

A Detailed Study of the Attachment Strategies of New Autonomous Systems in the AS Connectivity Graph

Ingrid Daubechies, Konstantinos Drakakis, and Tanya Khovanova

Abstract. The connectivity of the autonomous systems (ASs) in the Internet can be modeled as a time-evolving random graph, whose nodes represent ASs and whose edges represent direct connections between them. Even though this graph has some random aspects, its properties show it to be fundamentally different from “traditional” random graphs. In the first part of this paper, we use real BGP data to study some properties of the AS connectivity graph and its evolution in time. In the second part, we build a simple model that is inspired by observations made in the first part, and we discuss simulations of this model.

I. Introduction

The growth of the Internet and the connectivity of the autonomous systems (ASs) is not governed by one central authority. The Internet has grown rather in a more or less haphazard way, governed mostly by local decisions, insofar as they are consistent with generally prevailing rules or protocols. The first models of the AS connectivity graph, in which nodes correspond to either ASs themselves or to individual routers within the ASs (depending on the desired level of detail) and edges correspond to direct (physical or virtual) connections between them, were therefore quite naturally based on random graphs. (See [Zegura et al. 97] for a summary of graph models up to 1997.) However, it was observed [Siganos et al. 03] that several interesting characteristics of the AS

connectivity graph exhibit a strikingly “non-random” behavior. In particular, the *degree distribution*, i.e., the number of nodes with k neighbors, decays like a negative power of k for large values of k . Likewise, the eigenvalues of the graph’s incidence matrix show a power law decay [Siganos et al. 03]. Random graphs typically exhibit much faster decay for these quantities [Bollobás 01]. It follows that the Internet graph is fundamentally different from the random graphs studied traditionally in mathematics (see [Bollobás 01] for a very detailed study of “traditional” random graphs).

Subsequent research focused on both an explanation of the unusual characteristics of the graph and a critical evaluation of the measurements and the experiments used to determine these characteristics. On the latter issue, it has been pointed out that the data used to establish the power decay laws are often derived from summaries of BGP-traffic exchange, and it has been argued that such BGP-based measurements give an incomplete picture of the connectivity [Chen et al. 02]; additionally, BGP data may contain all sorts of irregularities, such as loops, tangles, and ramifications, which obscure to some extent the perception of the AS connectivity graph (see [Broido et al. 02] for a very comprehensive description). Therefore, one needs to be careful with the use of measurements, and tools for their evaluation need to be developed [Floyd and Kohler 03]. Nevertheless, a consensus seems to have been reached that the observed power laws appear to hold at least approximately [Albert and Barabási 02, Barabási and Albert 99, Chen et al. 02, Siganos et al. 03, Medina et al. 00].

This paper presents a phenomenological approach to the issue of understanding the power law. It was shown by Barabási and Albert [Barabási and Albert 99] that the combination of two simple principles, *preferential attachment* (PA) and *incremental growth*, leads to graphs that exhibit a power law decay for the node degree distribution, with exponent -3 (i.e., the number of nodes with degree d is inversely proportional to d^3). These principles appear reasonably plausible at first sight: the Internet is certainly growing incrementally, and it is likely that as new ASs are created, they take into account the success of the existing nodes in the setup of their own connectivity. On the other hand, it is clear that the true picture must be much more complex; the actual mechanisms responsible for the occurrence of power law decay in the Internet graph are surely more complex and remain undetermined. That Barabási and Albert’s explanation is not sufficient is already apparent from the value of the exponent (closer to -2 than -3) in the approximate power law observed in the degree distribution of the AS connectivity graph and also from the deviations from a pure power law behavior for both very large and very low degrees. A more complex explanation than [Barabási and Albert 99] is proposed in [Chen et al. 02]; it is based on financial considerations.

This paper consists of several parts. We started by extracting in a systematic way consecutive weekly summaries of BGP data from <http://moat.nlanr.net/AS> (which we shall abbreviate as *nlanr-routeviews*), spanning a period of slightly over two years; these BGP data have since been archived at <http://archive.routeviews.org/oix-route-views/> (we thank one of our reviewers for pointing this out to us). The archive *nlanr-routeviews* gives BGP snapshot data, observed in principle every two hours. Occasionally such a two-hourly snapshot is missing. Downloading all the two-hourly snapshots for every week was not feasible when we undertook this study; for this reason we used a sampling strategy to compile weekly summaries. For each weekly summary we downloaded several snapshots spread over the week. To make the weekly summary, we simply took the union of the sampled snapshots for that week; an edge is present in the summary graph if and only if it is present in at least one of the corresponding snapshots.

By looking at the weekly summaries, we hoped to obtain a less fragmented view of the AS connectivity graph than is typically given by the BGP snapshot data, in which the momentary silence of servers that happen to be down at the time of the snapshot may obscure parts of the connectivity graph. (Of course, if a server is down for the whole week over which we carry out a summary, we have the same problem; this is less likely, however.)

In the first part of this paper, we examine the fairly detailed picture of the time-evolving AS connectivity graph that emerges from this wealth of data. We shall consider two different ways to “slice” the data. On the one hand, we determine, week by week, the empirical degree distribution; this turns out not to change much over the observation period, and therefore we shall also consider the summary of all these weekly degree distributions. We discern some fine structure, superposed on the power law behavior and evident at first glance in these plots, that suggests that different parts of the plot exhibit different regimes, i.e., that there seem to be different classes of AS. Inspired by a classification proposed in [Zegura et al. 97], we propose to distinguish three distinct regimes.

This first way of looking at the data corresponds to an overview of the whole population at single-time instants. On the other hand, we can take the “opposite” point of view and look at an overview in time of a single AS. More precisely, we look at the evolution (as a function of time) of the degree of individual ASs; we observe that there is a marked difference in this behavior for the ASs in the three classes first identified in single-time population overviews. (The distinction of exactly three different classes is to some extent arbitrary, but is in fact adequate. Based on the arguments in [Zegura et al. 97], it is reasonable to distinguish at least three qualitatively different classes of AS. A larger number of classes turns out not to be necessary: we shall show later that introducing finer distinctions does not lead to significantly different results.) Finally, we also

derive from the data the “popularity” of the different ASs, measured by how many new ASs choose to connect to them each week. Again, this popularity measure indicates the existence of different regimes, and we study in particular the interaction between and among the three proposed classes.

In the second part of the paper, inspired by the observations in the first part, we formulate a model that is slightly more elaborate than a simple preferential attachment, in that it allows each of the three different regimes to follow its own preferential attachment strategy; to set the parameters for these strategies, we use values determined by our “popularity” measurements as derived from the weekly BGP summaries in the first part of the paper. A crude mathematical analysis of the model, with these parameter values, shows that it does indeed lead to a degree distribution that exhibits the observed deviation from Barabási’s power law decay. Apart from this “back of the envelope” calculation, we also present results of simulations that are less approximate. Finally, we discuss our approach and test its validity in several ways. In particular, we compare the distribution of the minimum path length between any two nodes for our model with the corresponding distribution for the measured graph. We find them to be very close to each other. We also compare the largest two eigenvalues of the incidence matrix for simulated and observed graphs, with again very similar results.

As a warm-up, we give in Section 2 a quick review of the Linear Preferential Attachment Model (LPAM) of Barabási [Barabási and Albert 99]. Numerous alternative models based on the principle of preferential attachment have been proposed in the literature (see [Albert and Barabási 02] for more details); since we give LPAM here only as a reference point, we shall not further discuss these models here.

2. The Linear Preferential Attachment Model (LPAM)

Consider a graph G_t that “grows” in time as follows. At time $t = 0$, the graph G_0 is a clique of m nodes. The passage of time is modeled by adding constant discrete time increments, so that the values of t can be taken to be simply the positive integers. For each $t \in \mathbb{N}$, the graph G_{t+1} contains one more node than G_t ; thus, edges in G_{t+1} consist of all the edges that were present in G_t , augmented by m edges that all connect the new node with m (not necessarily different) nodes in the old graph G_t ; these edges model the connectivity with which the new node is “born.” The nodes of G_t with which the new node connects are chosen probabilistically; for each of the m new edges and for each of the nodes in G_t , the probability that the edge connects to that particular

node is proportional to the degree of that node in G_t . Accordingly, new nodes “prefer” to “attach” themselves to nodes that have a high degree already, and that preference is expressed by an increase in connection probability that varies linearly with the degree, hence the name *Linear Preferential Attachment Model*.

2.1. Derivation of the Degree Distribution for LPAM

Let us denote by $E_i(t)$ the expected number of edges adjacent to the i th node at time t . Up till time t , a total of t nodes and mt edges have been added, and thus $\sum_i E_i(t) = [m(m-1) + 2mt]$, since we had a clique at $t = 0$, and each edge gives a contribution of two to the sum of degrees (one for each of its endpoints). If we assume that the initial m nodes are numbered from 1 to m and subsequent nodes continue to be numbered by the successive integers, then we have, for any $j \in \mathbb{N}$, that $E_{j+m}(j) = m$; it follows that

$$\begin{aligned} \Delta E_i(t) &= E_i(t+1) - E_i(t) = m \frac{E_i(t)}{\sum_i E_i(t)} = \frac{E_i(t)}{m-1+2t} \\ \implies E_i(t+1) &= \left(1 + \frac{1}{m+2t-1}\right) E_i(t) \\ \implies E_i(t) &= m \prod_{s=i-m}^{t-1} \left(1 + \frac{1}{m-1+2s}\right) \text{ for } i > m \text{ and } t > i-m. \end{aligned}$$

From this, one can determine the dependence on i of the behavior of $E_i(T)$ for large T :

$$\begin{aligned} E_i(T) &= m \exp \left[\sum_{s=i-m}^{T-1} \ln \left(1 + \frac{1}{m-1+2s}\right) \right] \\ &\approx m \exp \left(\sum_{s=i-m}^{T-1} \frac{1}{m-1+2s} \right) \\ &\approx m \exp \left(\frac{1}{2} \ln \left[\frac{m-1+2T}{m-1+2(i-m)} \right] \right) = m \left(\frac{2T+m-1}{2i-m-1} \right)^{1/2}. \end{aligned}$$

An even simpler way to guesstimate this behavior is to interpret $\Delta E_i(t) = E_i(t+1) - E_i(t)$ as an approximation to the derivative of $E_i(t)$ and to transform the difference equation into a differential equation:

$$\frac{dE_i(t)}{dt} \approx \Delta E_i(t) = \frac{E_i(t)}{m-1+2t} \implies \ln(E_i(t)) \approx C + \frac{1}{2} \ln\left(\frac{m-1}{2} + t\right).$$

Combining this with the initial condition $E_i(i-m) = m$ (we implicitly assume $i > m$ again) leads to

$$E_i(t) \approx m \left(\frac{2t+m-1}{2i-m-1} \right)^{1/2}.$$

This means that if we pick, at some time T , a random node i from the collection of nodes $\{m+1, m+2, \dots, m+T\}$, the probability (with respect to this process of selecting i randomly) that its expected degree (expected with respect to the random process of graph generation) does not exceed some positive real number x is

$$\begin{aligned} \mathbf{P}(E_i(T) \leq x) &= \mathbf{P}\left(m \left(\frac{2T+m-1}{2i-m-1}\right)^{1/2} \leq x\right) \\ &= \mathbf{P}\left(\frac{2T+m-1}{2i-m-1} \leq \left(\frac{x}{m}\right)^2\right) \\ &= \mathbf{P}\left(i \geq \frac{m^2(2T+m-1)}{2x^2} + \frac{m+1}{2}\right) \\ &= 1 - \left(\frac{m}{x}\right)^2 \left(1 + \frac{m-1}{2T}\right) + \frac{m+1}{2T}. \end{aligned}$$

For large T , this converges to a limit independent of T , namely $\mathbf{P}(E_i(T) \leq x) \approx 1 - \left(\frac{m}{x}\right)^2$; the corresponding probability density p_T is given by the derivative with respect to x of this expression, so that $p_T(x) \approx 2m^2x^{-3}$ for $T \rightarrow \infty$.

We have of course made some hand-waving approximations in this derivation; however, a more careful mathematical analysis leads to the same strict power law x^{-3} .

Note that preferential attachment plays a crucial role in our derivation of the power law degree distribution: if the same incremental growth procedure is followed for G_t , but every new node connects to m preexisting nodes that are picked uniformly randomly instead of according to LPA, the corresponding graph exhibits an exponential degree distribution [Callaway et al. 01].

2.2. Discussion

Although the model in Section 2.1 leads to a power law independent of time, it has some very serious limitations, which render it unrealistic:

- The decay exponent is fixed to 3. However, a similar argument for the occurrence of a PA (Preferential Attachment) principle can be made for incrementally growing graphs in simplified descriptions of many different systems, such as, for instance, the http link graph (two web addresses are linked if one contains a hyperlink to the other), the actors' graph (two actors are linked if they have played in the same movie), etc. A closer look at the distributions for these graphs show that they do indeed decay as a power law, but they typically each have their own exponent, often with a value between 2 and 4 but not exactly equal to 3 [Barabási and Albert 99, Albert and Barabási 02].
- The degree of a node grows according to the time of its initial appearance; in this model, the earlier a node has appeared, the more probable it is that its degree will be larger at a given time, and the time of first appearance is a dominating factor determining that probability. It is therefore exceedingly unlikely in LPAM that a later born AS “overtakes” ASs that appeared much earlier. This is not true in reality, though: some ASs start out with a small degree and remain so, whereas others, born later, receive connections extremely fast and soon overtake older and smaller ASs. We will return to this point later. As an example, consider Google: it was one of the last search engines to be introduced, yet today it essentially monopolizes its field.
- It is unrealistic to assume that the number of initial connections made by a new node, upon its entry in the graph, is deterministic and constant.
- The model allows for more than one connection between two nodes, but this situation hardly ever arises in practice.
- The model contains no mechanism allowing for edge or node dropouts from the graph; these dropouts are fairly frequent in practice.
- New edges appear when new nodes appear, connecting a new and an old node. In practice though, connections between old nodes also show some dynamics: new links between two nodes may be created, or existing links may be deleted.

3. Analysis of the AS Connectivity Graph from Experimental Observations.

3.1. Description of the Data Set

Our main data source was the archive of publicly available BGP measurements that can be found at <http://moat.nlanr.net/AS>, which we call *nlanr-routeviews* for short. The data at *nlanr-routeviews* span a period from 1997 to 2000, with measurements made on a two-hourly basis (with some time slots missing). However, miscommunications and misconfigurations (see [Broido et al. 02] for an extensive discussion) occasionally cause some ASs or some edges to disappear from one measurement to the next, to reappear again shortly after. In order to reduce this “flickering,” we merged the data from *nlanr-routeviews* into weekly summaries. Each summary consists of a list of pairs of AS numbers. To compile a weekly summary, we queried the *nlanr-routeviews* site for seven snapshots (corresponding to independently-picked, random timings on each of the seven days of that week) and made the union of the snapshots thus obtained; each link that occurs on one (or more) of the requested snapshots is listed by giving the pair of its end nodes. There is significantly less flicker in the weekly summaries than in the individual snapshots. (There probably would have been even less if we had downloaded all the snapshots for each week, but the size of the files rendered this infeasible; moreover, experimentation showed that the additional accuracy gained would have been slight.) We decided to discard the last four weeks for 2000 listed in *nlanr-routeviews* because the measurements for these weeks were clearly less complete than for the preceding weeks. We also observed that week 61 was similarly anomalous; this we addressed by “patching”: the data in week 61 were replaced by a copy of week 60 (which was thus repeated). To check that our analysis was robust for the choice of patching, we carried out the same analysis for other patch choices, such as replacing week 61 by a copy of week 62 (instead of 60) or just dropping week 61 and renumbering all the weeks after 61 by subtracting 1 from their label. The results of our experimental analysis were the same in all cases.

We shall give the name RV to the data set that we obtained in this way, giving weekly summaries of the Internet (or at least, of the part of the Internet that was “visible” to the measurers who put together *nlanr-routeviews*). The data set RV is publicly available at <http://www.pacm.princeton.edu/pacmrig>; the web site contains not only the list of edges in the AS connectivity graph for each of the 109 weeks in the RV observation period but also many other files that give insight to the structure of RV. All the files are given in a simple text format that makes it easy to graph them with a variety of graphics packages. Because it is not clear to us, as of the writing of this paper, how accessible *nlanr-routeviews*

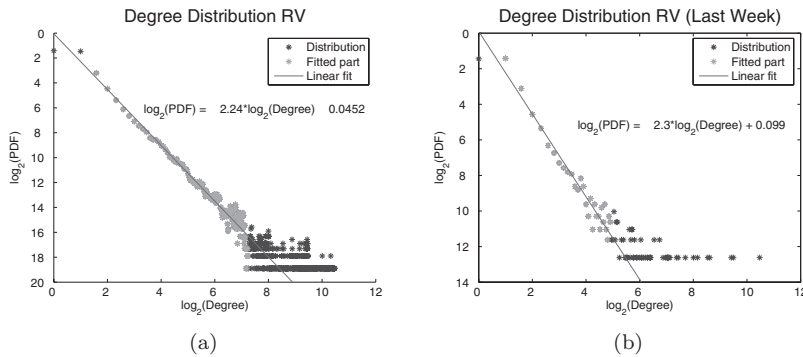


Figure 1. The degree distribution for the data set RV over (a) all weeks and (b) the last week only. Observe the approximate power law behavior in both graphs, indicated by the linear fit in these log-log plots.

will remain, the web site also gives the original snapshots (obtained by querying nlanr-routeviews in late 2001) from which RV was compiled.

We next describe several ways in which we analyzed RV.

3.2. The Degree Distribution

The degree distribution of data set RV is shown in Figure 1. In Figure 1(b), we show the degree distribution of the AS connectivity graph according to the summary of the last week in RV; to obtain Figure 1(a), we merged the degree distribution tables of all the individual weeks. (The number of points in Figure 1(a) is accordingly larger.) The distributions are plotted on a log-log scale, and in each case the best linear fit for the central part of the graph (more about this choice later) is also shown. In both cases the best linear fit gives a reasonable approximation, although it is clear that the observations show a structured deviation from this fit. The slopes of the best fit lines, respectively -2.24 and -2.3 , show that the power in the approximate power law exhibited by these distributions lies between 2 and 3; this is in accordance with the previous literature [Albert and Barabási 02, Barabási and Albert 99].

