# Small Cycles in Small Worlds

Petra M. Gleiss,[1] Peter F. Stadler[1],[2] Andreas Wagner[2],[3] and David A. Fell[4]

[1]*Institut für Theoretische Chemie, Universität Wien Währingerstraße 17, A-1090 Wien, Austria*
[2]*The Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe NM 87501, USA*
[3]*Dept. of Biology, University of New Mexico, 167A Castetter Hall, Albuquerque, NM-817131-1091,*
[4]*School of Biological & Molecular Sciences Oxford Brookes University, Oxford OX3 0BP, U.K.*

We characterize the distributions of short cycles in a large metabolic network previously shown to have small world characteristics and a power law degree distribution. Compared with three classes of random networks, including Erdős-Rényi random graphs and synthetic small world networks of the same connectivity, the metabolic network has a particularly large number of triangles and a deficit in large cycles. Short cycles reduce the length of detours when a connection is clipped, so we propose that long cycles in metabolism may have been selected against in order to shorten transition times and reduce the likelihood of oscillations in response to external perturbations.

Systems as diverse as the Western US power grid, metabolic networks of a cell, or the World Wide Web are well described as graphs with characteristic topology. Small world networks have received considerable attention since the seminal paper by Watts and Strogatz [1].

Most of the existing literature discusses small world networks in terms of the average path length between two vertices [2] or of the network's clustering coefficient [3, 4] which measures how close the neighborhood of a each vertex comes on average to being a complete subgraph (clique) [1]. Barabási *et al.* [5, 6] focussed on the degree distributions, finding a power law in a suite of real world examples including the world wide web or the US power-grid. Recent work on the spread of epidemics on a small world network [7] emphasizes the importance of "far-reaching" edges. The idea is that clipping a far edge will force a (relatively) long detour in the network. Hence it is these edges that are responsible for the small diameter of the graph $G$.

Let us look at detours in graphs in more systematic way. Throughout this paper we will represent a network as a simple (unweighted, undirected) graph $G(V, E)$ with vertex set $V$ and edge set $E$. A *cycle* in $G$ is a closed path which meets each of its vertices and edges exactly once. The length of a cycle $C$, i.e., the number of its vertices or edges, is denoted by $|C|$. With each edge $e \in E$ we can associate the set $\mathcal{S}(e)$ containing the shortest cycles in $G$ that go through $e$. It is easily verified that a far edge in the sense of [7] is an edge that is not contained in a triangle. In other words, $e$ is a far edge if and only if $\mathcal{S}(e)$ does not contain a triangle. The cycles $C \in \mathcal{S}(e)$ determine the shortest detours (which have length $|C|-1$) when $e$ is removed from the graph.

It seems natural to consider the set $\mathcal{S}(G) = \bigcup_{e \in E} \mathcal{S}(e)$ of shortest cycles of all edges in $G$ and to study e.g. their length distribution. However, as the example in Fig.1 shows, the shortest cycles $\mathcal{S}(G)$ do not convey the complete information about the graph. Additional cycles appear to be relevant, such as the hexagon in Figure 1.
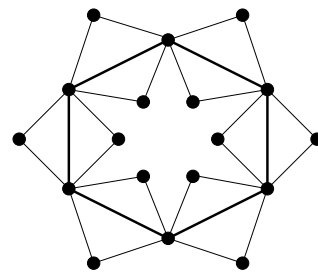


FIG. 1: $\mathcal{S}(G)$ consists of the twelve triangle only. The hexagon (bold edges), however, is obviously crucial for the network structure.

