficient with respect to the null hypothesis is *significant*.

However, apart from some straightforward cases, the statistical testing of network properties holds several potential pitfalls and possible sources of misinterpretations.

In the following,
complex network

### 3.5.1   Generati

The most basic nu
but lacking any o
network of the sa
links within the n
model.

Usually more
empirical networ
degree distributio
schematically dep
$(a \rightarrow b)$ and $(c -$
provided that the
that is, such that
network has a pre
the initial empiric
features of compl
makes use of a si
motif distribution
is preserved. Clo
of complex netwo
each vertex is assi
assigned edges is
randomly chosen

In any case, a
values found for t
to a significance s

**FIGURE 3.13**  Ge
iteration, two edges
and $(c \rightarrow b)$, provid
networks, there are

In the following, we briefly outline some of the most widely used null models for complex network analysis and point out possible ambiguities in their interpretation.

### 3.5.1    Generating Networks and Null Models

The most basic null model is a network of identical size (number of vertices and edges) but lacking any other internal structure. Conceptually equivalent to an Erdös–Rényi network of the same size, such an ensemble can be constructed by randomly rewiring links within the network—as already done in the construction of the Watts–Strogatz model.

Usually more appropriate, however, is to preserve the degree distribution of the empirical network. An ensemble of randomized surrogate networks with preserved degree distribution is obtained by iteratively swapping randomly selected edges, as schematically depicted in Fig. 3.13: For a directed network, at each iteration two edges $(a \rightarrow b)$ and $(c \rightarrow d)$ are selected at random and rewired as $(a \rightarrow d)$ and $(c \rightarrow b)$, provided that the respective edges do not already exist. Repeating this sufficient times, that is, such that most edges have a statistical chance to be selected, the resulting network has a preserved degree distribution, but lacks any other internal structure of the initial empirical network. The approach can be generalized to account for other features of complex networks. For example, the analysis of network motifs [44,59] makes use of a similar approach to generate networks with a preserved three-vertex motif distribution, swapping two edges if and only if the resulting motif distribution is preserved. Closely related to network randomization is the *configuration model* of complex networks: To construct a network with a specified degree distribution, each vertex is assigned a number $k_i$ of adjacent edges, such that the total number of assigned edges is even. Subsequently, pairs of the, as yet unconnected, "stubs" are randomly chosen and connected [52,59].

In any case, a network index $Q$ of interest is subsequently compared against the values found for the ensemble of surrogate networks and can be evaluated according to a significance score

$$S = \left| \frac{Q^{\text{network}} - \langle Q^{\text{surrogate}} \rangle}{\sigma_{\text{surrogate}}} \right|, \tag{3.7}$$
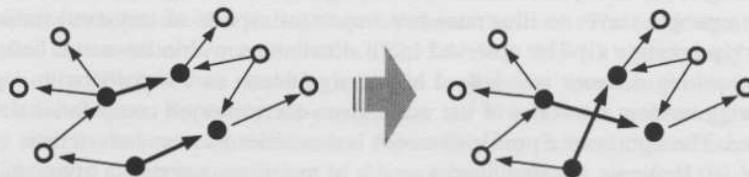


**FIGURE 3.13**    Generating random networks with preserved degree distribution: At each iteration, two edges $(a \rightarrow b)$ and $(c \rightarrow d)$ are selected at random and rewired as $(a \rightarrow d)$ and $(c \rightarrow b)$, provided that the respective edges do not already exist. Note that for undirected networks, there are two possible ways to rewire the links.

dynamic function of a network motif. While such a correspondence between topol-
ogy and function, that is, each motif has also a specific "hard-wired" function, was
claimed in several early studies on network motifs [44,59], later results indicate that
this is not the case. Though there is increasing interest in the analysis of dynamic
aspects of complex networks, the precise relationship between structure and function
of complex networks is still largely unclear [30,37,55,63,68]. Additional aspects in
the interpretation of motif distributions, including possible pitfalls in their statistical
estimation, are discussed in the next section. More details on network motifs are given
in Chapter 5.

## 3.5    STATISTICAL TESTING OF NETWORK PROPERTIES

The most crucial and probably most widely underestimated aspect of complex network
analysis is the statistical testing of network properties. As yet, all network indices were
described as numerical quantities that can usually be straightforwardly estimated from
any given network topology. However, in most applications, network indices are also
associated with a biological meaning or interpretation—we seek to uncover those
features of the network that are characteristic of the underlying system. Thus, for
example, in the case of a metabolic network, the question is not whether the clustering
coefficient takes a specific numerical value, but rather whether this value distinguishes
the network from other networks of similar size, that is, whether the metabolic network
can be regarded as "highly clustered." Only in the latter case, that is, if the clustering
coefficient indeed deviates from what could be expected for networks of similar size,
it represents a distinctive feature of the respective network. But then, how should such
a deviation be detected or quantified? What values of clustering coefficient should be
considered "usual" or "typical" for a network of given size?