3.3. Incompleteness of the RV Data

The weekly summaries in the RV data set, agglomerated from daily snapshots (as described in Section 3.1) give a more complete picture of the AS network than each snapshot by itself. Nevertheless, given the sheer size of the system, it would be surprising if even these summaries were completely accurate. In fact, it has been documented that the picture is not complete: see [Willinger et al. 02],

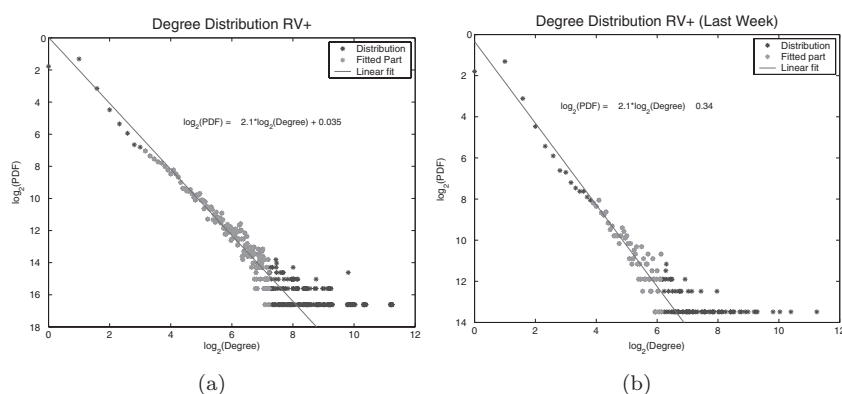


Figure 2. The degree distribution for the data set RV+ over (a) all weeks and (b) the last week only. Observe the approximate power law behavior in both graphs, indicated by the linear fit in these log-log plots, and how these figures are similar to the ones in Figure 1.

which also includes a discussion of different effects that may have obscured part of the network from the nlanr-routeviews snapshots. Different snapshots of the network, incorporating all the information from nlanr-routeviews but also from other sources, are given on S. Jamin's web page <http://irl.eecs.umich.edu/jamin/>; we shall refer to this more complete data set as RV+.

The dates for which this more complete information was compiled (nine consecutive weeks whose end dates are from March 31 to May 26, 2001) by S. Jamin fall unfortunately several months after the period covered by RV, and the two can thus not be directly compared. At the Oregon web site <http://archive.routeviews.org> (which we shall call Oregon-routeviews), more recent BGP measurements can be found, again taken at the two-hourly intervals familiar from nlanr-routeviews; they start in October 2001, however, several months after the end of Jamin's data set, making it also impossible to compare the Oregon-routeviews datasets directly with Jamin's. Although Jamin's data set contains probably significantly more links than the routeview BGP tables, the degree distributions for RV and for RV+ turn out to be very similar, although a detailed comparison does show differences (highlighted in [Chang et al. 04, Willinger et al. 02]). For the characteristics that we study in this paper, the similarity between the behavior of the two data sets, as seen by comparing, e.g., Figures 2 and 1, is sufficiently close that we felt confident that the RV data set was already capturing the essential relevant features.

In what follows, we shall study the behavior in time of AS degrees, the connection behavior of new ASs as they are born, and the relationship of these

with the shape of the degree distribution curve. The more complete snapshots at <http://irl.eecs.umich.edu/jamin/> are given for a few dates only, so that it is not possible to extract from them the temporal behavior that we can find in either nlanr-routeviews or in Oregon-routeviews. We shall therefore stick to a systematic study of the RV data only. Since the lack of completeness of RV doesn't prevent it from giving us a correct picture of the degree distribution, we shall (with some leap of faith) assume that it likewise does not affect the (fairly coarse) temporal features that we study.

3.4. Degree Evolution Patterns

The degree distribution already showed that the observations deviate from the LPAM model in two crucial ways: the degree distribution does not follow a “pure” power law, and the power of its best power law approximation does not equal 3. We now turn to an inspection of the evolution in time of the degree of individual ASs. The (continuous-time approximation of the) degree evolution according to LPAM (see Section 2.1) is proportional to $t^{1/2}$, for large t . This is not what we observe in the data.

Based on the evolution patterns of their degree, we distinguish three different types of AS:

- By far, the largest group is constituted by ASs that start out with a small degree and that never connect to many other ASs: more than 70% of the ASs present at the end of the observation period in RV have a maximum degree that does not exceed 3 during the entire observation period. For these ASs, the degree can fluctuate; for most, this fluctuation consists in at most short-lived deviations from an otherwise constant plateau (see Figure 3(e)), but for some (about 5% of this group of ASs) the fluctuations are more frequent and serious, to the point of looking random. (See Figure 3(f) and also Table 1 for some statistics.)
- At the other end of the “degree spectrum,” we have ASs with a degree that is large throughout the whole interval under consideration and that seems to increase linearly in time over long time stretches. Two examples are shown in Figure 3(a) and (b). There are few ASs of this type, and they are highly connected among themselves.
- The remaining ASs have a maximum degree that is larger than 3 but never become as large as the giants in the previous class. Their degree also seems to increase more or less linearly in time after the AS is “born,” albeit not as quickly as for the giants. See Figure 3(c) and (d).

		Number n of demises + rebirths						
		$n = 0$	$n = +1$	$n = -1$	$n = 2$	$2 < n \leq 9$	$9 < n \leq 29$	$n > 29$
Degree d	$1 \leq d \leq 3$	1333	3074	350	483	592	22	0
	$4 \leq d \leq 24$	897	346	50	35	69	14	1
	$25 \leq d \leq 299$	116	7	3	2	1	0	0
	$300 \leq d \leq 499$	2	0	0	0	0	0	0
	$d \geq 500$	3	0	0	0	0	0	0

Table 1. AS demise + rebirth number according to (maximum) degree; the category $n = -1$ consists of all the ASs in the data set that are present at the beginning but that disappear never to reappear again (at least during the observation time interval); $n = 0$ consists of all the ASs that are present in the beginning and never disappear; $n = 1$ consists of all the ASs that are born during the observation time interval and never disappear; all the remaining ASs experience at least two transitions from/to limbo, and n counts the total number of them (each demise contributes 1, as does each birth/rebirth). For some representative ASs of these categories, see Figure 3. The sum of all table entries gives the total number of ASs considered; there were 3069 ASs present in the first week.

The degree fluctuation, which can be so pronounced that an AS may disappear from the data set RV at some point in time and reappear again some time later, is an interesting feature. This type of “flickering” is less prevalent in the RV data set (because they are weekly summaries, in the sense described in Section 3.1) than in the daily snapshots in the nlanr-routeviews themselves, yet it is not completely eliminated. The demise + rebirth of an AS, as observed in the RV data, may correspond to a temporary suspension of activity by that AS, or it may reflect a shortcoming of the data set. Table 1 gives an overview of the demises and rebirths in the RV data set. To compile this table, the evolution through time of every node i was checked. If the node was present at the start of RV and then disappeared at some time, never to be reborn, then it was given a $n = -1$ flag; if the node was present at the start of RV and never disappeared, then it was given the value $n = 0$; if the node didn’t exist yet at the start of RV but never disappeared from the RV data after being born, then it was given the value $n = 1$. In all other cases, the node underwent a certain number of deaths and (re)births, and n was the total number of transitions from or to “nonexistence” that took place during the observation interval. Table 1 tabulates the numbers of nodes i according to the maximum degree $d_i = \max_t (E_i(t))$ that they attain during the observation period and their demise + rebirth number n_i .

The table shows that the phenomenon is more prominent among low degree ASs. At this point, we do not know whether this is an actual phenomenon or just an artifact of incomplete or defective BGP measurements (see the related

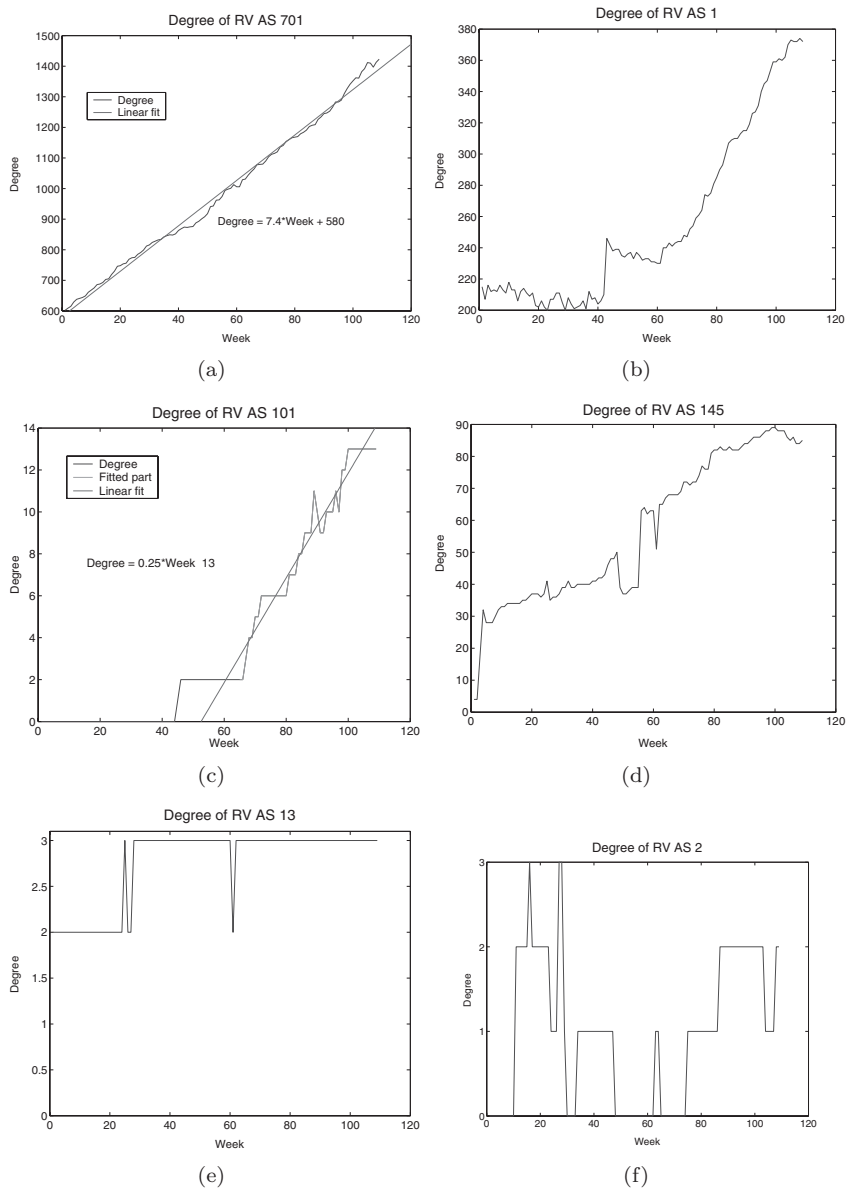


Figure 3. The degree evolution of the several types of AS: (a) the largest degree AS in RV, (b) a large degree AS in RV, (c)–(d) medium degree ASs in RV, (e)–(f) small degree ASs in RV.

discussion in [Chen et al. 02, Floyd and Kohler 03]); further study would be needed to clarify this. For our purposes here, the important result of Table 1 is that it shows that only a minority of the ASs is affected; we shall see later that our observations are robust to whether we exclude these ASs or not.

3.5. The Empirical Attachment Probability

We have seen that LPA, along with incremental growth, cannot be the complete explanation for the observed power law behavior of Internet growth. On the other hand, the prediction made by LPA is not far wrong. This suggests that there may indeed be some preferential attachment in play. In this subsection we use the RV data set to derive an Empirical Attachment Probability from the observations. We can then assess to what extent the probability, for each existing node, of being picked as the endpoint by a new node establishing its connections after being “born,” is indeed linked to the degree of the old node; if this quantity does make sense, we can also assess to what extent it deviates from LPA.

To compute the probability with which new nodes attach themselves to different types of old nodes, we must first agree on a definition of “new.” If it were the case that every AS present in RV at week t is automatically in RV for all weeks $t' > t$, then we could simply count an AS as new in the week where

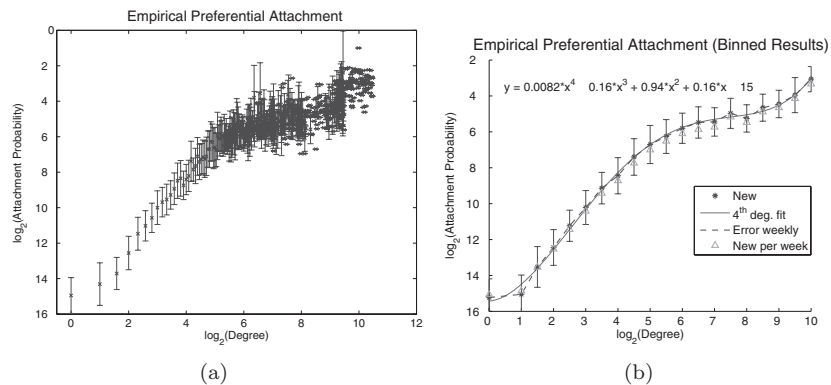


Figure 4. Log-log plot of the Empirical Attachment Probability as a function of d , unbinned in (a) and binned geometrically, with ratio $\sqrt{2}$, in (b). Figure (b) shows the result for both the generous and the cautious definitions of a “new” AS. (See text: each error bar shows the first and third quartile for the d -bin, with the mean in the middle, for the generous definition. The mean of the cautious definition is also indicated.) The resulting curves are virtually identical. This figure also shows the best fourth order polynomial fit for the data.

it first appears and as old in subsequent weeks. However, we saw earlier that there are ASs in RV that disappear and then reappear some weeks later (see also Figure 3(e)), so that the situation is not that simple. Because we identify an AS only by its AS number, and AS numbers that have been decommissioned can be reused, some of the disappearing and reappearing AS numbers may well correspond to two different autonomous systems, each of which should be counted as “new” when it appears. On the other hand, some spurious appearances and disappearances of ASs in the data set are surely due to either an outage that takes an AS out of the network for a while or to the imperfections of the BGP snapshots at nlanr-routeviews that may well miss observing an AS for even a whole week: in neither of these cases would it be reasonable to consider the AS as “new” when it reappears in RV. Since we had no means to distinguish all these different possibilities, we decided to carry out the empirical analysis for the two extreme cases:

1. in the *generous analysis* we count every (re)appearing AS as new on its “(re)birth” (i.e., node i is new at time t if $E_i(t) > 0$ and $E_i(t-1) = 0$, where in this context, we use $E_i(t)$ for the degree of node i observed in RV at time t),
2. in the *cautious analysis* an AS is counted as new only at its first appearance (i.e., node i is new at time t if $E_i(t) > 0$ and $\max_{t' < t} E_i(t') = 0$).

Of course, all ASs already in the network in the first week are considered as “old,” i.e., already existing.

Having decided what it means for an AS to be new, we can find, for each week t and each node i that is old in week t , the number $a(i, t)$ of connections made to node i by nodes that are new that week. We can then define the *Empirical Attachment (EA) Probability*, depending on degree d and in week t , as the ratio

$$\mathbf{P}_{\text{EA}}(d, t) = \frac{\sum_{i; i \text{ is old in week } t \text{ and } E_i(t)=d} a(i, t)}{\left(\sum_{j; j \text{ is old in week } t} a(j, t) \right) \#\{i; i \text{ is old in week } t \text{ and } E_i(t) = d\}}.$$

This gives, indeed, the (empirical) probability with which a specific old AS of degree d attracts new connections in week t , as inferred from the data.

We computed $\mathbf{P}_{\text{EA}}(d, t)$ for all degrees d that occurred in the data set and for all t . For each d , the results for different weeks t gave us a distribution of “measured” values for the Empirical Attachment Probability for degree d . The results are shown in Figure 4(a); for each d we indicate both the mean and the first and third quartiles of $\log(\mathbf{P}_{\text{EA}})(d, t)$, as t ranges over the whole observation time interval. It is clear from the figure that there is indeed a link between

the Empirical Attachment Probability and the degree of the “attachee” node. It is also clear that this link does not appear to be quite linear: if $\mathbf{P}(d)$ were proportional to d , then we would have $\log(\mathbf{P}(d)) = c + \log(d)$, i.e., the graph in Figure 4(a) would have been a straight line with slope 1. The graph clearly shows some deviation from this behavior, but its average slope is not far from 1: $\log(\mathbf{P})$ increases by about 13 for an increase in $\log(d)$ of about 11.

Note that each degree d in Figure 4(a) need not be represented in each week t . This is especially the case for some of the larger degrees, which belong to only a few ASs, each with a steadily (and fairly fast) increasing degree. These large degrees may therefore occur in a single week only, explaining why the means and first and third quartiles for each d collapse at the extreme right of the graph in Figure 4(a). To get a better reading of the empirical probability at these high degrees as well as in the middle degree range, where the crowding in Figure 4(a) makes it rather confusing, we binned the data geometrically with a ratio of $\sqrt{2}$. More precisely, for each k , we collected all the $\mathbf{P}_{\text{EA}}(d, t)$ with $2^{k/2-1/4} \leq d < 2^{k/2+1/4}$ in one group; the means and first and third quartiles of these groups are plotted as a function of $k/2$ in Figure 4(b).

All this was done for both the *generous* and *cautious* definitions of new ASs. The error bars (indicating first and third quartiles) in both Figure 4(a) and Figure 4(b) correspond to the generous analysis; in the more readable Figure 4(b) the means obtained from the cautious analysis are given as well. The curves showing the dependence on d of the Empirical Preferential Attachment, whether given by the generous or the cautious definitions, respectively, are virtually identical. It follows that the conclusions of our empirical attachment analysis do not depend critically on whether we are generous or cautious in our definition of newness, so that we can reasonably expect that for the “correct” but unknown identification of new ASs, which lies somewhere in between, the same conclusion will hold.

The plots of \mathbf{P}_{EA} incorporate an inflation effect, due to the increase in range of possible node degrees (as illustrated by Figure 3, ASs of type T1 and T2 (see Section 3.6) have degrees that grow in time) and of the sum of the degrees of all the nodes at time t . One can remove this inflationary effect by discounting, similarly to what is done in finance. Taking this effect into account affects the corresponding plots (of, e.g., the Experimental Attachment Probability) only slightly and has no effect on the conclusions that we draw; in order to not interrupt the flow of our argument, we therefore explain this subtle point separately, in the Appendix.