In order to extend $\mathcal{S}(G)$ to a more complete collection of cycles we need some more information on the cycle structure of graphs. Recall that the set of all subsets of $E$ forms an $|E|$-dimensional vector space over $\{0,1\}$ (with addition and multiplication modulo 2). Vector addition in this *edge space* is given by symmetric difference $X \oplus Y = (X \cup Y) \setminus (X \cap Y)$. The *cycle space* $\mathcal{C}$ consisting of all cycles and edge-disjoint unions of cycles in $G$ is a particularly important subspace of the edge space [8]. Its dimension is the *cyclomatic number* $\nu(G) = |E| - |V| + c(G)$, where $c(G)$ is the number of connected components of $G$. The length $\ell(\mathcal{B})$ of a basis $\mathcal{B}$ of the cycle space (*cycle basis* for short) is the sum of the lengths of its cycles: $\ell(\mathcal{B}) = \sum_{C \in \mathcal{B}} |C|$. A *minimum cycle basis* (MCB) is a cycle basis with minimum length. MCBs have the property that their longest cycle is at most as long as the longest cycle of any basis of $\mathcal{C}$ [9]. A MCB therefore contains the salient information about the cycle structure of a graph in its most compressed form. Most graphs, however, do not have a unique MCB. On the other hand, the distribution of cycle lengths is the same in all MCBs of a given graph [10]. The way to avoid ambiguities is to consider the union of all minimum cycles bases, also known as the set $\mathcal{R}(G)$ of *relevant cy-*

*cles.* The term "relevant" is justified by two important properties of $\mathcal{R}(G)$: (i) a cycle is relevant if and only if it cannot be written as an $\oplus$-sum of shorter cycles [11], and (ii) the shortest cycles through an edge are relevant, i.e., $\mathcal{S}(G) \subseteq \mathcal{R}(G)$ [10, 12]. Consequently, the composition of $\mathcal{R}(G)$ in terms of number and length distribution of cycles is an important characteristic of a graph. The numerical studies below make use of Vismara's [11] algorithm for computing $\mathcal{R}(G)$, which is based on Horton's MCB algorithm [13].

The most common model of graph evolution, introduced by Erdős and Rényi [14], assumes a fixed number $n = |V|$ of vertices and assigns edges independently with a certain probability $p$ [15]. In many cases ER random graphs turn out the be quite different from a network of interest. The Watts-Strogatz [1] model of small world networks starts with a deterministic graph, usually a circular arrangement of vertices in which each vertex is connected to $k$ nearest neighbors on each side. Then edges are "rewired" (in the original version) or added [2, 16] with probability $p$. We shall consider the latter model for $k = 1$, denoted SW1 below, which corresponds to adding random edges to a Hamiltonian cycle. Both ER and SW1 graphs exhibit an approximately Gaussian degree distribution.

In many real networks, however, the degree distribution follows a power law. Barabási *et al.* [5, 6] show that the scale invariant behavior of the degree distributions can be explained in terms of simple graph evolution model (AB model): Starting from a small core graph, at each time step a vertex is added together with $m$ edges that are connected to each previously present vertex $k$ with probability $\Pi(k) = d(k)/\sum_j d(j)$, where $d(j)$ is the degree of vertex $j$. In this contribution we will focus mostly on the AB model instead of Watt's original construction, because we will apply the analysis of the cycle structure to an empirical network for which a power-law like degree distribution has been established. This network is the system of all chemical reactions required for the synthesis of small-molecule building blocks and energy in the bacterium *Escherichia coli*. Its structure described in ref. [17]. Such chemical reaction networks are often referred to as metabolic networks.

It is clear that all triangles in a graph are relevant, since a triangle is for sure a shortest cycle through each edge. Hence $|\mathcal{R}(G)| \geq \Delta$, where $\Delta$ denotes the number of triangles in $G$. We expect $\langle \Delta \rangle_{\mathrm{ER}} = \binom{n}{3} p^3$ triangles in an ER random graph with edge-drawing probability $p$. For the SW1 graphs we obtain a similar expression:

$$\langle \Delta \rangle_{\mathrm{SW1}} = np + n(n-4)p^2 + \frac{1}{6}n(n^2 - 9n + 20)p^3 . \quad (1)$$