One answer to these questions, in addition to the prototype models described in
the last section, lies in the formulation of appropriately randomized null models of
complex networks. We create an ensemble of surrogate networks, usually of identi-
cal size and density, and compare the values of network indices obtained from the
empirical network with those obtained from the ensemble of randomized surrogate
networks.

More general, statistical testing always means to set up a null hypothesis, that
is, a process or mechanism that is assumed to account for an observed feature, and
a subsequent test whether the observed feature is actually consistent with the null
hypothesis. In the context of complex network analysis, a possible null hypothesis is,
for example, that an observed clustering coefficient is consistent with values arising
from Erdös–Rényi networks of the similar size. Given a certain probability threshold,
the null hypothesis is rejected if the probability to actually find the observed clus-
tering coefficient for Erdös–Rényi networks is below the defined threshold. In this
case, the deviation of the clustering coefficient with respect to the null hypothesis
is *significant*.

However, apart from some straightforward cases, the statistical testing of network
properties holds several potential pitfalls and possible sources of misinterpretations.

In the following, 
complex network a

### 3.5.1   Generating

The most basic null
but lacking any oth
network of the same
links within the net
model.

Usually more ap
empirical network. 
degree distribution i
schematically depict
$(a \rightarrow b)$ and $(c \rightarrow a$
provided that the resp
that is, such that mo
network has a preser
the initial empirical 
features of complex 
makes use of a simila
motif distribution, sw
is preserved. Closely
of complex networks
each vertex is assigne
assigned edges is eve
randomly chosen and

In any case, a netw
values found for the e
to a significance score

**FIGURE 3.13**   Genera
iteration, two edges $(a$
and $(c \rightarrow b)$, provided th
networks, there are two p

respondence between topol-
"hard-wired" function, was
9], later results indicate that
in the analysis of dynamic
tween structure and function
3,68]. Additional aspects in
ble pitfalls in their statistical
s on network motifs are given


OPERTIES

ed aspect of complex network
yet, all network indices were
ightforwardly estimated from
ons, network indices are also
—we seek to uncover those
nderlying system. Thus, for
is not whether the clustering
hether this value distinguishes
hether the metabolic network
case, that is, if the clustering
for networks of similar size,
rk. But then, how should such
ustering coefficient should be
size?

rototype models described in
y randomized null models of
e networks, usually of identi-
rk indices obtained from the
nble of randomized surrogate

set up a null hypothesis, that
for an observed feature, and
ually consistent with the null
, a possible null hypothesis is,
consistent with values arising
certain probability threshold,
tually find the observed clus-
the defined threshold. In this
respect to the null hypothesis

e statistical testing of network
sources of misinterpretations.

In the following, we briefly outline some of the most widely used null models for complex network analysis and point out possible ambiguities in their interpretation.

### 3.5.1   Generating Networks and Null Models

The most basic null model is a network of identical size (number of vertices and edges) but lacking any other internal structure. Conceptually equivalent to an Erdös–Rényi network of the same size, such an ensemble can be constructed by randomly rewiring links within the network—as already done in the construction of the Watts–Strogatz model.

Usually more appropriate, however, is to preserve the degree distribution of the empirical network. An ensemble of randomized surrogate networks with preserved degree distribution is obtained by iteratively swapping randomly selected edges, as schematically depicted in Fig. 3.13: For a directed network, at each iteration two edges $(a \rightarrow b)$ and $(c \rightarrow d)$ are selected at random and rewired as $(a \rightarrow d)$ and $(c \rightarrow b)$, provided that the respective edges do not already exist. Repeating this sufficient times, that is, such that most edges have a statistical chance to be selected, the resulting network has a preserved degree distribution, but lacks any other internal structure of the initial empirical network. The approach can be generalized to account for other features of complex networks. For example, the analysis of network motifs [44,59] makes use of a similar approach to generate networks with a preserved three-vertex motif distribution, swapping two edges if and only if the resulting motif distribution is preserved. Closely related to network randomization is the *configuration model* of complex networks: To construct a network with a specified degree distribution, each vertex is assigned a number $k_i$ of adjacent edges, such that the total number of assigned edges is even. Subsequently, pairs of the, as yet unconnected, "stubs" are randomly chosen and connected [52,59].