We can now compare this Experimental Attachment Probability with the results of a similar experiment for LPAM, that is, with the analogous plot, starting from a simulation of LPAM rather than from the data set RV. To construct our

simulation, we tried to build a parallel to the observed data as much as possible. We first grew a graph up to 3,500 nodes. (In order to avoid long-lasting problems caused by random fluctuations at the start, which would force us to either average over an enormous ensemble or to simulate for extremely long times, and also to be able to make a fair comparison between the LPAM simulation and our simulations later, we grew this initial graph in the same way as for the simulations of our own model; see Section 4.) Then, we added 100 increments of 60 more nodes each (so as to mimic the 100 weeks in RV, which saw the addition of 60 nodes per week on average; we will, by analogy, give a label t to the incremental groups of 60 nodes and their elements). For each t , each of the 60 nodes connected to the graph that it found at its birth according to the LPA principle, and we recorded $\mathbf{P}_{\text{EA}}(d, t)$ for each t ; we used $m = 1.5$. Figure 5(b) gives then the same figure as Figure 5(a), with the same binning and averaging procedures, but now for the LPAM simulation instead of the RV data.

A comparison between Figures 5(a) and 5(b) leads to several interesting conclusions.

- The LPAM simulation shows a very nice linear behavior (such as one would expect for a *Linear* Attachment Model) for small values of $k/2$.
- For the bins corresponding to low degrees, the behavior of RV deviates significantly from that of the LPAM simulation.
- For the middle range of degrees, corresponding (more or less) to $k/2$ between 3 and 8, the two plots, for RV and for LPAM, are surprisingly similar (see Figure 5(c)).
- For very high degrees, the LPAM simulation plot deviates significantly from linear behavior. This is caused by a minority of nodes (fewer than 10 %) that have very large degrees. One might believe that this is due to the inflation effect that we mentioned earlier, which will weigh heavier on high degrees than on low degrees; this is not the case, however (see Appendix). This “aberrant” behavior is probably due to exponential amplification of small deviations through the LPA mechanism itself.
- Unfortunately, this lack of reliability of the LPAM plot in Figure 5(b), as to “reading off” the underlying probability law in the high degree range, means that we cannot really deduce anything in this range from the similarity with or difference from its shape as compared to the RV plot.
- It is to be noted, however, that the RV data set has nodes of much higher degree than the LPAM simulation, even though they started from comparable initial conditions, indicating that the attachment probability behavior

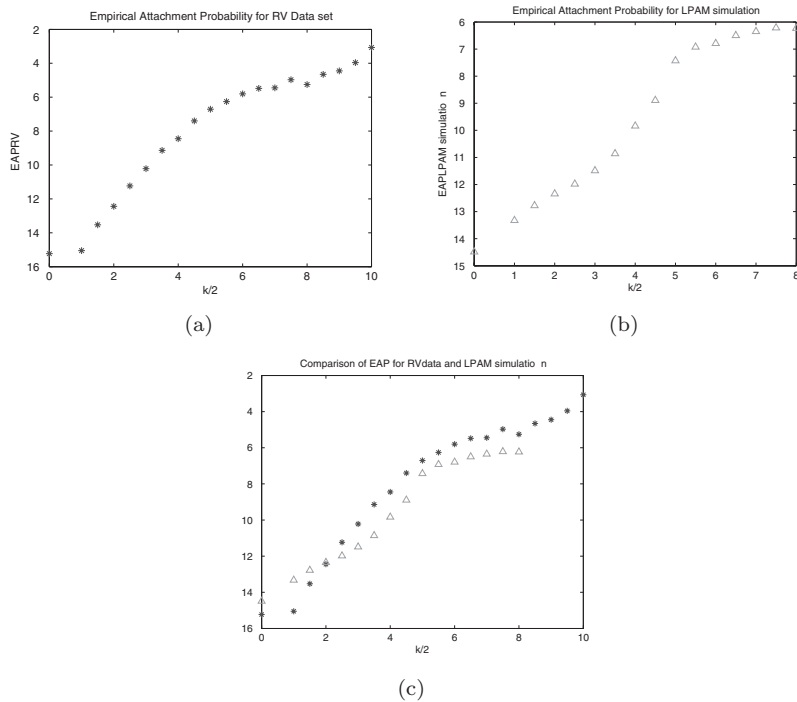


Figure 5. Graphs of the experimental preferential attachment probability versus the degree of the node to which attachment is made: for the RV data in (a) and for an LPAM simulation in (b). In both cases the probabilities are binned by degree in each week, discounted so that the different weeks can be compared reliably and averaged over weeks (see text). To compare the two better, they are plotted in the same graph in (c).

of very large degree nodes in RV must deviate from the LPAM simulation behavior.

- The comparison suggests therefore that RV shows the presence (roughly) of at least two and probably three different regimes, consistent with the observations made earlier.

Note that we have implicitly assumed, throughout this section, that a connection is always between a new node (which we interpret as initiating the connection) and an old node (which we view as the target of the connection). As mentioned earlier, there are connections in the real-world AS connectivity graph that do not fit into this pattern, because they connect two nodes that are either both old, at the time the connection is first established, or both new. When a connection is made between two preexisting (i.e., old) nodes, it doesn't come

into this subsection's discussion of connectivity behavior of new nodes at "birth." Connections between two new nodes do have a place in this discussion, however. A natural interpretation for these is that they are new-old connections in which the old node was born just a few days before the new node (after all, our RV summary lumps all the nodes born in a given week into to the same t -bin); because RV does not allow us to distinguish which node is newer among two nodes born in the same week, we decided to simply disregard these few connections in our analysis. They fortunately constitute a small minority among all the connections made by new nodes.

From now on, we will adopt the cautious definition of newness. Note that throughout this subsection, including in Figure 4, the degree of the old AS to which a new AS was connecting was taken to be the one at the time of connection. We will return to this point in Section 3.7.1.

3.6. Revisiting the Preferential Attachment Model

As pointed out above, different regimes can be distinguished in the degree distribution (see Figure 1) of the ASs in the Internet. In particular, the distribution has a middle region that nicely follows a linear log-log relation (or, equivalently, a pure power law); both the very large (above 300) and very small (3 or below) degree regions deviate from this behavior, with an almost flat region at small degrees. A similar distinction between small, middle, and high degree regions can be made in Figure 5(c), although the transitions appear to happen at slightly different degree values. This experimentally-observed behavior is in marked contrast with the behavior of LPAM simulations.

Implicitly, LPAM assumes that all nodes are equivalent, in the sense that they all have the same attributes and potential, and that only the time at which they join the graph affects their evolution. In practice, though, different ASs can function in fundamentally different ways. There exist some very large degree ASs, called "Tier 1," that form the core of the Internet (the "Internet backbone") [NRC 01]; next, one can distinguish ASs with a degree that is not quite as huge and that represent Internet Service Providers (ISPs) which offer Internet services to end user ASs through their own connection with Tier 1 ASs; finally, there are ASs that do not provide any Internet services to the outside world through themselves and that can be viewed as end users. Let us call these three categories of ASs by the names Type 1, 2, and 3, respectively (or T1, T2, and T3 in shorthand). From the interpretation we give them, as described above, we would expect Type 1 ASs to be very few in number, to have very large degree, and to be highly interconnected (maybe to the extent of forming a clique—see Section 3.4); new Type 1 ASs should appear rarely. In contrast, Type 2 ASs

should appear more often; at their “birth” we would expect them to Connect mostly to the backbone (Type 1 ASs) as well as to other ISPs (Type 2 ASs); it is unlikely that such a newly created Type 2 AS would choose to spend any of its own connection potential on end users (i.e., Type 3 ASs). Finally, interpreting the Type 3 ASs as just end users, we expect them to make very few connections, and those only to Type 1 or 2 ASs, not to other end users (Type 3 ASs). (Note that, in reality, there are performance and economic reasons that lead some end users to set up a private direct connection between them. These connections would not show up in the BGP tables and are therefore outside the scope of this paper.)

Our next task is to revisit the RV data and check whether this intuition is indeed borne out. Note that our classification follows closely the one proposed earlier in [Zegura et al. 97], with the Types 1, 2, and 3 corresponding to Transits, Stubs, and Leaves, respectively.

3.7. Study of the Empirical Preferential Attachment Strategies

3.7.1. The Kismet Assumption. The discussion at the end of the preceding subsection suggests that we consider different Preferential Attachment strategies for an AS at birth, depending on its own Type as well as on the Types of the targets of its possible connections. In order to extract these different behaviors from the data, we must therefore make a distinction between the births of Type 1, 2, and 3 ASs. As illustrated by, e.g., the graphs Figures 3(c) and (d), a newborn AS typically starts small: its degree at birth does not (necessarily) indicate its Type. Therefore, we must classify an AS, at the time it is born, based on its future behavior.

In reality, as opposed to our model, ASs are of course not able to predict the future and cannot base their connection strategies on such considerations. However, it is entirely reasonable to assume that most ASs know at their birth whether they are a new backbone or Type 1 node (the emergence of which we expect to be extremely rare; none was in fact observed in RV), a new ISP provider or Type 2 node, or a new end user or Type 3 AS. Moreover, it is also reasonable to expect that if a new node is of either Type 1 or Type 2, its birth and availability will be advertised immediately (or even in advance of the event) and widely to its possible customer basis, i.e., to the ASs seeking new connections (at their birth). Based on these arguments, we shall therefore assume that *at the time of its introduction an AS is already labeled with its type* and that its own connection strategy as well as the connection strategies of other ASs are all in accordance with this label, rather than with the label one might tentatively assign to the AS based on the value of its degree at that time. Once the *strategy*

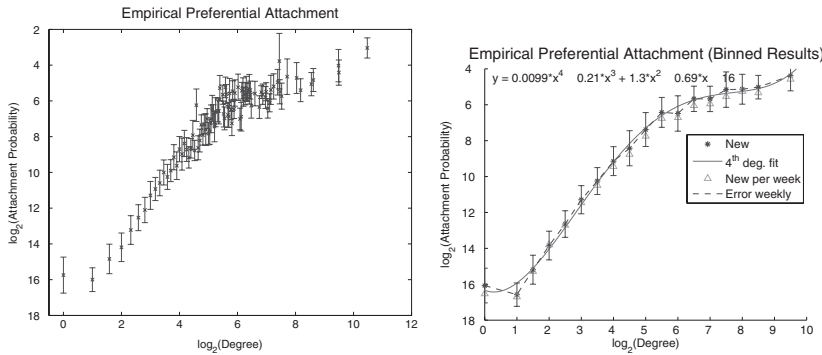


Figure 6. Empirical attachment probability as a function of the *maximum* degree of the connectee AS, to test the effect of the Kismet assumption (see text). As in Figure 4, the graph on the left shows the results unbinned; on the right the results are binned geometrically, with ratio $\sqrt{2}$.

for the preferential attachment is thus fixed, the attachment probability itself takes into account the actual degree of the connectee AS in order to determine its probability of being picked as endpoint for a new connection. We will call this assumption the *Kismet assumption*, after the Turkish/Persian/Arabic word for *fate*.

A similar argument is made in [Tangmunarunkit et al. 01], where the AS degree power law was analyzed in detail, distinguishing ASs by their size (in terms of number of servers, routers, etc.). Note also that the Kismet assumption makes it possible for the degree of an AS to grow faster (if it is of a different Type that turns it into a more desirable connectee) than that of other ASs, even if they were born earlier. The Kismet assumption will play a significant role in our theoretical derivations below.

The careful reader may have noticed that we have argued ourselves into somewhat of a paradox here: we based our distinction of the ASs into (three) different Types in part on the existence of different regimes in the dependence of the Empirical Attachment Probability on the degree of the connectee, as illustrated in Figures 4 and 5. Yet this dependence was computed with the (implicit) assumption that all the ASs of the same degree were equivalent, an assumption that is no longer valid if we adopt the three Types and the Kismet assumption. Because we expect most Type 1 and Type 2 ASs to have rapidly growing degrees, we can hope that there are at any time a negligible number of small-degree new T1s and T2s, so that the resulting overall effect of the change of attribution to their correct Type will be small. Figure 6 shows that this intuition is correct. This figure is the equivalent of Figures 4(a) and (b), in which the abscissa denotes the

maximum degree of the connectee AS (and therefore implicitly its Type) instead of the degree of that connectee AS at the time of connection. In other words, the graph is obtained in exactly the same way, except that we replace the definition of $\mathbf{P}_{\text{EA}}(d, t)$ by

$$\mathbf{P}_{\text{EA}}^{\text{K}}(d, t) = \frac{\sum_{i; i \text{ is old in week } t \text{ and } \max_{t'} E_i(t')=d} a(i, t)}{\left(\sum_{j; j \text{ is old in week } t} a(j, t) \right) \#\{i; i \text{ is old in week } t \text{ and } \max_{t'} E_i(t')=d\}},$$

where the superscript K stands for “Kismet.” Note that whereas the computation leading to $\mathbf{P}_{\text{EA}}(d, t)$ did misclassify some old ASs, possibly assigning them to a Type to which a connection is less attractive (i.e., Type 3 instead of 2, or 2 instead of 1), the computation for $\mathbf{P}_{\text{EA}}^{\text{K}}(d, t)$ makes the opposite mistake, making old ASs more attractive than in reality, by linking them to their maximal degree. The similarity of Figure 4 and Figure 6 shows that the effect of these misassignments is negligible. We can therefore safely assume that the graphs give an accurate impression of the postulated “true” Empirical Attachment Probability (according to the Kismet assumption), which takes into account both the Type label of the candidate connectee ASs and their actual degree at the time of connection, and therefore must lie in between the two extremes given by Figure 4 and Figure 6. It follows that our observation of the different regimes in the attachment strategies persists.

(Note that when we replace \mathbf{P}_{EA} by $\mathbf{P}_{\text{EA}}^{\text{K}}$, in which old ASs are characterized by their maximum degree, the inflation phenomenon that we mentioned earlier does not occur.)

From now on, we shall assume that the Kismet assumption holds, unless we explicitly state otherwise.

3.7.2. The three AS types. In our discussion in Section 3.7.1, we implicitly assume that the degree- and Type-dependent behavior of the connection strategies remain constant in time. This can be tested by comparing the estimates of $\mathbf{P}_{\text{EA}}^{\text{K}}(d)$ computed over different time periods. By comparing weekly estimates of this empirical probability, we checked that this was indeed the case. Figure 7 illustrates the similarity for two different weeks during the observation period. On each of the 100 weekly graphs, we also plotted the corresponding line with slope 1, corresponding with the theoretical LPA probability, completely fixed by the degree distribution in the week under consideration. We then observed in particular that

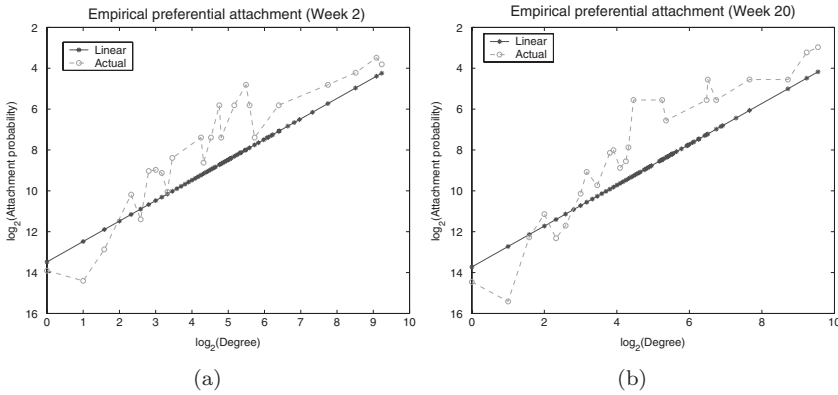


Figure 7. Empirical attachment probabilities for (a) week 2 and (b) week 20, along with the corresponding linear LPAM prediction.

- the overall behavior of the empirical attachment probability does not deviate very much from the slope 1 line corresponding to the simple LPA principle;
- in the Type 3 region (ASs with maximum degree below 4), the LPA theoretical line overestimates the number of connections that these ASs receive;
- in the Type 2 region (ASs with maximum degree between 4 and 300), the LPA theoretical line systematically underestimates the number of new connections made to these ASs;
- in the last region (ASs with maximum degree over 300), the LPA theoretical line again underestimates the EAP. Direct analysis of the data shows that these ASs form a clique, i.e., each one of them is connected to all the others.

We propose here to approximate the preferential attachment behavior for the different regions by a linear preferential attachment that is adapted to the region. In order to test how reasonable this hypothesis is, we plotted the attachment probability for each region separately and computed, by least-squares approximation, the straight lines that best fit the weekly EAP graph for each type. Figure 8 plots the histogram of the slopes corresponding to the first 100 weeks of RV for each AS type. Except for T1s, the slopes cluster near 1, so that we are indeed not far from the LPA principle (in which the slope is 1); otherwise, the three histograms are quite different.

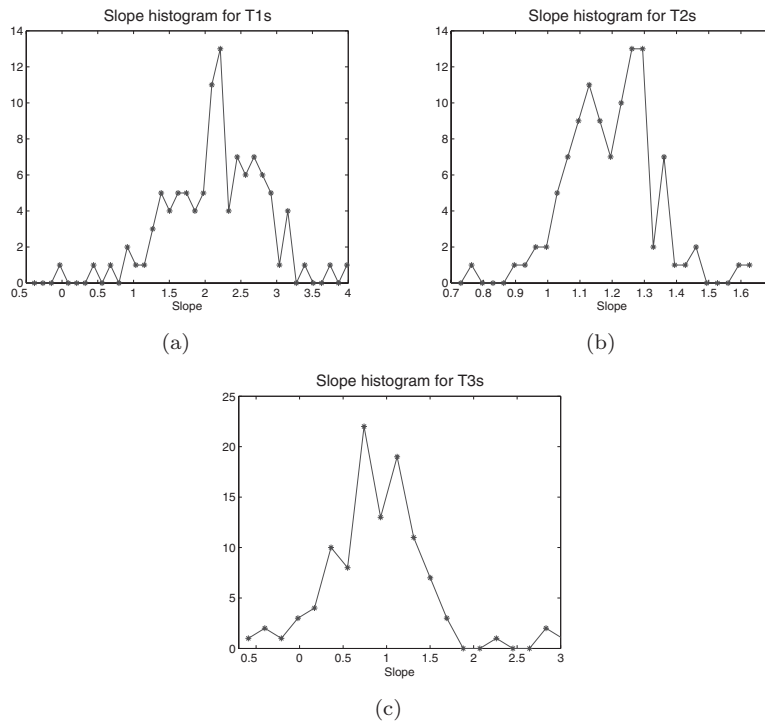


Figure 8. Histograms of the slopes of the least-squares best linear approximations to the weekly measured EAP, shown for each of the three regions: (a) T1, (b) T2, and (c) T3. In the T3 and T1 cases there are very few possible values for the degree and very few points, respectively, so that the best linear approximation is less reliable. The slope histogram is correspondingly more spread out for T3 and T1 (note the difference in abscissa scale in the 3 graphs); in fact, there were a few scattered outliers for T3 and T1 that fall outside the plotted scale.

