

Scale-free Geometry in Object-Oriented Programs

Alex Potanin, James Noble, Marcus Frean, Robert Biddle
School of Mathematical and Computing Sciences
Victoria University of Wellington, New Zealand
{alex, kjax, marcus, robert}@mcs.vuw.ac.nz

Introduction

Object-oriented programs, when executed, produce a complex web of objects that can be thought of as a graph with objects as nodes and references as edges. In recent years interest has grown in the geometry of networks (or graphs), particularly those of human origin, many of which show a rather striking property: their structure is *scale-free*. In the case of the World Wide Web, for example, the number of web pages with 1 incoming link is about twice the number with 2 incoming links, and *that* is twice the number with 4 links, and so on all the way up to Google and other massively referenced sites [2]. The phrase ‘scale-free’ relates to the fact that if we double the number of links n , the number of pages is always halved (or other fixed ratio) regardless of what n is. Compare this to what happens if a graph is constructed by simply adding links at random [4]. Doing this leads to nearly all nodes having around the same number of links (i.e. the number of links divided by the number of nodes), and hence such random graphs have a ‘typical scale’ about them. By contrast, the web has *no* typical scale to its connectivity — a remarkable and somewhat counterintuitive property closely related to fractals.

Other scale-free graphs include the network formed by co-authors of papers in medical journals, the physical connections forming the Internet, the network of airports connected by airline flights, networks of sexual contacts, and even the patterns of connectivity between neurons in the human brain [1]. Well before being noticed in real-world graph structures, scale-free distributions were found in other contexts, such as the relative frequencies of English words, the distribution of personal wealth, the sizes of cities, and the number of earthquakes of given strength [12].

In this article we examine the graphs formed by object-oriented programs written in a variety of languages, and show that these turn out to be scale-free networks as well. Apart from its considerable intrinsic interest, this unexpected facet of the geometry of real programs may help us optimize language runtime systems, improve the design of future object-oriented languages, and reexamine modern approaches to software design.

Power Laws

The way to detect a scale-free phenomenon is to see if it shows up statistically in the form of a *power law*. In power law distributions the number of occurrences N_k of some event of size k is proportional to k raised to some power. One drawback is that very rare events are by their nature “noisy” (there may be one node with, say 1000 connections, and another with 1005, but none with 1002). For this reason an alternative approach is often adopted in which we first rank the event sizes by how often they occur, and then look for a power law in the relationship between the number of occurrences N_k , and rank R_k of the form:

$$N_k \propto R_k^s$$

The easiest way to see this is to take logarithms of both sides, or to plot N versus R on logarithmic scales — if the distribution follows a power law we expect to see a straight line with slope s .

For example, consider how often a particular word appears in any English novel. Common words like “*the*”, “*of*”, or “*and*” can be found many orders more times than the majority of other words, while at the other extreme there are a huge number of words that are used only rarely. In 1925, George Kingsley Zipf, a Harvard linguistics professor, conducted empirical studies [11] of word occurrences and made a remarkable observation that if we rank the words by the number of times they can be found in a text of a particular novel, then their rank will be proportional to their number of occurrences. Hence, if you take your favourite novel and draw a logarithmic plot of the number of times you can find each word against the rank of such a word, then you will see a straight line

Object Graphs

An object graph — the object instances created by a program and the links between them — is the skeleton of the execution of an object-oriented program. Because each node in the graph represents an object, the graph grows and changes as the program runs. It contains just a few objects when the program is started, gains more objects as they are created, and loses objects when they are no longer required. The structure of the graph (the links between objects) changes too, as every assignment statement to an object’s field may create, modify, or remove an edge in the graph.

Figure 1 illustrates the object graph of a simple part of a program — in this case a doubly-linked list of `Student` objects. The list itself is represented by a `LinkedList` object which has two references to `Link` objects representing the head and tail of the list. Each `Link` object has two references to other `Link` objects — the previous and the next links in the list, plus a third reference to one of the `Student` objects contained in the list.