In any case, a network index $Q$ of interest is subsequently compared against the values found for the ensemble of surrogate networks and can be evaluated according to a significance score

$$S = \left| \frac{Q^{\text{network}} - \langle Q^{\text{surrogate}} \rangle}{\sigma_{\text{surrogate}}} \right|, \qquad (3.7)$$
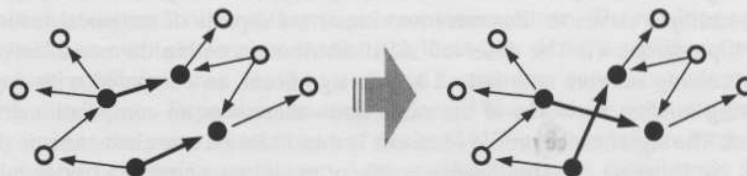


**FIGURE 3.13**   Generating random networks with preserved degree distribution: At each iteration, two edges $(a \rightarrow b)$ and $(c \rightarrow d)$ are selected at random and rewired as $(a \rightarrow d)$ and $(c \rightarrow b)$, provided that the respective edges do not already exist. Note that for undirected networks, there are two possible ways to rewire the links.

where $\langle Q^{\text{surrogate}} \rangle$ denotes the average found within the ensemble of surrogate networks and $\sigma_{\text{surrogate}}$ the standard deviation. Unfortunately, the construction of randomized networks that preserve features other than the degree or motif distribution is far from straightforward.

### 3.5.2   The Conceptualization of Cellular Networks

The most difficult problem of complex network analysis is often the choice of an appropriate null model or null hypothesis. By definition, all statistical tests rely on the definition of a null model, and any notion of "significance" is defined with respect to this null model only. Thus, if the null model is erroneous or trivial, so will be any result obtained from an evaluation of the significance of network properties against this null model. In particular, the choice of the null model is additionally complicated by the fact that networks are usually abstract representations of more complex biological processes. While this "reduction in complexity" is often necessary to make biological questions mathematically tractable, it also holds the temptation to neglect properties of the underlying system—resulting in an erroneous or misleading interpretation of network properties.

An illustrative example of such a case was discussed in the context of a recent study of motif distributions in complex networks [43]. Therein, the significance profile of small subgraphs (motifs) within a neuronal network was evaluated and compared with simple degree-preserving randomized networks (see also Chapter 5). The study concluded that the neural information processing networks exhibits a highly characteristic significance profile for its motif distribution, suggesting evolutionary mechanisms that result in key circuit elements to perform specific tasks. However, as pointed out later [5], a neural network, that is, a network of neurons connected by synapses, is not just a network of vertices connected by edges. Rather, neurons have a spatial position and a tendency to form local clusters, hence neighboring neurons have greater chance of forming connections than distant neurons. As the spatial properties are not reflected in the null hypothesis, the statistical test misclassifies a completely random but spatially clustered network as one that is nonrandom and exhibits significant network motifs. Indeed, a simple toy model that preserves the spatial position of neurons and connects neurons preferentially to nearby neurons is able to reproduce an almost identical significance profile for the motif distribution, without the need to invoke any evolutionary mechanisms to select for specific functional tasks [5].

The example serves to illustrates two important aspects of statistical testing of network properties: (i) The observed motif distribution within the neural information processing network was indeed highly significant, as compared with degree-preserving random networks of the same size—there was no computational error involved. The significance profile of motifs is thus indeed a true characteristic of the system. (ii) However, the significance profile of motifs tests against a trivial null hypothesis: the assumption of a completely random network. Strictly speaking, rejecting this null hypothesis only trivially proves that neural networks are not random. However, only little can be learned about the possible biological function of network motifs from the sole fact that their distribution is not random. As exemplified here for the

the ensemble of surrogate net-
ately, the construction of ran-
he degree or motif distribution

:s

lysis is often the choice of an
tion, all statistical tests rely on
ficance" is defined with respect
oneous or trivial, so will be any
f network properties against this
l is additionally complicated by
ons of more complex biological
en necessary to make biological
temptation to neglect properties
 or misleading interpretation of

ssed in the context of a recent
]. Therein, the significance pro-
etwork was evaluated and com-
works (see also Chapter 5). The
sing networks exhibits a highly
bution, suggesting evolutionary
erform specific tasks. However,
etwork of neurons connected by
by edges. Rather, neurons have a
hence neighboring neurons have
eurons. As the spatial properties
l test misclassifies a completely
 nonrandom and exhibits signifi-
 preserves the spatial position of
rby neurons is able to reproduce
 distribution, without the need to
ecific functional tasks [5].
 aspects of statistical testing of
tion within the neural informa-
cant, as compared with degree-
ere was no computational error
ndeed a true characteristic of the
tifs tests against a trivial null hy-
work. Strictly speaking, rejecting
 networks are not random. How-
ogical function of network motifs
om. As exemplified here for the

neural network, the significant deviation from random networks is most likely a simple
and straightforward consequence of the (neglected) spatial structure of the system.

Thus, as a general rule, network properties that are found to be significant with re-
spect to simple randomized networks must not necessarily be important with respect
to function. Even though the statistical estimation of significance might be techni-
cally correct, significance here only implies deviations from randomness—which is
often a trivial consequence of the underlying process or system itself. Two important
classes of complex networks where the construction of the networks itself implies
deviations from randomness, and thus implies highly significant network properties,
are discussed below.