3.8. Empirical Probability Parameters

We have presented a number of arguments to distinguish different classes among the ASs, and we have proposed to make a (fairly coarse) classification into three classes. We shall discuss later how satisfactory or acceptable this classification is. If we accept, at least for the time being, the categorization of the ASs into the three types described in Section 3.6, then we next have to find the values of the appropriate quantitative parameters. Consistent with the basic phenomenological approach of this paper, we shall derive them from the observed data set RV. As suggested by Figures 5 and 1, we put the cut-offs between our categories at 4 and 300, respectively; i.e., we will consider as T1s the ASs whose

(maximum) degree is above 300 and as T3s the ASs whose (maximum) degree is 3 or below. The others will be considered to be T2s.

Figure 9 shows each type's percentage among all the (then) existing ASs, weekly for the first 70 weeks. (We exclude weeks at the end because, not knowing the Kismet label of each AS, we have to deduce its type from its behavior after birth; this becomes less reliable if too few weeks are available for the AS to develop its "adult" behavior.) We see that the majority of ASs are T3s; T2s are fewer, and T1s are extremely rare. This is also what common sense suggests. In all cases, there appears to be a trend: the percentage of T3s increases, whereas the percentage of T1s and T2s decreases. Such a behavior is consistent, at least qualitatively, with the interpretation of T1s as "users" and T2s and T3s as "providers": the Internet is experiencing a rapid expansion during the observation time period, with an enormous growth of the number of users; the number of providers grows as well, but because of technological advances, it is natural to expect that, on average, providers can each take care of steadily more users as time progresses. (This is only qualitative because it is clear that the user-provider distinction is not so clear cut: one would expect a "true" ratio of users to providers to be larger than 7/3.)

Even so, the rate of change is not very high, so that constant percentages given by the averages over the whole time period are reasonable approximations. For instance, one can take the values attained by the data around week 50, the middle of the RV time period, when approximately 70% of the ASs are T3s and 30% T2s, while the T1s make up a negligible percentage. (See Figure 9.)

The determination of the "observed" probabilities by which ASs of different types connect to each other was based on the Kismet assumption. That is, for each new connection that involves a new AS, we interpret the connection as one that is "made" by the new AS, and we take note of the Kismet labels of the new AS and the old AS at the other end of the connection, as determined by the maximum degrees of both these ASs. For each type of AS, and over the first 70 weeks in the RV data set, we tally the total numbers of connections made to old ASs, according to their type; the percentages of these three types of connections are given by the three rows in Table 2. (We restricted the tallying to the first 70 weeks, again, because for ASs that are born near the end of the observation period, we may lack sufficient information to determine their type label. The 70-week cut-off is arbitrary and fairly conservative; moving it up or down by, say, 10 weeks doesn't change the table.) Note that the row for new T1s is empty: no new T1s were introduced into the Internet within the observation period; this is consistent with our assumption that such ASs are rare. The table also shows that, according to the data and our interpretation, 3.6% of new T3s connected to other T3s, and 9.2% of T2s connected to T3s. This is *not* consistent with

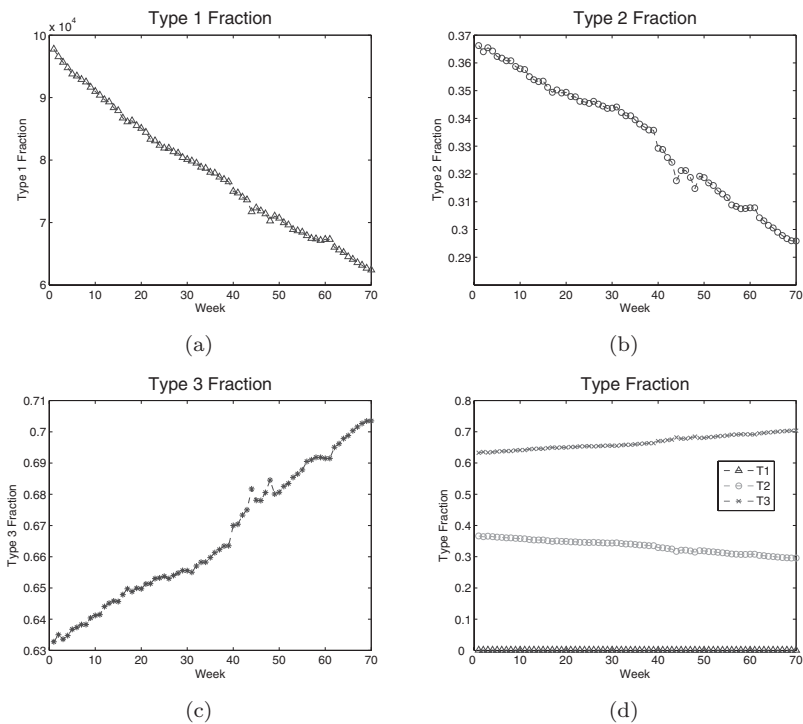


Figure 9. Fraction of ASs of Type 1, 2, and 3 per week, graphs (a), (b), and (c), respectively; (d) shows all 3 graphs together on the same scale.

our presuppositions: our discussion in Section 3.6 assumed that neither of these would happen. However, both percentages are small, and we shall neglect them in our analysis (and come back to their possible impact in our discussion at the end).

Within our framework, we could explain these “erroneous connections” as follows. The percentage of T3-to-T3 connections could be just an artifact of our (simplistic) definition of a T3; the percentage of T2-to-T3 connections results from only 47 connections in total, for most of which the initiating AS had a small degree for quite a long time after birth before growing into the T2 range. This could mean that these few ASs trans-typed: they might have been intended as T3s and then later changed their role, morphing into T2s. Under this interpretation, these 47 connections should be counted as T3-to-T3 connections as well. If we do this, then the number of T3-to-T3 connections as well as the total number of connections from T3s to other nodes goes up by 47, and the number of T2-to-T3 connections drops to 0. Under this correction, the entries of the table

		Old		
		3	2	1
New	3	3.6%	69%	27.3%
	2	9.2%	71.5%	19.3%
	1	—	—	—

		Old		
		3	2	1
New	3	5.4	67.7%	26.8%
	2	—	77.7%	22.3%
	1	—	—	—

		Old		
		3	2	1
New	3	—	71.6%	28.4%
	2	—	77.7%	22.3%
	1	—	—	—

Table 2. Connection probabilities between various types of ASs: as observed (on the left), reconsidering the 47 connections of T2-to-T3 nodes as T3-to-T3 nodes (middle), and neglecting the T3-to-T3 connections altogether (right). (See text.)

change and become as shown in the middle table of Table 2; the change is small. We can also decide to ignore all observed T3-to-T3 and T2-to-T3 connections altogether, in which case we obtain the right table of Table 2; the change is again small. (Of course, the unexpected connections may also simply be an artifact of the interpretation of the T3s as “users” and the T2s as “providers,” which was pointed out earlier to be overly simplistic.)

Finally, for later use we need to extract from the observed evolution data the degrees to be assigned to the T2s and T3s at their birth. These degrees are not always the same, and we have in fact a (fairly narrow) distribution for each of the two types. Taking the mean over all the new ASs of Kismet type 2, over all the weeks 1 to 70, we find that T2s have an average degree of 1.44 as they are introduced; the degree of each T2 becomes larger as time progresses and other newborn T2s or T3s connect to it. For T3s, we shall, in the simplified version where we neglect T3-to-T3 connections, assign a degree (which is their initial as well as their final degree under this simplifying assumption) that is the average degree of all T3s; from the data we compute that this average is 1.5. It is because these numbers are fairly close that it is very hard, in general, to distinguish a T3 from a T2 right at birth. We will come back to the exact distributions in Section 5.2.

3.9. Summary of Our Empirical Findings

The RV data set has the following striking features:

- The observed degree distribution seems to remain the same (up to inflation) over the full time period.
- The shape of the degree distribution suggests the existence of different types of AS; a first approximation suggests three classes.
- The evolution in time of the degree of an individual AS can take on different forms; to a first approximation, we can distinguish the same three classes by the pattern of their time behavior.

- As new ASs “are born,” one can define an Empirical Attachment Probability (EAP) as a function of the degree of candidate attachee nodes; this EAP has a stable profile over the full time period (again, up to inflation). One can again distinguish the same three classes in the EAP.

These conclusions are of course drastic simplifications of a much more complex reality. In formulating them, we paid attention to several issues, to ensure that the observation was robust and held for several different possible interpretations of the data. In particular, we acknowledge that it is impossible, from the RV data, to decide whether an AS is truly new if it shows up one week in the RV data set after being absent the previous week; we addressed the possible influence of mistakes by introducing a *cautious* as well as a *generous* definition of newness, and found that both extreme definitions led to the same observations. We also took note of several other ways in which the complexity of the data differs from the simplified picture: the EAP can be viewed only very approximately as given by a piecewise LPA, and the distinction into only three classes is surely a gross simplification.

In the next section, we shall propose a simple model that extends LPAM by incorporating into it the observations above.

4. The Reformed Preferential Attachment Model (RPAM)

4.1. Description of the Model

In this section we propose a new model, which we call the Reformed Preferential Attachment Model (RPAM). To build it, we incorporate (a simplified version of) the quantitative and qualitative knowledge drawn from the analysis of RV. In particular, we enrich the Linear Preferential Model with the existence of three classes of AS, each of which has, at birth, its own characteristic probabilities of connection to old ASs of the three types.

Before we start on the details of the model, we would like to emphasize that, as announced earlier, this is a *phenomenological* approach to model building. Although we don’t start from first principles (such as the dynamics induced by network protocols), our approach is not just a “fit” to the data either. We build a model with few parameters that each have a simple interpretation in terms of the observed average “connection” behavior of ASs at their emergence in the RV table. We use this interpretation to deduce the values of these parameters for the data set that we are studying; then, we study the dynamics of the model, with these (now fixed) values of the parameters, and compare the predictions of the model for the degree distribution (i.e., a *different* quantity) with the ob-

servations. (Straightforward fitting would correspond to leaving the parameters free throughout the process and fixing them at the end only, assigning to them the values that provide the best fit with *all* the observations.) If the model predictions are in good agreement with the data, then this indicates that the model and its parameters capture some “truth.”

Although complete understanding requires derivations from first principles, phenomenological studies are a good way to learn how to see the forest rather than the trees, often a necessary preliminary stage on the road to full understanding of a complex system such as the Internet and its growth.

In RPAM, the Internet is modeled as a random graph, the nodes of which represent ASs; the edges represent links between them. The nodes are divided into three categories, as explained in Section 3.6: Type 1, 2, and 3 (or T1, T2, and T3, in shorthand).

The simplest version of the model is built along the following lines:

1. At time $t = 0$, only T1s exist; we shall typically assume that they form a clique.
2. The time is discrete; at each new time instant, one new node is introduced into the network, and it can be T1, T2, or T3, according to fixed probabilities p_1 , p_2 , or p_3 , respectively (where $p_1 + p_2 + p_3 = 1$).
3. If the new node is a T1, then it connects to all existing T1s.
4. If the new node is a T2, then it connects to exactly m_2 nodes, which can be only of type T1 or T2, with probabilities p_{21} and p_{22} , respectively (with $p_{21} + p_{22} = 1$).
5. If the new node is a T3, then it connects to exactly m_3 nodes, which can be only of type T1 or T2, with probabilities p_{31} and p_{32} , respectively (with $p_{31} + p_{32} = 1$).
6. The m_2 (respectively m_3) choices of the type of AS to which a new T2 (respectively T3) node connects are all independent.
7. Multiple connections between nodes are allowed.
8. For each connection made by a new T2 to an (old) T1, the choice of the particular AS of type T1 (to which the connection will be made) obeys the Preferential Attachment Principle: for each old T1, the probability of being selected is proportional to its degree prior to the birth of the new T2 node.

9. *Mutatis mutandis*, the same applies to each new connection from a new T2 to an old T2, from a new T3 to an old T1, or from a new T3 to an old T2.

In the earlier sections, we have given an extensive discussion of the data set RV, and we can derive (approximate) values for the parameters m_2 , m_3 , p_1 , p_2 , p_3 , p_{21} , p_{22} , p_{31} , and p_{32} from the data set. We thus obtain, from our observation that about 70% of the nodes are T3s and 30% are T2s, from the entries in Table 2, from the average initial degrees 1.44 for T2s, and from the average degree 1.5 for T3s, the following parameter settings:

- $p_1 = 0$ (since no new T1s were born in the RV observation period),
- $p_2 = 0.3$,
- $p_3 = 0.7$,
- $p_{22} = 0.78$,
- $p_{21} = 0.22$,
- $p_{32} = 0.72$,
- $p_{31} = 0.28$ (again from Table 2),
- $m_2 = 1.44$,
- $m_3 = 1.5$.

Although this list has eight nonzero items (we disregard p_1 as a “parameter”), the three normalizations reduce the number of different choices to five.

Before we actually begin the derivation, we would like to bring the following points to the attention of the reader:

- It is of course completely unrealistic to assume that at time t only one AS is created: in every week of RV there are many new ASs. On the other hand, RV gives a summary of what happened during the whole week, and one could assume that at each particular instant at which an AS is created, it is the only one that gets born at exactly that time. In this view, going from one week to the next in RV corresponds to an increment in t that is much larger than 1. However, this means that if we simply transpose the trends, we shall deduce from linear or power behavior in t , to similar behavior in true time, as measured by the label of the weeks in RV, then we implicitly assume that the numbering of the weeks corresponds (more or less) to a linear increase in t , i.e., a linear increase in the total number of

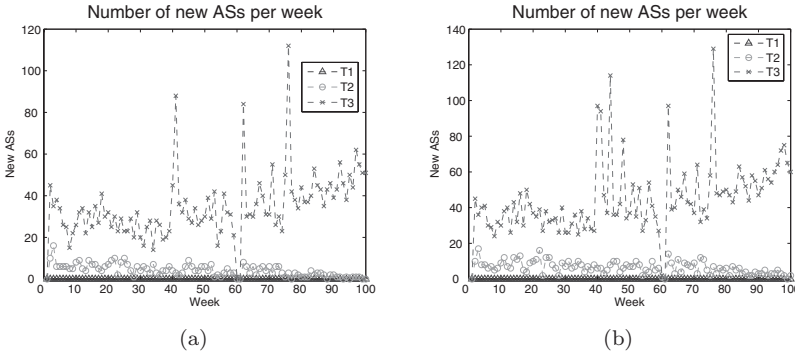


Figure 10. Number of new ASs per week for the first 100 weeks, where “new” means (a) appearing for the first time and (b) not present the week before. (Note that the classification of the nodes into T1, T2, or T3 should be taken seriously only up to week 70 or so; for ASs emerging near the end of the RV data set, we could not as accurately decide on their Kismet.)

ASs, week by week (this means that we implicitly assume that the number of new ASs created per week is (more or less) constant). The data show that such an assumption is not very far from reality (see Figure 10).

- The noninteger values of m_2 and m_3 are averages. They make sense in an approximate analytic discussion (as in the next subsection); in simulations, it is of course not possible to generate a new node with a fractional number of new connections. In this case, the new nodes will have a probabilistic connection strategy (see below).
- Although it is acceptable for a node to connect to another node through more than one link, and this can in fact occur in actual networks to a limited extent, this is not reflected in the RV data set, since the links themselves are not labeled: they are identified only by their two endpoints. The main reason for allowing multiple connections between two nodes in the model here is that excluding it makes the mathematical analysis much harder. For the numbers of nodes and with the probabilities that we will discuss, the number of multiple links in simulations turns out to be negligibly small.
- Under the extremely simple assumptions listed above, and with some reasonable approximations, it turns out that this model can still be discussed analytically. We shall do so below, mainly to see how the differences of this model with the LPAM are reflected in the end result.

- Alternative and arguably more realistic results, resulting from slight modifications of the hypotheses above, will be discussed after the presentation of the main model in its simplest form. In particular, we shall consider the case in which the T1s do not start as a clique and new T1s need not automatically connect to all existing T1s; we shall make the number of potential links that a new T2 or T3 can make a random variable; and we shall look at what happens if multiple connections are prohibited.
- We shall carry out our derivation both for $p_1 = 0$ and for the case where p_1 is allowed to be different from zero (albeit small).

4.2. The Equations

In what follows, N will denote the number of nodes of a particular AS type, and E the degree of a particular node. The superscripts ^[1], ^[2], and ^[3] will denote the AS type (an exception will be made for the parameters p_2 , p_3 , m_2 , and m_3 introduced in Section 4.1, which keep their index as a subscript); within each AS type, the nodes will be numbered sequentially by an integer subscript, according to the order of introduction into the graph. Finally, the letter t will be used to denote time in general, but when a superscript and subscript is used, it will denote the time a particular node of a certain AS type was introduced in the graph: $t_i^{[*]}$ is the time at which the i th node of type $*$ is born.

As in Section 2.1, we will use continuous time models as approximations of discrete time models when this is convenient.