The MCB will therefore consist almost exclusively of triangles if $\Delta \gg \nu(G)$. The average vertex degree is $d = 2|E|/n = p(n-1)$ for ER and $d = 2 + p(n-3)$ for SW1, resp. Assuming that $n$ is large we expect to find only triangles in $\mathcal{R}(G)$ for $d \gg \sqrt{3n}$. Numerical simulations show that this is indeed the case, Fig 2. In this
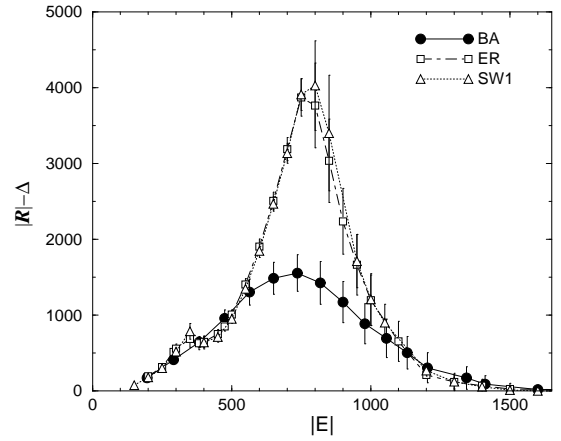


FIG. 2: Relevant non-triangles in ER ($\square$), SW1 ($\triangle$), and AB ($\bullet$) random graphs with $n = 100$.

regime, we have $|\mathcal{R}(G)| \sim d^3/6$, and the graph contains no far edges. Not surprisingly, there is little difference between SW1 and ER random graphs for large $n$.

Since the AB model is based on a fixed vertex degree $d$, it should be compared to random graph models with given vertex degree $d$, not with given edge drawing probabilities $p$. We have an asymptotically constant number of triangles for both ER and SW1: $\Delta_{\mathrm{ER}} \to d^3/6$ and $\Delta_{\mathrm{SW1}} \to d^3/6 - d + 2/3$, resp. Note that as a consequence the clustering coefficient vanishes asymptotically. In SW networks with *a priori* connectivity $k > 1$ we find of course a number of triangles that grows at least linearly with $n$, since the initial ($p = 0$) networks already contains $(k-1)n$ triangles. The clustering coefficient stays finite for large $n$ [18].

The large vertex degree of the "early" vertices in the AB model suggests that there should be many more triangles than in ER or SW1 models. The expected degree of vertex $s$ at "time" $t$ is known [19]: $d(s|t) = m[\sqrt{t/s} - 1]$. The probability of an edge between $s$ and $t$, $t > s$, is therefore $p_{st} = md(s|t-1)/2(t-1)m$, where $2(t-1)m$ is the sum of the vertex degrees at "time" $t-1$. Thus $\langle \Delta \rangle = \sum_{r<s<t} p_{rs}p_{st}p_{rt}$. This can be approximated by

$$\langle \Delta \rangle \approx \frac{m^3}{8} \int_{1<r<s<t}^{n} (1/st^2) \left( \sqrt{\frac{s}{r}} - 1 \right) \left( \sqrt{\frac{t}{r}} - 1 \right) \left( \sqrt{\frac{t}{s}} - 1 \right)$$

$$\sim Cm^3 \ln^3 n + \mathcal{O}(\ln^2 n) \quad (2)$$

Fig. 3 shows $\Delta$ for typical AB-random graphs with $m = 2, \ldots, 8$ as a function of "time'. The behavior of $\Delta$ in a individual growing network is well represented by equ.(2).