Object graphs are the most fundamental structure in object-orientation. The primary aim of object-oriented analysis is to model the real world in terms of

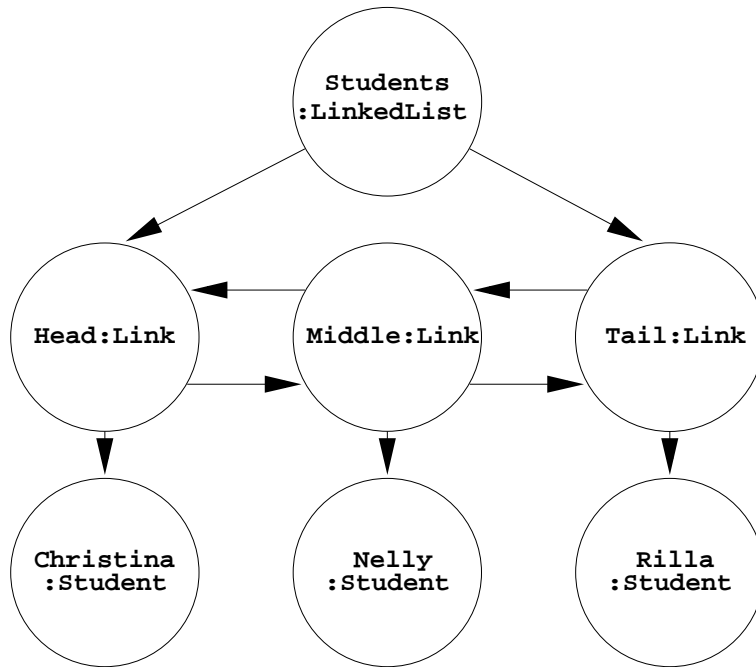


Figure 1: A simple object graph of a linked list. Each link object has two references to other link objects, except for the head and tail of the list. The student objects stored inside the list are pointed at by the link objects that store them.

communicating objects (that is, in terms of an object graph), while object-oriented design produces a description of an object graph that will eventually be embodied in a program. The artifacts and methods of object-orientation (classes, associations, interfaces, inheritance, packages, patterns, UML, CRC Cards, and so on) are ultimately techniques for defining object graphs by describing the contents of the objects and the structure of the links between them.

Given that object graphs are so basic to object-oriented programs, it is somewhat surprising that so little attention has been paid to their global structure. Some temporal properties of object graphs have been analysed to support garbage collection — such as the time performance of garbage collection algorithms and the distributions of object lifetimes [7]. Visualisation of object graphs is used to support debugging [10]. Programming language designers work on controlling object graph structures using type systems [8], and compiler developers analyse parts of the graph to find ways to improve program performance [6].

Concerning scale-free structure in programs, it has been shown [9] that class diagrams of the Java Development Kit 1.2 have a scale-free aspect when *classes*

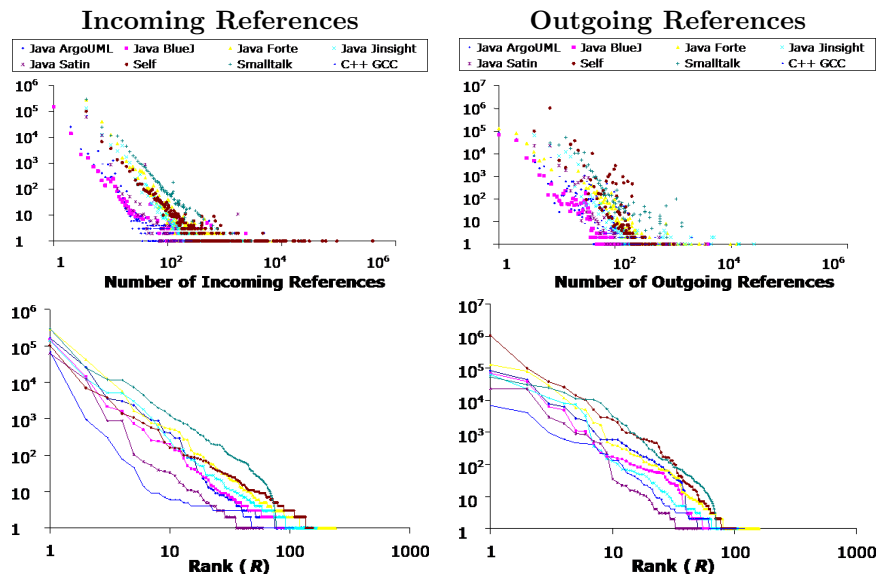


Figure 2: Power laws in object graphs. The upper two figures plot the number of objects with k references versus the number of references k , for incoming and outgoing references respectively. The lower two figures plot the number of occurrences of each number of incoming (and outgoing) references against their rank. All exhibit clear linearity on log-log scales, the characteristic feature of scale-free networks.

are considered as nodes (not the actual instantiation of objects *of* those classes at run-time). Similar structure has also been observed in the distribution of pointers to atoms in Lisp [3].