According to the description of the model, the mean number (i.e., the average over a large ensemble of realizations) of nodes of a particular type will evolve in time as

$$\begin{aligned} N^{[*]}(t+1) - N^{[*]}(t) &= p^{[*]} \\ \text{so that } N^{[*]}(t) &= p_* t + N_0^{[*]}, \\ \text{where } * &\text{ stands for 1, 2, or 3} \end{aligned}$$

(take $N_0^{[3]} = N_3(0) = 0$ and $N_0^{[2]} = N_2(0) = 0$, according to the assumptions). Denoting by $S^{[*]}(t)$ the sum, at time t , of all the degrees of nodes of type T*, we get

$$S^{[*]}(t) = \sum_{i=1}^{N^{[*]}(t)} E_i^{[*]}(t),$$

we can write the evolution equation for the degree of each node as follows:

$$\begin{aligned}
E_i^{[1]}(t+1) - E_i^{[1]}(t) &= p_1 + (m_3 p_{31} p_3 + m_2 p_{21} p_2) \frac{E_i^{[1]}(t)}{\mathcal{S}^{[1]}(t)}, \\
&\quad i = 1, \dots, N^{[1]}(t), \\
E_i^{[2]}(t+1) - E_i^{[2]}(t) &= (m_3 p_{32} p_3 + m_2 p_{22} p_2) \frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)}, \\
&\quad i = 1, \dots, N^{[2]}(t), \\
E_i^{[3]}(t+1) - E_i^{[3]}(t) &= 0, \\
&\quad i = 1, \dots, N^{[3]}(t).
\end{aligned} \tag{4.1}$$

We first solve for $\mathcal{S}^{[*]}(t)$. Adding all the equations in (4.1) gives us the difference between the total degrees at time $t+1$ and t , summed over all the nodes *that already existed at time t* ; to obtain $\mathcal{S}^{[*]}(t+1) - \mathcal{S}^{[*]}(t)$, we must add $p_* m_*$, the degree of the newly created AS, to this total. For the case of a T1, m_1 , the number of links made by a new T1 as it is born, depends on time and is given by $N^{[1]}(t)$, the number of T1s that exist at that time. All this leads to

$$\begin{aligned}
\mathcal{S}^{[1]}(t+1) - \mathcal{S}^{[1]}(t) &= 2p_1 N^{[1]}(t) + (m_3 p_{31} p_3 + m_2 p_{21} p_2), \\
\mathcal{S}^{[2]}(t+1) - \mathcal{S}^{[2]}(t) &= p_2 m_2 + m_3 p_{32} p_3 + m_2 p_{22} p_2, \\
\mathcal{S}^{[3]}(t+1) - \mathcal{S}^{[3]}(t) &= p_3 m_3,
\end{aligned} \tag{4.2}$$

which can be solved immediately (use the explicit expression for $N^{[1]}(t)$ derived above):

$$\begin{aligned}
\mathcal{S}^{[1]}(t) &= p_1^2 t^2 + \left(2p_1 N_0^{[1]} + p_{21} p_2 m_2 + p_{31} p_3 m_3 - p_1^2 \right) t + \frac{N_0^{[1]}(N_0^{[1]} - 1)}{2}, \\
\mathcal{S}^{[2]}(t) &= (p_2 m_2 + m_3 p_{32} p_3 + m_2 p_{22} p_2) t, \\
\mathcal{S}^{[3]}(t) &= p_3 m_3 t,
\end{aligned} \tag{4.3}$$

where we have taken into account the fact that at $t = 0$ the graph is a clique of $N_0^{[1]}$ nodes of type T1.

For the time being, we will keep p_1 in our equations, even though the observed value for p_1 in the RV dataset equals 0; we can, of course, always set $p_1 = 0$ again in the final solution. Substituting (4.3) back into (4.1) leads to

$$\begin{aligned}
E_i^{[2]}(t+1) - E_i^{[2]}(t) &= E_i^{[2]}(t) \frac{F_2}{t}, \\
E_i^{[1]}(t+1) - E_i^{[1]}(t) &= p_1 + \frac{K}{p_1^2 t^2 + at + b} E_i^{[1]}(t),
\end{aligned} \tag{4.4}$$

where

$$\begin{aligned} F_2 &= \frac{m_3 p_{32} p_3 + m_2 p_{22} p_2}{p_2 m_2 + m_3 p_{32} p_3 + m_2 p_{22} p_2}, \\ K &= m_3 p_{31} p_3 + m_2 p_{21} p_2, \\ a &= 2p_1 N_0^{[1]} + p_{21} p_2 m_2 + p_{31} p_3 m_3 - p_1^2, \\ b &= \frac{N_0^{[1]}(N_0^{[1]} - 1)}{2}. \end{aligned}$$

The equation for the degrees of T2s can be (approximatively) solved as follows:

$$\begin{aligned} E_i^{[2]}(t) &= \left[\prod_{\ell=t_i^{[2]}}^{t-1} \left(1 + \frac{F_2}{\ell} \right) \right] E_i^{[2]}(t_i^{[2]}) \\ &= \exp \left[\sum_{\ell=t_i^{[2]}}^{t-1} \log \left(1 + \frac{F_2}{\ell} \right) \right] E_i^{[2]}(t_i^{[2]}) \end{aligned} \quad (4.5)$$

$$\begin{aligned} &\approx \exp \left[\sum_{\ell=t_i^{[2]}}^{t-1} \frac{F_2}{\ell} \right] E_i^{[2]}(t_i^{[2]}) \\ &\approx \exp \left[F_2 \log \left(\frac{t}{t_i^{[2]}} \right) \right] E_i^{[2]}(t_i^{[2]}) = m_2 \left(\frac{t}{t_i^{[2]}} \right)^{F_2}. \end{aligned} \quad (4.6)$$

The approximations made are similar to the substitution of an integral for a sum in Section 2.1; in fact (4.5) is the solution that we would have found if we had replaced the difference equation for $E_i^{[2]}$ in (4.4) by a differential equation. One can derive an approximate solution for the degrees of T1s in a similar way; here we show how to solve the corresponding differential equation, i.e., the equation

$$\frac{dE_i^{[1]}}{dt} = p_1 + \frac{K}{p_1^2 t^2 + at + b} E_i^{[1]}(t).$$

For the special case where $p_1 = 0$, this is straightforward to solve; in this case we have moreover $a = K$, so that the solution is

$$E_i^{[1]}(t) = \left(1 + \frac{Kt}{b} \right) (N_0^{[1]} - 1).$$

When $p_1 \neq 0$, the solution is given by

$$E_i^{[1]}(t) = \left(\frac{t - t_+}{t_i^{[1]} - t_+} \right)^{K/\delta} \left(\frac{t_i^{[1]} - t_-}{t - t_-} \right)^{K/\delta} [N_i^{[1]}(t_i^{[1]}) - 1] \\ + p_1 \left(\frac{t - t_+}{t - t_-} \right)^{K/\delta} \int_{t_i^{[1]}}^t \left(\frac{t' - t_-}{t' - t_+} \right)^{K/\delta} dt' ,$$

$$\text{where } \delta = \sqrt{a^2 - 4bp_1^2},$$

$$t_+ = \frac{-a + \delta}{2p_1^2},$$

$$t_- = \frac{-a - \delta}{2p_1^2}, \text{ and we assume that}$$

$$t > t_i^{[1]}.$$

(Note that we implicitly assume that $a^2 > 4bp_1^2$; since p_1 is supposed to be very small, this is a reasonable assumption.) Under the assumption that p_1 is small, we can approximate this to keep only the highest order terms. We have then

$$t_+ = -\frac{b}{a} + O(p_1^2), \\ t_- = -\frac{a}{p_1^2} (1 + O(p_1^2)), \\ \frac{K}{\delta} = 1 - 2p_1 \frac{N_0^{[1]}}{K} + O(p_1^2),$$

and for times t that are at most of the order $O(p_1^{-2+\epsilon})$, the solution reduces to

$$E_i^{[1]}(t) = \begin{cases} \left(1 + \frac{Kt}{b}\right) (N_0^{[1]} - 1) + p_1 \left(t + \frac{b}{K}\right) \ln \left(\frac{K}{b}t + 1\right) & \text{if } i \leq N_0 \\ \left(\frac{Kt+b}{Kt_i^{[1]}+b}\right) (N_0^{[1]} + p_1 t_i^{[1]}) + p_1 \left(t + \frac{b}{K}\right) \ln \left(\frac{Kt+b}{Kt_i^{[1]}+b}\right) & \text{if } i > N_0^{[1]} \text{ and } t > t_i^{[1]}. \end{cases}$$

Finally, we also have

$$E_i^{[3]}(t) = m_3, \quad t > t_i^{[3]}. \quad (4.7)$$

4.3. Computation of the Degree Distribution

The cumulative node degree probability distribution at time t is

$$\begin{aligned}
 \mathbf{P}(E < x) &= \mathbf{P}(E < x \mid T1) \frac{p_1 t + N_0^{[1]}}{t + N_0^{[1]}} + \mathbf{P}(E < x \mid T2) \frac{p_2 t}{t + N_0^{[1]}} \\
 &\quad + \mathbf{P}(E < x \mid T3) \frac{p_3 t}{t + N_0^{[1]}} \\
 &\approx \left(p_1 + \frac{N_0^{[1]}}{t} \right) \mathbf{P}(E^{[1]} < x) + p_2 \mathbf{P}(E^{[2]} < x) \\
 &\quad + p_3 \mathbf{P}(E^{[3]} < x).
 \end{aligned} \tag{4.8}$$

In this formula, we just decomposed the total probability into probabilities conditioned on type; for each type the probability that the node picked at random at time t is of type T^* is given by the ratio between $N^{[*]}(t)$ and the total number of nodes at that time, $N_0^{[1]} + t$. To simplify the expressions, we have assumed that we consider $t \gg N_0^{[1]}$. It remains to compute these conditional probabilities, and for this, we will use that the quantities $\{t_i^{[1]}\}$, $\{t_i^{[2]}\}$, and $\{t_i^{[3]}\}$ are uniformly distributed within $(0, t)$. This is a reasonable assumption for all three sequences if $\min(p_1, p_2, p_3)t = p_1 t \gg 1$; if $t \gg 1$, but $p_1 t$ isn't, then it is only approximately true for the $t_i^{[1]}$. For simplicity, let us assume that $p_1 t = O(p_1^{-1+\epsilon}) \gg 1$; in that case, we can neglect even the $N_0^{[1]}/t$ contribution to the T1 term. Then,

$$\begin{aligned}
 \mathbf{P}(E^{[2]} < x) &= \mathbf{P} \left(m_2 \left(\frac{t}{t^{[2]}} \right)^{F_2} < x \right) = \mathbf{P} \left(t^{[2]} > t \left(\frac{m_2}{x} \right)^{1/F_2} \right) \\
 &= 1 - \left(\frac{m_2}{x} \right)^{1/F_2},
 \end{aligned} \tag{4.9}$$

$$\begin{aligned}
 \mathbf{P}(E^{[1]} < x) &\approx \mathbf{P} \left(p_1 + \frac{N_0^{[1]}}{t^{[1]}} + p_1 \ln \left(\frac{t}{t^{[1]}} \right) < \frac{x}{t} \right) \\
 &= \mathbf{P} \left(\frac{N_0^{[1]}}{p_1 t^{[1]}} - \ln(p_1 t^{[1]}) < \frac{x}{p_1 t} - 1 - \ln(p_1 t) \right) \\
 &\approx \mathbf{P} \left(\ln(p_1 t^{[1]}) > 1 - \frac{x}{p_1 t} + \ln(p_1 t) \right) = \mathbf{P} \left(p_1 t^{[1]} > p_1 t e^{1-x/(p_1 t)} \right) \\
 &= 1 - \exp \left(\frac{p_1 t - x}{p_1 t} \right),
 \end{aligned} \tag{4.10}$$

$$\mathbf{P}(E^{[3]} < x) = 1_{\{x \geq m_3\}}. \tag{4.11}$$

(Note that we have implicitly assumed that $x > p_1 t$ in the approximation for the T1-probability.) Substituting the equations above into (4.8), we obtain

$$\begin{aligned} \mathbf{P}(E < x) &\approx p_1 \left(1 - e^{-\frac{x-p_1 t}{p_1 t}}\right) \mathbf{1}_{\{x > p_1 t\}}(x) \\ &\quad + p_2 \left(1 - \left(\frac{m_2}{x}\right)^{1/F_2}\right) \mathbf{1}_{\{x \geq m_2\}}(x) \\ &\quad + p_3 \mathbf{1}_{\{x \geq m_3\}}(x). \end{aligned} \quad (4.12)$$

This derivation was carried out in the assumption that p_1 is small, but not zero, and that t is sufficiently large (of order $O(p_1^{-2+\epsilon})$) to justify the approximations. When $p_1 = 0$, the derivation is simpler, and we obtain

$$\begin{aligned} \mathbf{P}(E < x) &= \frac{N_0^{[1]}}{t + N_0^{[1]}} \mathbf{1}_{\{x > (b+Kt)/[b(N_{0,1}-1)]\}}(x) \\ &\quad + \frac{p_2 t}{t + N_0^{[1]}} \left(1 - \left(\frac{m_2}{x}\right)^{1/F_2}\right) \mathbf{1}_{\{x \geq m_2\}}(x) + \frac{p_3 t}{t + N_0^{[1]}} \mathbf{1}_{\{x \geq m_3\}}(x) \\ &\approx p_2 \left(1 - \left(\frac{m_2}{x}\right)^{1/F_2}\right) \mathbf{1}_{\{x \geq m_2\}}(x) + p_3 \mathbf{1}_{\{x \geq m_3\}}(x) \quad \text{for large } t. \end{aligned} \quad (4.13)$$

4.4. A First Discussion of the Revised Preferential Attachment Model

Although the computations above are only approximations, albeit to roughly the same extent as the estimate in Section 2.1, we can at this point make a first evaluation of the RPAM.

- The distributions, contrary to the case in [Barabási and Albert 99], are *not* completely independent of time. Some of this time dependence (in the case where $p_1 = 0$, all of it) is simply due to the initial condition, in which the different types are not distributed according to the probabilities with which they are created later on; this causes a transient time dependence at small t , which dies out as t increases. In the case $p_1 \neq 0$, there is some time dependence even for large t : the T1s contribute something like a probability wave, moving with constant speed; it has a very small amplitude, however, proportional to p_1 . This is probably an artifact of the extreme simplifications in our derivations rather than a true “physical” phenomenon.
- To obtain the degree distribution, we must take the derivative of the cumulative probability distribution derived in Section 4.3. If we disregard

the T1s, then we find again a power distribution, with exponent

$$\begin{aligned}
 \alpha &= 1 + F_2^{-1} \\
 &= 1 + \frac{m_3 p_{32} p_3 + m_2 p_2 (p_{22} + 1)}{m_3 p_{32} p_3 + m_2 p_{22} p_2} \\
 &= 2 + \frac{m_2 p_2}{m_3 p_{32} p_3 + m_2 p_{22} p_2}.
 \end{aligned} \tag{4.14}$$

With the values for these parameters that were proposed in Section 4.1, this leads to $\alpha \approx 2.39$, which accords reasonably well with the observed power law for the middle region (see Figures 1(a) and (b), which show slopes of 2.24 and 2.3, respectively). At this point, with only this “back-of-the-envelope” calculation to go by, we shouldn’t read too much into the correspondence; on the other hand, it is interesting that the RPAM mechanism has generated a power law with a slope different from 3 and close to the value observed in the degree distribution of the AS graph, when we use the attachment parameters derived from the (different) observation of the temporal behavior of the degrees of different AS classes.

- The LPAM had the undesirable feature that the degree of a node was expected to be directly proportional to its age—something definitely not observed in practice. The RPAM no longer has this feature: the existence of different types, with their own attachment policies, now makes it possible for later-born T2s to exceed in degree earlier-born T3s.
- The RPAM still has some of the (rather unrealistic) assumptions of the LPAM, such as the following:
 - Multiple connections between two nodes are allowed. We discuss briefly in Section 4.5.2 how this could be corrected; moreover, simulations prove that it does not affect the results, at least for the parameter values of interest (see Figures 11 and 14).
 - We have set the initial degrees of the nodes, at birth, to be deterministic and constant. It is straightforward to randomize m_3 : if m_3 takes the value n_i with probability q_i , then (4.11) becomes $\mathbf{P}(E^{[3]} < x) = \sum_{\{i: n_i < x\}} q_i \mathbf{1}_{\{x \geq n_i\}}(x)$. Randomization of m_2 , however, will make the equations too involved for even approximate analytic treatment.
 - This model still does not address the issues of node and edge dropout. An indirect way to incorporate these features in RPAM is by adjusting the total number of nodes, the initial node degrees, and the interconnection probabilities accordingly.

- The model still does not address any dynamics of already established connections or new connections between two preexisting nodes.

4.5. Some Extensions of the Computation

This section will discuss some possible modifications of the original RPAM assumptions, which were set in Section 4.1. As expected, these can quickly make the model analytically intractable.

4.5.1. The Transits do not form a clique. Suppose that a newly introduced T1 does not connect to every existing T1 but rather to a fraction λ of them; assume also that the original T1s do not form a clique either but that only a fraction β of all the edges of the complete graph connecting them is present. Then, the resulting model can still be solved, by following step by step the derivation above. One finds

$$\mathcal{S}^1(t) = \lambda p_1^2 t^2 + (\lambda p_1 N_0^{[1]} + m_3 p_3 p_{31} + m_2 p_2 p_{21} - \lambda p_1^2) t + \beta \frac{N_0^{[1]}(N_0^{[1]} - 1)}{2}. \quad (4.15)$$

This results in a modified differential equation for $E_i^{[1]}(t)$:

$$\frac{d}{dt} E_i^{[1]} = p_1 \lambda + \frac{\widehat{K}}{p_1^2 t^2 + \hat{a} t + \hat{b}} E_i^{[1]}, \quad (4.16)$$

where

$$\begin{aligned} \widehat{K} &= \frac{m_3 p_3 p_{31} + m_2 p_2 p_{21}}{\lambda}, \\ \hat{a} &= \frac{\lambda p_1 N_0^{[1]} - \lambda p_1^2 + m_3 p_3 p_{31} + m_2 p_2 p_{21}}{\lambda}, \\ \hat{b} &= \beta \frac{N_0^{[1]}(N_0^{[1]} - 1)}{\lambda}. \end{aligned}$$

The remainder of the derivation then proceeds as before, with the earlier parameters replaced by their adjusted versions.

4.5.2. At most one link between two nodes is allowed. Instead of allowing multiple links between two nodes and having a fixed number m_* of edges be initiated by each new node of type T*, assume that a newly introduced node of type T* will make m_* attempts to form a connection with existing nodes, following the PA principle, but that an attempt to link to a node succeeds only if there is not yet a connection with that other node; a failed attempt, i.e., an attempt by a new node to link with an old node to which it already has a connection, results in the new node losing that linking opportunity. It follows that at the end of the

round in which it is introduced, a T2 (or T3) node may end up with fewer than m_2 (or m_3) connections; in fact, its number of connections (lying between 1 and m_2 (or m_3)) will be given by the difference between m_2 and the number of its failed attempts. The successive linking attempts are, however, still independent, so that the probability for a new T* node to establish a connection with an existing old node with its k th attempt is given by $(1-p)^{k-1}p$, where p stands for the probability of establishing this link with the first of the new node's m_* attempts. The equivalent of equation (4.1) for T2 nodes thus becomes

$$\begin{aligned} E_i^{[2]}(t+1) - E_i^{[2]}(t) &= \frac{E_i^{[2]}(t)}{\mathcal{S}_2(t)} \left(p_2 p_{22} \sum_{j=1}^{m_2} \left(1 - \frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)} \right)^{j-1} + p_3 p_{32} \sum_{j=1}^{m_3} \left(1 - \frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)} \right)^{j-1} \right) \\ &= p_2 p_{22} \left(1 - \left(1 - \frac{E_i^{[2]}(t)}{\mathcal{S}_2(t)} \right)^{m_2} \right) + p_3 p_{32} \left(1 - \left(1 - \frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)} \right)^{m_3} \right). \end{aligned} \quad (4.17)$$

Note that, unlike the simple case discussed earlier, this derivation makes sense only if m_2 and m_3 are integers. These equations will be very close to (4.1) if $\frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)} \ll 1$ for all i , because $(1+x)^n \approx 1-nx$ if $x \ll 1$; on the other hand, they will be very different if the graph contains nodes that have attracted a considerable fraction of the edges of the graph.