The number $|\mathcal{R}| - \Delta$ of non-trivial relevant cycles has its maximum around $|E| \approx 0.74n^{3/2}$ independent of the model. The scaling of their number is consistent with $|\mathcal{R}| - \Delta \sim Cn^{5/2}$, where the constant $C \approx 0.036$ is the same for ER and SW1 random graphs and $C \approx 0.016$ for the AB models. For small vertex degrees, $d \ll |V|^{1/2}$ we find $\mathcal{R}(G) \approx \nu(G)$, i.e., the MCB is (almost) unique.
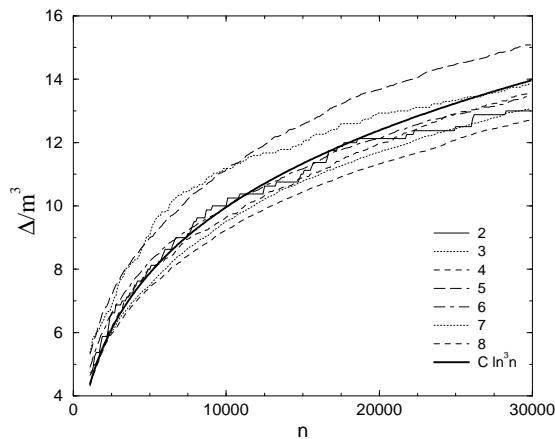
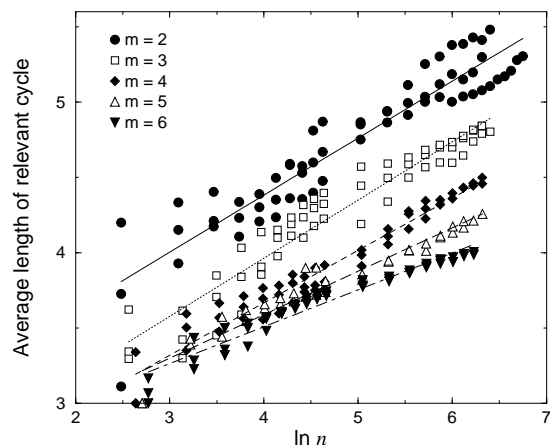FIG. 3: Triangles in AB models with different values of $m$.



FIG. 4: Mean length of a relevant cycle in AB networks.

TABLE I: Cycle Structure of Metabolic Networks.

| Model | $|C|$ | 3 | 4 | 5 | 6 | 7 | 8 | 9 | $\sum$ |
|---|---|---|---|---|---|---|---|---|---|
| Ecoli1 | MCB | 282 | 51 | 19 | 20 | 3 | 5 | 1 | 381 |
| | $\mathcal{R}$ | 379 | 114 | 90 | 83 | 5 | 36 | 16 | 723 |
| | $\mathcal{S}$ | 379 | 56 | 24 | 42 | 2 | 14 | 16 | 533 |
| AB | MCB | 78 | 158 | 124 | 20 | 0.4 | 0.01 | 0 | 380 |
| | $\mathcal{R}$ | 81 | 285 | 527 | 161 | 5.5 | 0.4 | 0 | 1060 |
| | $\mathcal{S}$ | 81 | 273 | 414 | 144 | 5.5 | 0.4 | 0 | 918 |
| ER | MCB | 18 | 58 | 163 | 131 | 11 | 0.4 | 0 | 381 |
| | $\mathcal{R}$ | 18 | 61 | 212 | 528 | 82 | 3.2 | 0 | 904 |
| | $\mathcal{S}$ | 18 | 61 | 205 | 311 | 68 | 3.2 | 0 | 666 |
| SW1 | MCB | 15 | 46 | 131 | 167 | 21 | 1.1 | 0.03 | 381 |
| | $\mathcal{R}$ | 15 | 48 | 157 | 427 | 151 | 7.1 | 0.2 | 805 |
| | $\mathcal{S}$ | 15 | 48 | 155 | 301 | 108 | 6.5 | 0.2 | 634 |

substrate concentration. Thus, perturbations may travel backwards even from irreversible reactions. A similar argument for considering undirected graphs can be derived from metabolic control theory [20].

The key ingredient of MFA is the *stoichiometric matrix* $\mathbf{S}$. Its entries are the stoichiometric coefficients $s_{kr}$, i.e., the number of molecules of species $k$ produced ($s_{kr} > 0$) or consumed ($s_{kr} < 0$) in each reaction $r$. Reversible reactions are entered as two separate reactions in most references. In general, additional "pseudo-reactions" are added to describe the interface of the metabolic reaction network with its environment. Stationary flux vectors $\vec{f}$ in the network satisfy $\mathbf{S}f = \vec{o}$ and $f_r \geq 0$ for each reaction $r$, see e.g. [21, 22, 23, 24, 25, 26]. It is not hard to see that if all reactions are mono-molecular, then $\mathbf{S}$ is the incidence matrix of a directed graph; The stationary flux vectors span the cycle space of this graph. The close connection between the cycle space of a directed graph and its underlying undirected graph [27] allows us to use the relevant cycles of the substrate graph $\Sigma$ to describe the structure of the metabolic network in a way complementary to that provided by MFA.