Power Laws in Object Graphs

To analyse the geometry of object graphs in Java programs we used the facilities of the Heap Profiler (HPROF) Library and the Java Virtual Machine Profiler Interface (JVMPi) to collect a corpus of 60 object graphs from 35 programs, encoded as binary snapshots of the Java heap. These are instantaneous static snapshots of the objects in the programs, together with the topology of the references between them — exactly the kind of information shown in figure 1. To analyse this corpus, we extended the Java Heap Analysis Tool (HAT) [5] that parses these snapshots to determine the properties of the program’s object graphs.

For each graph, we first count the number of objects with k references, for k from 1 upwards. If the object graph is scale-free we should see a straight line when the number of objects is plotted against k on log scales, or when it is

plotted against their rank ordering. Without exception, all the object graphs in our corpus demonstrated this phenomenon: the object graphs are *scale free*. The same general effect applies to both incoming references (which reflect how *popular* objects are), outgoing references (which reflect object *sizes*), and can be detected in multiple snapshots taken during a single run of some programs. It appears that the world of object graphs is indeed scale-free, just like the World Wide Web, the Internet, and many other networks around us.

Figure 2 shows five large Java snapshots and three additional object graphs from programs in other object-oriented languages, as described in Table 1. These were chosen for their size, popularity, and diversity. These plots show all objects inside each program’s memory at the moment the snapshot was taken. Although not shown here, excluding pointers from local variables and other references from the stack further improves the scale-free structure.

Perhaps the most intriguing aspect of the ranked graphs shown in figure 2 is that all eight plots have similar slopes. This is surprising because they come from run-time snapshots of separate programs written for entirely different ends in different languages! For incoming references the slope of the line is close to -2.5 while for outgoing references it is around -3 .

Figure 3 shows the number of objects having a given combination of incoming and outgoing references for the Forte data — the largest Java heap snapshot in our study. Notably there are no objects with both high in-degree and high out-degree: on the contrary, those with many incoming references have very few outgoing ones, and vice-versa. This may be a consequence of widely shared data structures with many outgoing references, such as arrays, having a proxy object that hides the actual reference to the array from the other objects who use it.

Discussion

If object-oriented programs were constructed out of completely independent components, like Lego bricks, then we would expect the distribution of the size and popularity of objects to stay the same, no matter how large the program: just as fixed-size Lego bricks can form buildings of any size. The rhetoric of object-oriented design is that large programs can be constructed in just the same way as small programs, by encapsulating complexity within objects at one level of abstraction, and then composing those objects together at the next. Thus all objects should appear to be the same size and complexity: larger programs merely use more objects and more levels of abstraction. This is not what we have found in our corpus of object-oriented snapshots. The power law indicates the reverse of the Lego hypothesis: there is no evidence of a typical size (a Lego brick) to objects at all.

The relative steepness of the slope we obtained reflects the fact that there are an exceedingly large number of objects with very few references. We might take this to imply that programmers prefer simple objects to complex ones, avoiding complexity just as software engineering guidelines would suggest. The power law distribution indicates that this adage is not followed — instead large

Program	Description	Objects	Objects
		≥ 1 in-refs	≥ 1 out-refs
Java ArgoUML	A popular CASE tool.	203,875	153,106
Java BlueJ	Visual OO programming and learning environment.	171,666	123,701
Java Forte	A Java integrated development environment by Sun Microsystems.	358,279	267,755
Java Jinsight	A memory analysis tool by IBM.	76,312	118,272
Java Satin	A pen-based user interface research tool from Stanford University.	80,415	53,328
C++ GCC	GCC is a C++ compiler used by developers for Unix platforms.	71,990	15,064
Self	Self is a prototype-based OO language and environment.	120,748	1,259,668
Smalltalk	Smalltalk is one of the original OO languages, self-contained in an environment developed using Smalltalk.	375,529	188,031

Table 1: The object graphs shown in figure 2. We list the number of objects with at least one incoming or outgoing reference. We obtained them when the programs were most heavily used. While for C++ and Java we had to run particular programs, for Self and Smalltalk it was possible to obtain all objects for all the programs running inside these systems.

programs contain objects that are much more highly connected than one might expect. For example, for the unranked power law the Java programs all have a slope of approximately -2 , and it follows that for a given number of objects of size k there are about one quarter that number of size $2k$. Thus a program generating 10000 objects of size 1 will also involve about 2500 objects of size 2, 625 of size 4, 156 of size 8 and so on, leading to an expectation of one object of size roughly 100. In programs with twice as many objects altogether, we expect the number of very popular objects and the size of the largest object to go up by a factor of $\sqrt{2}$.