If $m_2 = m_3 = 2$ (i.e., the smallest nontrivial values for this model), then (4.17) becomes

$$E_i^{[2]}(t+1) - E_i^{[2]}(t) = 2(p_2 p_{22} + p_3 p_{32}) \frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)} - (p_2 p_{22} + p_3 p_{32}) \frac{\left(E_i^{[2]}(t)\right)^2}{\mathcal{S}^{[2]}(t)^2}.$$

Thus, we recovered (4.1) plus a corrective higher-order term. By summing on all T2, we obtain an equation for $\mathcal{S}^{[2]}(t) = \sum_{i=1}^{N^{[2]}(t)} E_i^{[2]}(t)$:

$$\begin{aligned} \mathcal{S}^{[2]}(t+1) - \mathcal{S}^{[2]}(t) &= p_2 \left(2 - \sum_{i=1}^{N^{[2]}(t)} \left[\frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)} \right]^2 \right) + 2(p_2 p_{22} + p_3 p_{32}) \\ &\quad - (p_2 p_{22} + p_3 p_{32}) \frac{\sum_{i=1}^{N^{[2]}(t)} \left(E_i^{[2]}(t)\right)^2}{\mathcal{S}^{[2]}(t)^2} \\ &= 2(P + p_2) - (P + p_2) \sum_{i=1}^{N^{[2]}(t)} \left[\frac{E_i^{[2]}(t)}{\mathcal{S}^{[2]}(t)} \right]^2, \end{aligned} \quad (4.18)$$

with $P = p_2 p_{22} + p_3 p_{32}$;

the quantity between large brackets multiplying p_2 (in the first term of the right-hand side of the first equation) is the average degree (in the ensemble sense) acquired at birth by a new T2 in this model. Let us estimate the impact of the corrective term

$$(P + p_2) \sum_{i=1}^{N^{[2]}(t)} \left(E_i^{[2]}(t) / \mathcal{S}^{[2]}(t) \right)^2.$$

If u_1, \dots, u_N are all positive and constrained by $\sum_{i=1}^N u_i = 1$, then the maximum and minimum values of $\sum_{i=1}^N u_i^2$ are given by respectively 1 (if all but one of the u_i equal 0) and N^{-1} (if all the u_i equal N^{-1}). It follows that the corrective term in (4.18) takes values between $(P + p_2)/N^{[2]}(t)$ and $(P + p_2)$.

This term is therefore smaller than the first, dominant term. This agrees with what one would expect intuitively: since the corrective term represents the (small) probability that a T2 node gets hit twice by the same new node born at t , it will have a small effect only, increasingly so as t becomes large, and the likelihood of a double hit decreases. If we compute its effect by perturbation analysis, then its first-order effect is to multiply the zeroth-order estimate for $E_i^{[2]}(t) \approx (t/t_i^{[2]})^{P/(P+p_2)}$ by $\left(1 + O((t/t_i^{[2]})^{-\min(p_2, P)/(p_2+P)})\right)$, a correction that does indeed diminish in importance as t becomes large. A similar analysis holds if we don't restrict ourselves to $m_2 = m_3 = 2$.

Thus, we expect that restricting the model by prohibiting multiple links will not influence the power law behavior or its exponent; this is also borne out by our simulations (see Section 5).

4.5.3. Nonlinear PA. So far we have considered *Linear* PA only, in which the probability for an old node to attract a connection from a new node is directly (linearly) proportional to this old node's degree. The Empirical Attachment Probability curves that we found (see Figures 6 and 4), as well as related discussions in the literature (see [Chen et al. 02]), suggest that a nonlinear PA might be closer to reality. The attachment probability could, for instance, be given by a nonlinear increasing function $(W) : \mathbb{R}_+ \rightarrow \mathbb{R}_+$, such as, for instance, $\mathcal{W}(x) = x^\alpha, \alpha > 0$. The equations (4.1) then become

$$\begin{aligned} \frac{dE_i^{[1]}(t)}{dt} &= p_1 + (m_3 p_{31} p_3 + m_2 p_{21} p_2) \mathcal{W}\left(\frac{E_i^{[1]}(t)}{\sum_{j=1}^{N^{[1]}(t)} E_j^{[1]}(t)}\right), \quad i = 1, \dots, N^{[1]}(t), \\ \frac{dE_i^{[2]}(t)}{dt} &= (m_3 p_{32} p_3 + m_2 p_{22} p_2) \mathcal{W}\left(\frac{E_i^{[2]}(t)}{\sum_{j=1}^{N^{[2]}(t)} E_j^{[2]}(t)}\right), \quad i = 1, \dots, N^{[2]}(t), \\ \frac{dE_i^{[3]}(t)}{dt} &= 0, \quad i = 1, \dots, N^{[3]}(t). \end{aligned}$$

However, the summation trick that allowed us to eliminate the denominators in the original analysis will not work for a general function \mathcal{W} , making an approximate analytic model much less tractable.

5. Simulations

Having explored the analytical models derived from drastic simplifications of the RPAM, we now turn to simulations of a richer version. In this richer version we shall in particular

- experiment with allowing or prohibiting multiple connections between a new AS and older AS, by carrying out simulations in both cases (leaving all the other simulation characteristics identical) to compare the results;
- let the initial degree of a T2 or T3 node be a random variable (with a PDF lifted from our observations in the RV data set) rather than a fixed number;
- extend the model to incorporate the (observed) dropping out of ASs and check the effect of this on the results;
- check whether allowing the amount of connecting to T3s that was observed in RV (but disregarded in our simple analytic model) affects the simulations;
- investigate to what extent the (crude) approximation to the Empirical Attachment Probability (EAP), which we basically assumed to be piecewise LPA in order to build our simple model, plays a role, by comparing simulations with this simple model with simulations run with the observed EAP.

5.1. An Initialization Issue

Although our analysis in Sections 4.2 and 4.3 assumed that $N_0^{[2]} = N_0^{[3]} = 0$, this was mainly for the sake of simplicity; one easily checks that setting the initial populations of the T2s and T3s to nonzero levels doesn't affect the asymptotic behavior (for large t) of the solutions to our equations. It is important, however, to bear in mind that the equations describe only *mean* behavior, in the ensemble sense. Depending on the initial populations, one finds in practice that the variation among different realizations of the solution depends heavily on the initial conditions, so that many more simulations need to be run for some initial conditions than for others, in order to capture the mean behavior with some

confidence. For instance, if one starts out with no T2s at all at $t = 0$, and the simulation happens to produce several T3s in a row after the first T2 is born, then this T2 will already have such a head start in degree on its later-born siblings that it will continue to attract a lion's share of the new connections to T2s at every later time. Such a large deviation from mean behavior is much less likely to happen when there are already several T2s at the start, so we then typically need many fewer simulations (to the extent that virtually every simulation already shows "typical" behavior) than when $N_0^{[2]} = 0$.

For this reason, we start from initial conditions in which we have not only a number of T1 nodes (forming a clique or not) but also several T2 nodes, connected among each other and with the T1s in a fairly homogeneous way. These T2s are connected in such a way as to form a ring, while each one connects additionally to exactly one T1; the degree of the initially present T2 is thus three. The initial T1s share the connections from the ring of T2s as uniformly as possible; this is made easy by choosing the number of initial T2s to be a multiple of the number of initial T1s.

It turns out that the exact parameters in this initialization or other corrections are not important: simulations (not shown here) demonstrate that the results are robust.

5.2. Simulation Results for the Simplest Case

Our first simulations closely correspond to the simplistic RPAM of Section 4.2. The initial condition has a clique of four T1s; in addition, it has eight T2s, connected as just described. Every new node as it is born is either a T2 or T3 (no new T1s), with the probabilities p_2 and p_3 ; it connects to preexisting T2s and/or T1s according to the probabilities p_{22} , p_{21} , p_{32} , or p_{31} , respectively, where we use for all these parameters the values determined in Section 3.8. Within each category T* of targets for a connection from a new node, the different possible targets for a connection from a new node "attract" that connection with the probability given by Linear Preferential Attachment, as explained earlier in the RPAM model setup. We allow a new node to connect multiple times to the same preexisting node, if that is how the random "dice" roll. In this simplistic version of the model, we are assuming that no new node connects to an existing T3 node, so that every T3 node keeps its birth degree for all later times. For this reason, we let the T3 nodes be born, in this simulation, with the distribution of their maximal degree shown in RV, which can be modeled well by a Poisson distribution with mean 2. For the initial degree distribution of the T2s, we mimic the degree distribution of T2s at birth observed in RV; to a first approximation this is an exponential distribution.

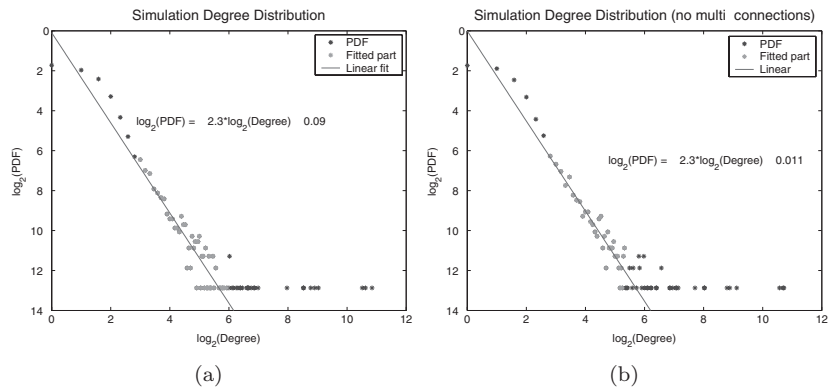


Figure 11. The degree distribution for the RPAM simulations, where multiple connections are (a) allowed and (b) not allowed. Observe the similarity with Figure 1.

The results of a typical simulation are shown in Figure 11(a). Comparing with Figure 1, we see qualitatively that the simulations are almost identical to reality and quantitatively that the simulations capture very accurately the decay exponent of the actual PDF: the RV data set provided an estimated value of 2.24 to 2.3, whereas the simulations give 2.3. In general, there is an uncertainty about the value of this exponent caused by two reasons: (a) the portion of the PDF representing the tail is chosen quite arbitrarily, and (b) the simulation is inherently random. Over a large number of simulations, we observed the decay exponent to vary between 2.1 and 2.4.

The theoretical prediction (4.14) for the decay exponent of the degree PDF of RPAM simulations for the values of the parameters chosen is 2.3. (Note that we substituted $m_3 = 2$ in the formula, rather than the mean value $m_3 = 1.5$, in accordance with the argument in the first paragraph of this section.) This value is close to both the simulation results and the exponent estimated from RV.

5.3. Forbidding Multiple Connections

The first enrichment that we consider is to forbid multiple connections of a new node to the same preexisting node. Otherwise, the simulation conditions are identical to the ones described in Section 5.2. As shown in Figure 11(b), simulations with this extra restriction give a result that is virtually identical to the case where multiple connections were allowed. This result also holds for all the other simulations that we carried out for more enriched models. Moreover, it appears that allowing multiple connections or not (see Section 4.5.2) does not really make a difference.

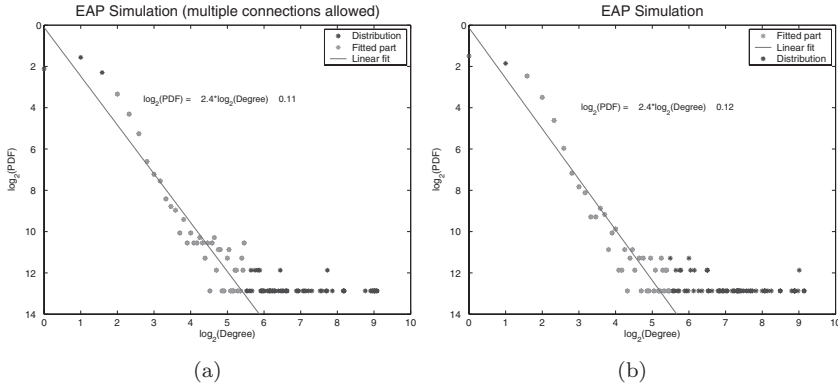


Figure 12. The degree distribution for the RPAM simulations, using the EAP (shown in Figure 4), where multiple connections are (a) allowed and (b) not allowed. Observe the similarity with Figures 1 and 11.

5.4. Using the Full Empirical Attachment Probability

In Section 3.5 we showed the Empirical Attachment Probability derived from the data. We argued that it was, in first approximation, consistent with a classification of ASs in three different classes, each of which observed its own Linear Attachment Probability strategy for connections to the three classes; this was the basis for the simple RPAM model that we then discussed in Section 4. In this subsection we shall not make this approximation: instead, we simulate a process in which nodes are born with the distribution (for their maximum degree, consistent with the Kismet assumption—see Section 3.7.1) that we observed in RV, and we let them connect with preexisting nodes according to the Empirical Attachment Probability observed in RV. (In this simulation, there is therefore no classification into three different types, neither for the newborn ASs nor for the targets of their connections.) The initial degree of a newborn AS follows the observed PDF for the initial degree of ASs of that Kismet. Figure 12 shows the result. Even though we have now used a much richer model, the results are very close to that of the simplest model, both qualitatively (shape of the degree distribution curve) and quantitatively (exponent of the power decay law for the PDF).

Note that this simulation effectively uses a continuum of different AS classes. We also carried out simulations in which we considered a finite number (but larger than three) of AS classes. We do not present these in detail here; they too gave results that are very close to the simplistic model, which is not surprising, given that the two extreme cases (three classes only and an infinite number of classes) are already so close.

5.5. AS Dropouts and Their Impact

Contrary to the RPAM assumptions, RV shows that ASs can disappear again after birth. In some cases the dropout appears to be permanent, in others temporary, with the AS reappearing after some time of absence (see Figure 3(d) for an example). In this subsection we want to explore what the effect is of allowing dropouts, to the extent observed in RV, in simulations of RPAM.

We can get a good estimate of how extensive AS dropout is by computing how many new ASs have been introduced into the graph since the beginning of our data set measurements and up to a given week and then computing what percentage of them is still present in this week. This percentage is plotted in

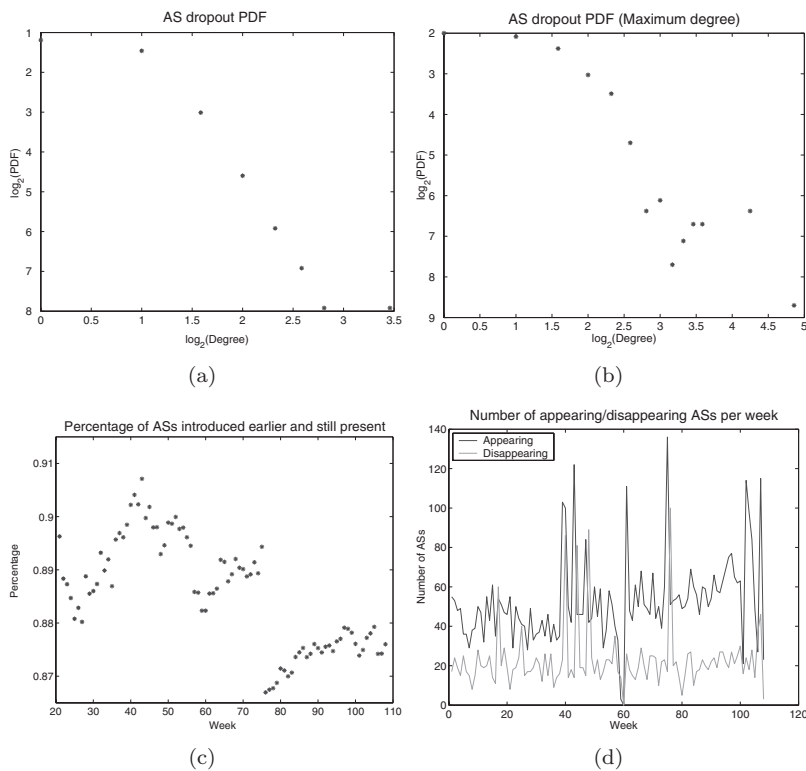


Figure 13. Dropout statistics: (a) shows the degree distribution of the dropout ASs the week immediately before dropping out; (b) shows the distribution of the maximum degree of the dropout ASs; (c) shows the percentage of the ASs which were introduced earlier than a given week, but within the data set period, and are still present in this week; (d) shows the number of appearing and disappearing ASs in every week.

Figure 13(c) (we omit the first 20 weeks); we see that it is fairly stable (despite an abrupt downward transition around week 75, which also appears as a sudden surge in dropouts in Figure 13(d), and that on the average 12% of the ASs introduced in the graph drop out some time later. Figure 13(a) shows the probability distribution, among ASs that drop out, of their degree immediately before dropping out, whereas Figure 13(b) shows this probability distribution according to the maximum degree of the dropout AS; the two distributions are slightly different, but in both cases ASs of small degree tend to drop out more frequently, which is precisely what one expects. Figure 13(d) shows the introduction and dropout of ASs in absolute numbers. (See also [Chang et al. 04, Chen et al. 02] for similar observations and a more extensive discussion.)