For our analysis of metabolic graphs, we use the substrate graph of the Ecoli1 core metabolism, a set of chemical reactions representing the central routes of energy metabolism and small-molecule building block synthesis. Similar to [17], we omit the following substrates from the graph: $CO_2$, $NH_3$, $SO_4$, AMP,ADP, and ATP, their deoxy-derivatives, both the oxidized and reduced form of thioredoxine, organic phosphate and pyrophosphate. The resulting graph has $n = 272$ vertices and $|E| = 652$ edges. It is analyzed below.

Table I shows that the three random models AB, SW1, and ER agree at least qualitatively with each other. The AB random graphs exhibit a much broader distribution of cycle sizes (not shown) than the ER and SW1 models. As a consequence, the average cycle numbers for ER and SW1 have statistical uncertainty of about 2%, while the uncertainty of the AB values is 5 to 10 times higher. Note that ER and SW1 have a similar number of relevant cycles, but the cycles are slightly longer in
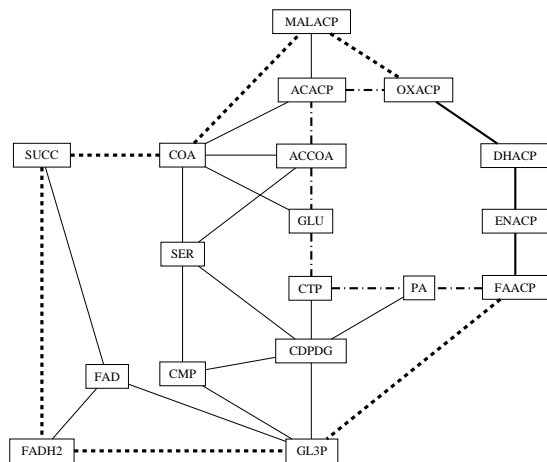
The cyclomatic number of a AB random graph is $\nu(G) \sim (m/2 - 1)n$; Hence eq.(2) implies that almost all relevant cycles must be long. Fig. 4 shows that the average length of a relevant cycle grows logarithmically with $n$. Not surprisingly, the slopes decrease with $m$.

Let us now turn to an example of metabolic networks. Because it is germane to their functional analysis, we first point out a nexus between graph representations of metabolic network, and metabolic flux analysis (MFA), the most generic framework to analyze the biological function of metabolic networks. A graph representation of metabolic networks was introduced as a *substrate graph* $\Sigma$ in [17]. Its vertices are the molecular compounds (substrates); two substrates $k$ and $l$ are adjacent in $\Sigma$ if they participate in the same reaction $r$. Substrate graphs are undirected because directed graphs would not properly represent the propagation of perturbations: even for irreversible reactions the product concentration may affect the the reaction rate, for instance by product occupancy of the enzyme's active site; this in turn affects the

FIG. 5: The subgraph of `Ecoli` spanned by the relevant cycles of length 9. Two of these long cycles are highlighted. The edges shown in bold are part of each of the 16 relevant 9-cycles.

SW1. Two features distinguish the metabolic network `Ecoli1` from all random networks: (1) The number $\Delta$ of triangles is almost 10 times larger than expected. This can be explained at least in part as a consequence of the substrate graph representation: multi-molecular reactions translate to cliques and hence a large number of triangles. The ratio $282/379 \approx 0.744$ indicates that in fact almost all triangles are contained in 4-cliques, since in each 4-clique we have three triangles that belong to a particular MCB, while the fourth face of the tetrahedron is their $\oplus$-sum [28]. (2) There is a much smaller number of relevant pentagons and hexagons, which results in an overall somewhat reduced number of relevant cycles: 723 compared to about 1060 (AB), 904 (ER), and 805 (SW1).