One aspect of scale-free networks is their robustness to damage. Because the vast majority of nodes are poorly connected to the rest of the graph, deleting them has a negligible effect on the connectivity of those remaining [2]. On the other hand, a small number of “hub” objects are very highly connected, and deleting *them* is far more destructive. An implication of this is that by concentrating our debugging methodologies on such popular well-connected objects, rather than the unpopular ones, we may be able to improve the reliability of code more efficiently: first eliminate bugs from the hubs, then deal with other objects.

Aside from their scale-free character, power laws are notable in that they have much longer “tails” than, say, exponential distributions. Thus larger programs will contain considerably larger objects and more popular objects than

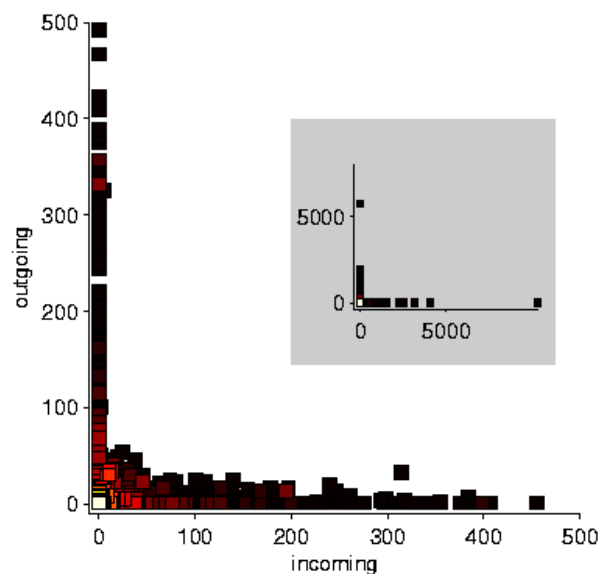


Figure 3: Distribution of incoming vs outgoing references in the Forte snapshot. Lighter squares correspond to a greater number of objects having that combination of references. The objects shown have up to 6,000 outgoing and 10,000 incoming references.

simpler models would predict. This may have consequences for both the design and implementation of object-oriented programming languages. For example, this explains why garbage collectors can improve their performance in search and traversal of object graphs by assuming that objects will probably have only one or two outgoing references.

To conclude, we have found that distributions of incoming and outgoing references in object graphs follow a power law. This unexpected result raises theoretical questions and has implications for debugging costs, program understanding, and garbage collection. More generally, it challenges the perceived wisdom of OO design: unlike Lego bricks, objects within large programs have no characteristic scale.

References

- [1] BARABASI, A.-L. *Linked: The New Science of Networks*. New York: Perseus Press, 2002.
- [2] BARABASI, A.-L., ALBERT, R., JEONG, H., AND BIANCONI, G. Power-law distribution of the world wide web. *Science* 287, 2115a (2000). In technical comments.

- [3] CLARK, D. W., AND GREEN, C. C. An empirical study of list structures in Lisp. *Communications of the ACM* 20, 2 (February 1977), 78–87.
- [4] ERDOS, P., AND RENYI, A. On the strength of connectedness of random graphs. *Acta Math. Acad. Sci. Hungary* 12, 35 (1961), 261–267.
- [5] FOOTE, B. Heap analysis tool. <http://java.sun.com/people/billf/heap/>, 2002.
- [6] GHIYA, R., AND HENDREN, L. J. Is it a tree, a DAG, or a cyclic graph? a shape analysis for heap-directed pointers in C. In *Symposium on Principles of Programming Languages* (1996), pp. 1–15.
- [7] JONES, R., AND LINS, R. *Garbage Collection*. Wiley, 1996.
- [8] NOBLE, J., VITEK, J., AND POTTER, J. Flexible alias protection. In *European Conference on Object-Oriented Programming (ECOOP)* (1998), Lecture Notes in Computer Science, Springer-Verlag.
- [9] VALVERDE, S., FERRER-CANCHO, R., AND SOLE, R. Scale-Free Networks from Optimal Design. *Europhysics Letters* 60 (2002), 512–517.
- [10] ZIMMERMANN, T., AND ZELLER, A. Visualizing memory graphs. In *Software Visualization*, vol. LNCS 2269 of *Lecture Notes in Computer Science*. Springer-Verlag, May 2001, pp. 191–204.
- [11] ZIPF, G. K. *Psycho-Biology of Languages*. Houghton-Mifflin, 1935.
- [12] ZIPF, G. K. *Human behavior and the principle of least effort : an introduction to human ecology*. New York : Hafner, 1965. Facsimile of 1949 edition.