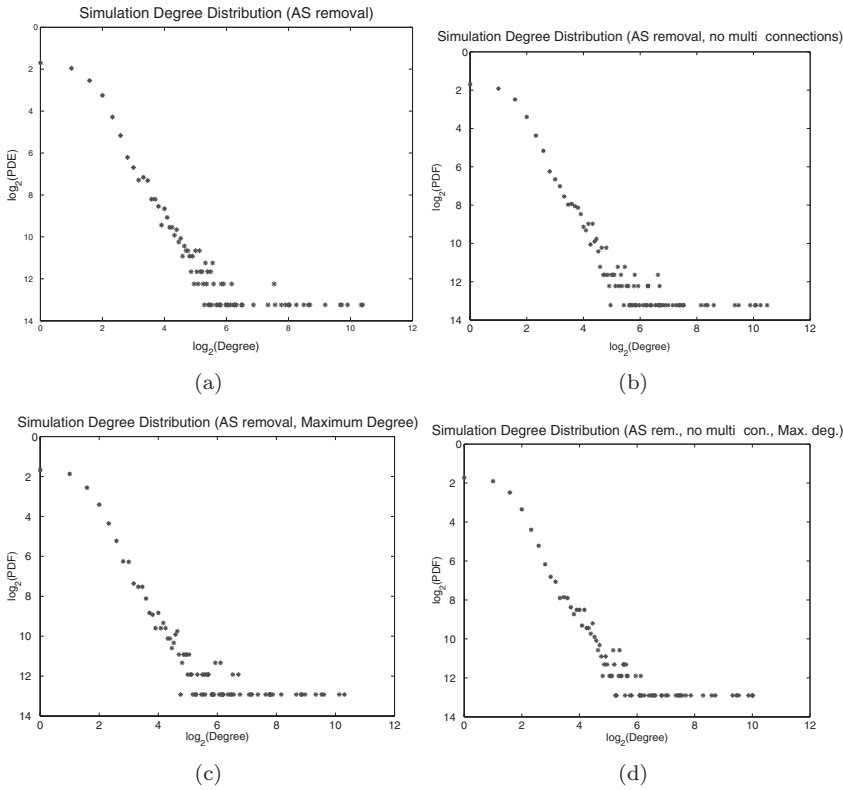


Figure 14. The degree distribution for the RPAM simulations with AS dropout. Top row: AS dropout according to degree immediately before the dropout (Figure 13(a)). Bottom row: AS dropout according to maximum degree (Figure 13(b)). Right (left) column: multiple connections (not) allowed. In all cases, observe the similarity with Figure 1.

In its original version, RPAM has no dropout mechanism at all. To run simulations with dropout, we shall randomly remove ASs in each week. The percentage of nodes that we remove is the average percentage per week that was observed in RV; we distribute this number according to the observed dropout distribution. Since we took note of two different distributions for the dropout nodes, one according to the degree just before dropout and the other according to maximum degree, we have a choice of two strategies. In both strategies, every time step consists of two stages: in the first stage, we run (for one time step) the simulation as before (without dropout), starting from the AS population achieved in the previous time step; in the second stage, we remove some AS nodes, consistent with

- (for strategy a) the distribution of dropout AS degrees immediately before dropouts, shown in Figure 13(a);
- (for strategy b) the distribution of the dropout AS maximum degrees, shown in Figure 13(b);

the remaining ASs then make up the simulated population for this time step. The results for both strategies are shown in Figure 14; for each of the two strategies, we ran simulations both without and with a prohibition against multiple connections between the same ASs; both are shown in Figure 14. We find that the two strategies (with or without multiple connections) lead to results that are very close to each other as well as to the observed degree distribution of Figure 1. From this we conclude that possible dropout of ASs, at least to the limited extent observed in RV, does not have a significant impact on the degree distribution of the ASs.

Remark 5.1.

1. One can also follow two different strategies for the attachment probability distribution at each time step: for each AS that survives the second stage, one can take the attachment probability as it would have been computed in the standard simulation (i.e., at the end of the first stage), with a simple renormalization so that the probabilities add to 1, or one can, on the contrary, first remove ASs, then reassess the numbers of AS present for the different types, and determine accordingly the attachment probabilities for all the existing ASs at that point. In Figure 14 the first of these strategies was followed; we also carried out the simulation with the other strategy, with virtually identical results (not shown here).

2. Note that we don't claim to have simulated the complex time evolution behavior of the degree of individual ASs, which can decrease as well as increase and even decrease to the extent that the AS can disappear. Instead, we have taken note of the fact that at any time only a small number of ASs seem to have such aberrant behavior, and we have tried to evaluate the effect of this more complex behavior by removing ASs "by hand" at every time step, as described above, and checking the effect of this "deus ex machina" intervention on the degree distribution in our simulations, as it evolves in time. At first order, the dropout of ASs does not seem to have any effect in this admittedly crude experiment; however, we take this observation as an indication that a more careful simulation of the actual degree oscillations would not have an impact either.

5.6. Conclusions from Our Simulations of Degree Distributions

The RPAM model is a drastic simplification of a very complex system, in which individual nodes have a more complex, richer connectivity behavior and degree evolution than RPAM allows. Nevertheless, this simple model seems already to capture the main structure of the observed degree distribution in the AS network observed in RV. In this section, we have simulated various "complicating factors" that do occur in the real-life AS network; our simulations showed that although these factors no doubt are indications of additional layers of complexity in the dynamic behavior of the network, they don't seem to have any impact on the coarse features that we observed in the degree distribution and empirical attachment probability distribution, i.e., the classification into three different types, each following (approximately) its own Linear Preferential Attachment strategy with respect to each of the classes and with "class attachment probabilities" given by Table 2. Despite its sweeping simplifications, the simple RPAM model seems already to capture the main structure of the observed degree distribution in the AS network observed in RV. In what follows we shall therefore stick to the simple model only.

6. Influence of the Main RPAM Parameters on the Behavior of the Degree Distribution

The degree distribution of simulations produced by RPAM is remarkably robust to the small changes in the model discussed in the previous section. In this section we discuss briefly how the degree distribution of the simulations changes if we modify the values of some of the parameters (m_2 , m_3 , p_2 , p_3 , p_{22} , p_{23} , p_{32} ,

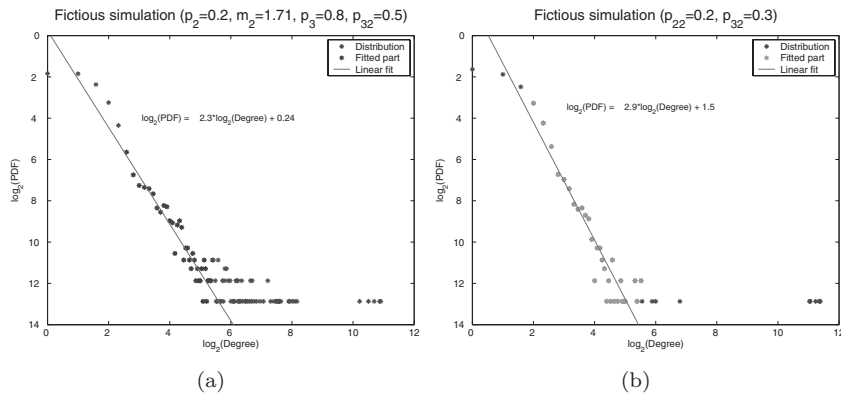


Figure 15. Two simulations with parameters different from those observed (see text).

and p_{33}). Formula (4.14) gives us an approximate formula for the slope of the “main trunk” in the log-log plot of the degree distribution; we saw earlier that for the values of the parameters corresponding to the RV observations, this slope does indeed correspond to the slope of the observed log-log degree distribution. Figure 15 shows the degree distribution and the slope α for two other choices of the parameters (which do not correspond to reality). All parameters remain the same as in the previous simulations, except the ones mentioned in the title of each figure.

In Figure 15(a), we take $p_2 = 0.2$, $m_2 = 1.71$, $p_3 = 0.8$, and $p_{32} = 0.5$; because the ratio $(m_2 p_2)/(m_3 p_3 p_{32} + m_2 p_2 p_{22})$ is the same as before, these new parameter values still lead to $\alpha = 2.3$, which is borne out by the simulation that led to Figure 15(a), which shows indeed the same slope $\alpha = 2.3$ as before. In Figure 15(b), we take $p_{22} = 0.2$ and $p_{32} = 0.3$; substituting this new choice into (4.14) leads to $\alpha = 2.9$, which does again correspond to the slope of the log-log degree distribution in simulations (Figure 15(b) gives $\alpha = 2.9$).

These experiments serve both as a validation of the back-of-the-envelope derivation (4.14) and as reassurance that the log-log degree distribution curve *can* be budged, so that the robustness of the simulated log-log degree distribution curve under all the perturbations in the previous section is indeed meaningful.

7. Comparing Other Features of the RPAM Simulations with the RV Data

The RPAM model incorporates data only from the (average) connecting behavior of ASs as they emerge in the network; it uses these to grow, from fairly small

initial graphs, simulated graphs that show a degree distribution similar to the one observed in reality. This similarity is moreover very robust, in the sense that the end result, after many (typically between 7,000 and 10,000) time steps, is the same for a wide range of initial conditions, with some limitations (see Section 5.1). Provided there is a sufficient number of T2s in our initial state, we find that we obtain, after the preset number of steps, a degree distribution very similar to Figure 1, even if the degree distribution of the initial condition did not have this same shape. This indicates that the degree distribution obtained is that of the equilibrium state for the RPAM dynamics and that 7,000 steps is sufficient to reach the equilibrium state from this range of initial conditions.

The degree distribution is, however, only one property of a graph; in fact, graphs can have the same degree distribution and be widely different. Given a finite sequence, one can formulate necessary and sufficient conditions for it to be the degree sequence of a connected graph; sequences that satisfy these conditions are called *realizable*. Given a realizable sequence, there exist various techniques to construct connected graphs that have this sequence as its degree sequence. In some of these techniques, high degree nodes are highly interlinked; in others, the lowest degree nodes are invariably linked to the highest degree nodes (see [Mihail et al. 02] for a nice description). To distinguish the resulting, very different-looking graphs, one needs to use other graph properties. The constructive methods for small graphs with given degree sequences can become unwieldy for graphs of the size of the AS network, but methods for large graphs have also been developed. For instance, the following gives a general method [Molloy and Reed 95] to build a “uniformly random” graph with a given degree distribution.

1. Start with a number of nodes n , and, for each node, choose its degree according to the given degree distribution.
2. Build a new list of nodes by making as many copies of each original node as its chosen degree.
3. Choose randomly a matching of pairs of nodes in this new list, and map the edges back on the original graph.

The result may have to be adapted slightly in order to achieve a simple graph; because the sum of the degrees is always even, it is always possible to do this by rewiring. Constructing a randomly uniform graph of this type, with the same degree distribution as that of the AS network in RV, leads typically to a graph that is very different, however; in particular, it has a much longer mean path. One can also try to construct particularly *dense core* or particularly *sparse core* graphs, by the following two procedures.

For the dense core graph,

1. construct the largest clique the degree sequence allows;
2. start connecting the other nodes in order from the highest to the lowest degree directly on the nodes of the clique: we always choose the node of the clique with the maximum possible available connections;
3. once the clique has no available connections, repeat the previous step, but this time with the nodes directly connected to nodes of the clique etc.

Ties are broken arbitrarily. It is clear that this procedure creates a graph with an extremely dense core, as we “push” the nodes as close to the center of the graph as possible.

For the sparse core graph,

1. sort the nodes by the available number of connections (available degree),
2. connect the lowest available degree node to the highest available degree node (ties are broken arbitrarily),
3. and then return to the first step.

(This sparse graph construction procedure is borrowed from [Mihail et al. 02].) The procedure ends when there are no positive available degrees left. This graph clearly has quite long paths, as the nodes of low degree connect, as a rule, directly to nodes of high degree, thus spreading high degree nodes further apart.

It is not a priori clear that such dense core and sparse core graphs exist for a given degree sequence, even if we ensure that the degree sequence is realizable (by, e.g., taking the true degree sequence of one of the weekly graphs of RV). For the second procedure, [Mihail et al. 02] proves that the sparse core graph can be constructed if the sequence is realizable. However, there is no such proof for the dense core graph procedure, for the reason that this is not true! Fortunately, the process fails to complete only very rarely, and even then a rerun (in which the random breaking of ties falls out differently) usually makes it. In any case, in the experiments that we ran, even failed runs were only about 10 edges short (out of a total of approximately 15,000).

In this section we shall compare the path length distributions and the largest two eigenvalues of the incidence matrix of the observed AS graph with the corresponding quantities of RPAM-simulated graphs, as well as with dense core and sparse core graphs that have the same degree distribution.

7.1. Minimum Path Length (MPL) Comparison

For both the AS network in the last week and RPAM simulations with the same number of nodes, we computed the distribution of the minimum path lengths between two nodes. That is, for each pair of different nodes in the graph, we compute the MPL for that pair, and we then looked at the distribution of these MPL over the $N(N-1)/2$ possible pairs of N nodes. These MPL distributions for the real and a simulated graph are shown in Figure 16: the agreement between the two is striking.

The MPL are typically very small in both the real and the simulated graphs, especially in view of their sparsity—after all, e.g., for the last week of RV, only

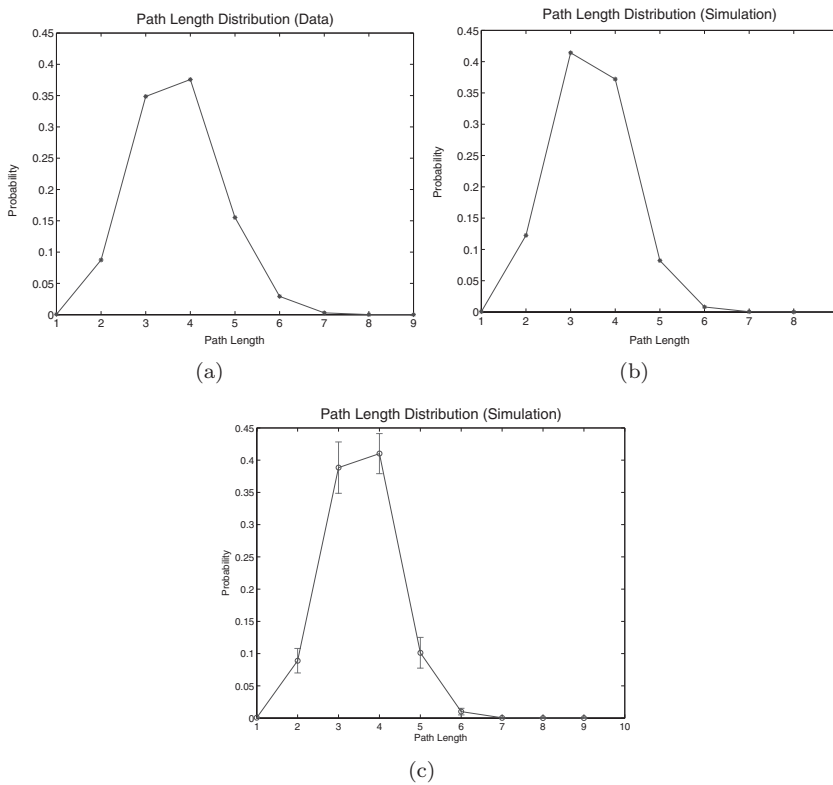


Figure 16. MPL distributions for (a) data and (b) a RPAM simulation: the two distributions are almost identical. (c) shows the averages over five simulations and gives the 99.9% confidence interval for each point in the discrete distribution: despite the randomness of the simulation, path length distribution results for RPAM-simulated graphs (with the same parameters) are almost deterministic.

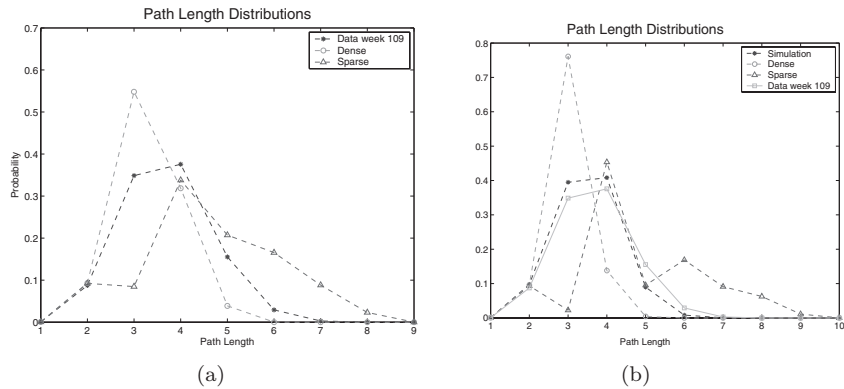


Figure 17. (a) Path length probabilities for three graphs with the same degree distribution: the AS graph during week 109 in RV and the results of dense core and sparse core graph constructions with the same degree distribution; (b) path length probability distributions for an RPAM-simulated graph (based on parameters derived from the AS graph for week 109) and for the sparse core and dense core constructions with the same degree distribution; in this plot, the path length probabilities of the AS graph for week 109 are also shown, for comparison purposes (even though its degree distribution is not absolutely identical to that of the other three, since RPAM simulation is not a replication of a given degree distribution; note that this difference of the degree distribution also means that the dense and sparse core graphs constructed for (b) are different than those for (a)—which is why the path length plots in (a) and (b) are different).

12,988 or 6.54% of the $N(N-1)/2$ possible connections between the $N = 6,304$ nodes are realized. This is due to the hierarchical nature of the graphs, in which the T1-clique at the core of the graph makes distances between points very short (this is the *small world* effect [Albert and Barabási 02]). Two typical path types, accounting for the majority of the paths, are

- node—(T1 or T2)—T2—node',
- node—T2—(T1 or T2)—T2—node',

leading to the peaks in the distribution at path lengths 3 and 4. A bit rarer, but still important, are the cases

- node—(T1 or T2)—node',
- node—T2—(T1 or T2)—(T1 or T2)—T2—node',

which lead to path length values of 2 and 5, respectively. As an example, for the last week of RV, the path types mentioned above account for 34.5%, 35.5%, 8.73%, and 13.5%, respectively, of the full collection of paths in the BGP table.

Figure 17 compares both the real-life AS graph and an RPAM-simulated graph with graphs that have the same degree sequence but are constructed according to the dense core, respectively sparse core, graph construction procedures. (Because we are dealing here with properties that are computed from one graph, the RPAM-simulation here was carried out with parameters that were adopted (where possible) from a particular week of RV, in this case week 109, rather than with the mean quantities listed earlier.) It is quite clear from the plots that the RPAM-simulation is much closer to the AS graph than either the dense core or the sparse core graph, suggesting that the RPAM construction captures more properties of the AS-graph reality than only the degree distribution.

7.2. Eigenvalues of the Incidence Matrix

The eigenvalues of a graph's incidence matrix reveal a lot about the structure of the graph and have been studied systematically for certain categories of graphs. For instance, for a classical random graph, it is known [Bollobás 01] that, as the number n of nodes increases to infinity, the largest eigenvalue of the incidence matrix of the graph is positive and increases proportionally to n , while the second eigenvalue is also positive and increases proportionally to \sqrt{n} .

We computed the largest (in absolute value) two eigenvalues of the incidence matrices corresponding to the data, to the RPAM-simulations, and to the artificially constructed dense core and sparse core graphs.