Strictly speaking, we do not know the biological significance of this relative paucity of longer cycles. However, we would like to venture a speculation. Organisms are constantly exposed to environmental fluctuations requiring transitions in metabolic states. That is, a metabolic network needs to produce different outputs depending on the environment. Environments may vary rapidly, requiring rapid transition between metabolic states. Possibly, networks with long cycles have longer transition times, because environmental perturbations may lead to prolonged oscillations in such networks. The dynamical system representation of metabolic networks required to test this idea rigorously lies beyond the scope of this article.

The longest relevant cycles in a metabolic networks are of particular interest since they reflect parts of the network that cannot easily be replaced by alternative routes. In Fig. 5 we show the largest such cycle in `Ecoli1`. We emphasize that the cycles in our analysis represent routes for transmission of perturbations, but not necessarily of mass, as it is commonly considered in MFA. This is apparent from Fig.5 , which does not correspond to a path-way from a biochemical chart, but links serval pathways together.

## REFERENCES

[1] D. J. Watts and H. S. Strogatz, Nature **393**, 440 (1998).

[2] M. E. J. Newman, C. Moore, and D. J. Watts, Phys. Rev. Lett. **84**, 3201 (2000).

[3] H. Herzel, Fractals **6**, 301 (1998).

[4] A. Barrat and M. Weigt, Europ. Phys. J. B **13**, 547 (2000).

[5] A.-L. Barabási and R. Albert, Science **286**, 509 (1999).

[6] A.-L. Barabási, R. Albert, and H. Jeong, Physica A **173-187**, 272 (1999).

[7] S. A. Pandit and R. E. Amritkar, Phys. Rev. E **60**, R1119 (1999), chao-dyn/9901017.

[8] W.-K. Chen, SIAM J. Appl. Math. **20**, 525 (1971).

[9] D. M. Chickering, D. Geiger, and D. Heckerman, Inform. Processing Let. **54**, 55 (1994).

[10] G. F. Stepanec, Uspekhi Mat. Nauk. 2 **19**, 171 (1964), (Russian).

[11] P. Vismara, Electr. J. Comb. **4**, 73 (1997).

[12] A. A. Zykov, *Theory of Finite Graphs* (Nauka, Novosibirsk, USSR, 1969), (Russian).

[13] J. D. Horton, SIAM J. Comput. **16**, 359 (1987).

[14] P. Erdős and A. Rényi, Publ. Math. Inst. Hung. Acad. Sci., Ser. A **5**, 17 (1960).

[15] B. Bollobás, *Random Graphs* (Academic Press, London UK, 1985).

[16] M. E. J. Newman and D. J. Watts, Phys. Lett. A **263**, 341 (1999).

[17] A. Wagner and D. A. Fell, *The small world inside large metabolic networks*, Tech. Rep. 00-07-041, Santa Fe Institute (2000).

[18] D. J. Watts, *Small Worlds* (Princeton University Press, Princeton NJ, 1999).

[19] S. N. Dorogovtsev, J. F. F. Mendes, and A. N. Samukhin (2000), cond-mat/0004434.

[20] A. K. Sen, Biochem. J. **275**, 253 (1991).

[21] B. L. Clarke, Cell Biophys. **12**, 237 (1988).

[22] R. Heinrich and S. Schuster, *The Regulation of Cellular Systems* (Chapman & Hall, New York, 1996).

[23] D. A. Fell, *Understanding the Control of Metabolism* (Portland Press, London, 1997).

[24] C. H. Schilling, D. Letscher, and B. Ø. Palsson, J. Theor. Biol. **203**, 229 (2000).

[25] J. D. Edwards and B. Ø. Palsson, Proc. Natl. Acad. Sci. USA **97**, 5528 (2000).

[26] S. Schuster, D. A. Fell, and T. Dandekar, Nature Biotechnol. **18**, 326 (2000).

[27] B. Bollobás, *Modern Graph Theory* (Springer, New York, 1998).

[28] P. M. Gleiss, J. Leydold, and P. F. Stadler, Elec. J. Comb. **7**, R16 [16pages] (2000).