The results are shown in Tables 3 and 4. In both tables the dense core graph's eigenvalues behave more like those of the AS-graph or of the RPAM-simulated graph than the sparse core graph's eigenvalues. In each case, the sparse core eigenvalues capture only the sign change and the order of magnitude of the two eigenvalues, while the dense core graph captures the sign change, approximates the eigenvalues quite well, and exhibits the gap between the positive and the (absolute value of the) negative eigenvalue, which is prominent in the AS-graph and in the RPAM-simulated graph and absent in the sparse core graph. This gap is systematically wider in the dense core graph than in either the AS-graph or the RPAM-simulated graph, but the discrepancy is less pronounced for the AS-graph than for the RPAM-simulated graph.

In any case, we are far from the behavior of the largest two eigenvalues of classical random graphs, yet another confirmation (although none was needed any longer) of the difference between the Internet graph and a classical random graph.

Before we close, we would like to point out that there are numerous other benchmarks that we could have used for graph comparison, especially some based on the small world phenomenon and the concept of clustering (see [Bollobás 01]);

Week		Eigenvalues		
		Data	Sparse Core	Dense Core
	1	33.1574 −26.3548	25.7407 −25.7375	35.9938 −24.2128
	56	40.1286 −33.3483	32.7347 −32.7346	40.1821 −31.1816
	109	46.4572 −39.1388	38.6270 −38.6170	53.8639 −37.3393

Table 3. For each of the weeks 1, 56, and 109 (last) in RV, the table lists the two eigenvalues with largest absolute values for the data and for the sparse core and dense core graphs constructed with the same degree sequence as the data.

Number of nodes		Eigenvalues		
		RPAM-sim.	Sparse Core	Dense Core
	3000	29.4344 −24.8770	25.1393 −25.1383	35.8797 −22.4648
	4500	35.5025 −31.6772	30.8917 −30.7656	41.9315 −27.5739
	6000	38.3280 −33.5019	33.4802 −33.4801	49.2500 −30.3405

Table 4. The two eigenvalues largest in absolute value for three RPAM-simulations, with 3,000, 4,500, and 6,000 nodes respectively (these are approximately the numbers of nodes in the AS graphs at weeks 1, 56, and 109) and, in each case, for the sparse core and dense core graphs with the same degree sequence as the simulation.

we believe, however, that the ones that we presented are the most straightforward and the simplest possible and that they were sufficient to prove our point.

8. Summary and Conclusion

We have presented a model that we dubbed RPAM (for Revised Preferential Attachment Model), in which nodes belong to a small number of different species that each follow a preferential attachment scheme with respect to their own and the other species. This model was inspired by observations that we made in the data set RV, which consists of weekly summaries of the daily snapshots from 1997 to 2000 available on the web site nlanr-routeviews. We discussed both analytical properties of this model and showed simulated properties; we discussed mostly the degree distribution but also compared the minimum path length distribution and the largest eigenvalues of the incidence matrix. The model turns out to lead

to results remarkably similar to the true Internet graph, at least in these three properties; it is moreover very robust in the sense that adding extra detail to the model, to make it richer and more similar to the less neatly organized real Internet graph, doesn't appreciably affect the degree distribution.

Does this mean that RPAM “explains” the AS connectivity graph in the Internet? Of course not—there are surely many features exhibited by the real-world AS connectivity graph that are not captured by the RPAM model. But it does indicate that the power decay observed in the degree distribution of the AS connectivity graph may be a “fairly coarse” feature of the graph, not as dependent on the finer details of the dynamics as other features.

Linear attachment has been cited as a possible explanation for power law behavior in other graphs. Maybe the idea behind RPAM, i.e., the categorization of the nodes into groups of different “potential” (rate of growth), can work with other graphs as well (like the http link graph, the actors graph, etc.), besides the Internet AS graph. This “potential” could again be subject to the Kismet assumption: for example, in the case of the actors graph, one could assume that an actor's first few appearances reveal his/her talent.

Appendix—A Discussion of Discounting and Its Effect on Our Argument

As mentioned in Section 3.5, the plots of \mathbf{P}_{EA} incorporate an inflation effect. As illustrated by Figure 3, ASs of type T1 and T2 have degrees that grow in time. As a result, the range on the abscissa of the EAP plot will continually grow, and one can expect a gradual move to the right of the EAP curve. To show that this does indeed happen, we computed versions of time-averages of $\mathbf{P}_{\text{EA}}(d, t)$ for truncated time periods, namely for the first and second halves (in time) of the RV data set; we repeated this exercise with a partition of the whole RV time period into three consecutive thirds. The results are shown in Figure 18.

The degree inflation can also be observed in another type of diagram, inspired by Figure 4 in [Chen et al. 02]. Start by ordering the ASs that appear as new somewhere in RV according to the order in time with which they appear in RV, and number them one by one, thus obtaining their *ordinal numbers*. Figure 19 plots the degrees (at the time of connection) of the old ASs to which a new AS connects versus the ordinal number of that new AS. One clearly observes lines of positive slope in this graph; the distinctive lines at the top correspond to the T1s, which once more stand out. On the other hand, lines of constant degree are more sparsely populated at the “large end” of the ordinal number range than at the “low end”; this demonstrates that new ASs tend to connect to existing ASs of larger and larger degree (see Figure 19(d)).

Analogous readjustments are carried out in finance, for the same reason, as in, e.g., the price evolution in the Black-Scholes model, where the transition to a risk-free measure causes inflation to disappear. One can remove the inflation effect in our \mathbf{P}_{EA} graphs by a “discounting” very similar to what is done in finance; we define the *discounted degree* $E_i^{\text{disct.}}(t)$ of node i at time t as $E_i^{\text{disct.}}(t) = E_i(t) \times \sum_j E_j(t_0) / \sum_j E_j(t)$, where t_0 is some arbitrary but fixed time (which could be, e.g., the start time of RV or the end time). The corresponding discounted preferential attachment, directly in binned form, is then

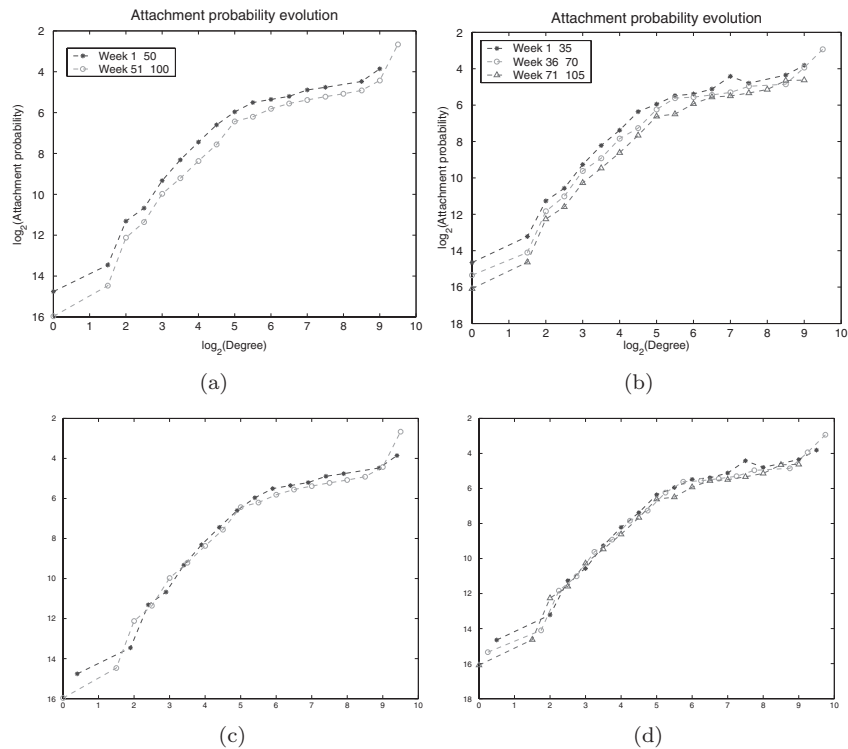


Figure 18. EAP evolution throughout the data set: (a) Evolution from the first to the second half of the time period covered by RV; the two curves almost coincide if the first is moved by 0.4 to the right, as shown in (c) on the lower left. (b) Evolution in two steps, from the first to the second 1/3, and then from the second to the third 1/3, of the time covered in RV; the three curves almost coincide if the first is moved by 0.5 and the second by 0.25 to the right, as shown in (d). (Note: the “slide” parameters are approximate and for illustration purposes only.)

$$\tilde{\mathbf{P}}_{\text{EA,bin}}^{\text{disct.}}(k,t) = \frac{\sum_i; i \text{ is old in week } t \text{ and } 2^{k/2-1/4} \leq E_i^{\text{disct.}}(t) < 2^{k/2+1/4} a(i,t)}{\left(\sum_j; j \text{ is old in week } t a(j,t)\right) \# S_k^{\text{old}}(t)}$$

where $S_k^{\text{old}}(t)$ is the set defined as follows:

$$S_k^{\text{old}}(t) = \{i; i \text{ is old in week } t \text{ and } 2^{k/2-1/4} \leq E_i^{\text{disct.}}(t) < 2^{k/2+1/4}\}.$$

The corresponding curves, for different weeks, are now no longer shifted with respect to each other. One can again compute the average (discounted) probability, over all the weeks, for each bin k ; this average can be taken uniformly (i.e.,

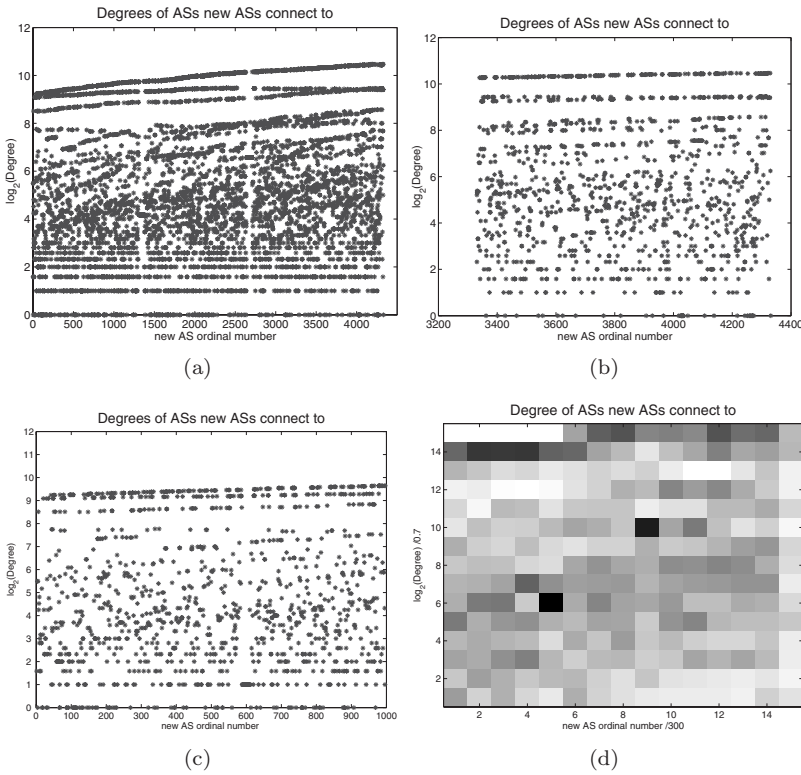


Figure 19. Degrees of the ASs to which new ASs connect, versus the ordinal number of those new ASs (see text): the main figure is (a), while (b) and (c) are just the end and the beginning of (a) (last and first 1,000 ASs), respectively. (d) is the same as (a) but in coarser scale (the darker a rectangle is, the more points it contains), and it makes it easier to visualize the degree growth: observe that the darker rectangles seem to be moving up and right.

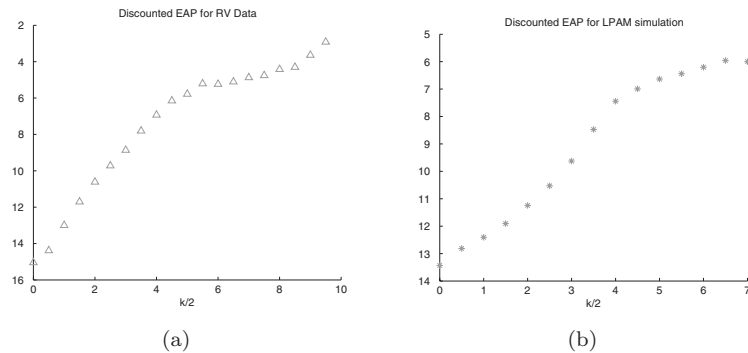


Figure 20. Graphs of the experimental preferential attachment probability versus the degree of the node to which attachment is made: (a) for the RV data and (b) for an LPAM simulation. In both cases the probabilities are binned by degree in each week, discounted so that the different weeks can be compared reliably, and averaged over weeks (see text).

equal weight, for each bin, to all weeks t) or weighted (i.e., each $\tilde{\mathbf{P}}_{\text{EA},\text{bin}}^{\text{disct.}}(k, t)$ is given weight proportional to $\#\{i; i \text{ is old in week } t \text{ and } 2^{k/2-1/4} \leq E_i^{\text{disct.}}(t) < 2^{k/2+1/4}\}$ in the computation of the average over t). It turns out that both averages give pretty much the same result. In Figure 20(a) we show the unweighted averaged result. A comparison with Figure 5 shows the qualitative behavior of the curve hasn't changed much as a result of discounting.

A similar discounting can, of course, be done with our (or any) simulation of LPAM; the result is shown in Figure 20(b), which is the analog of Figure 20(a), with the same discounting, binning, and averaging procedures, but now for the LPAM simulation instead of the RV data.

The discounted LPAM graph is surprisingly nonlinear, despite the linearity of the underlying probability rule with respect to discounted degree (the discounted degree of node i is after all, up to a multiplicative constant, equal to the ratio $\frac{d_i(t)}{\sum_j d_j(t)}$, which is the probability that a new connection will choose to attach to node i in LPAM). At this point we haven't completely understood this effect; its explanation will have to wait for a future paper. The metric used in [Chen et al. 02], which plots the cumulative Experimental Attachment Probability derived from an LPAM simulation versus the cumulative theoretical LPA probability (which is in fact the “cumulative discounted degree”), does look much more linear and might therefore be deemed more appropriate for a comparison with data. On the other hand, the closely parallel behavior of the RV attachment probability and the LPAM simulation probability in the “middle region” would have been masked by the integration over the less similar lower (or higher, depending in which direction one sorts to define the cumulant) degree regions.

Acknowledgments. The authors gratefully acknowledge support by NSF grant DMS-9872890, as well as by AFOSR and DARPA.

K. Drakakis would also like to thank the Lilian Boudouris foundation for the scholarship that it granted him. T. Khovanova thanks A. Broido for useful discussions and Alexey Radul for valuable assistance with Java programming issues. All three authors thank the reviewers for insightful comments that led to an improved revision of the paper.

References

- [Albert and Barabási 02] R. Albert and A.-L. Barabási. “Statistical Mechanics of Complex Networks.” *Reviews of Modern Physics* 74 (2002), 47.
- [Barabási and Albert 99] A.-L. Barabási and R. Albert. “Emergence of Scaling in Random Networks.” *Science* 286 (1999), 509–512.
- [Bollobás 01] B. Bollobás. *Random Graphs*, Second edition. Cambridge, UK: Cambridge University Press, 2001.
- [Broido et al. 02] A. Broido, E. Nemeth, and K. C. Claffy. “Internet Expansion, Refinement and Churn.” *European Transactions on Telecommunications* 13:1 (2002), 33–51.
- [Callaway et al. 01] D. Callaway, J. Hopcroft, J. Kleinberg, M. Newman, and S. Strogatz. “Are Randomly Grown Graphs Really Random?” *Physical Review E* 64 (2001), 041902.
- [Chang et al. 04] H. Chang, R. Govindan, S. Jamin, S. Shenker, and W. Willinger. “Towards Capturing Representative AS-Level Internet Topologies.” *Computer Networks* 44 (2004), 737–755.
- [Chen et al. 02] Q. Chen, H. Chang, R. Govindan, S. Jamin, S. Shenker, and W. Willinger. “The Origin of Power Laws in Internet Topologies Revisited.” In *Proceedings of IEEE INFOCOM 2002, the 21st Annual Joint Conference of the IEEE Computer and Communications Societies*, CD-ROM. IEEE Press, 2002.
- [Floyd and Kohler 03] S. Floyd and E. Kohler. “Internet Research Needs Better Models.” *ACM SIGGRAPH Computer Communication Review* 33:1 (2003), 29–34.
- [Medina et al. 00] A. Medina, I. Matta, and J. Byers. “On the Origin of Power Laws in Internet Topologies.” *ACM Computer Communication Review* 30:2 (2000), 18–28.
- [Mihail et al. 02] M. Mihail, Ch. Gkantsidis, A. Saberi, and E. Zegura. “On the Semantics of Internet Topologies.” Technical report GIT-CC-02-07, College of Computing, Georgia Institute of Technology, 2002.
- [Molloy and Reed 95] M. Molloy and B. Reed. “A Critical Point for Random Graphs with a Given Degree Sequence.” *Random Structures and Algorithms* 6 (1995), 161.
- [NRC 01] National Research Council, editor. *The Internet’s Coming of Age*. Washington, DC: National Academies Press, 2001.

- [Siganos et al. 03] G. Siganos, M. Faloutsos, P. Faloutsos, and C. Faloutsos. “Power-Laws and the AS-level Internet Topology.” *IEEE/ACM Transactions on Networking* 11:4 (2003), 514–524.
- [Tangmunarunkit et al. 01] H. Tangmunarunkit, J. Doyle, R. Govindan, S. Jamin, S. Shenker, and W. Willinger. “Does AS Size Determine Degree in AS Topology?” *ACM Computer Communication Review* 31:5 (2001), 7–8.
- [Willinger et al. 02] W. Willinger, R. Govindan, S. Jamin, V. Paxson, and S. Shenker. “Scaling Phenomena in the Internet: Critically Examining Criticality.” *PNAS* 99 Suppl. 1 (2002), 2573–2580.
- [Zegura et al. 97] E. W. Zegura, K. L. Calvert, M. J. Donahoo. “A Quantitative Comparison of Graph-based Models for Internet Topology.” *IEEE/ACM Transactions on Networking* 5:6 (1997), 770–783.

Ingrid Daubechies, Program in Applied and Computational Mathematics, Fine Hall,
Princeton University, Washington Road, Princeton, NJ 08544
(ingrid@math.princeton.edu)

Konstantinos Drakakis, JCMB, KB, Mayfield Road, Edinburgh EH9 3JZ, UK
(K.Drakakis@ed.ac.uk)

Tanya Khovanova, 125 Elm Street, Belmont, MA 02478 (tanyakh@yahoo.com)

Received September 30, 2004; accepted June 13, 2005.