

A SHORT HISTORY OF LINES OF CODE (LOC) METRICS

Version 1.0

May 9, 2008

Abstract

The oldest metric for software projects is that of “lines of code” (LOC). This metric was first introduced circa 1960 and was used for economic, productivity, and quality studies. The economics of software applications were measured using “dollars per LOC.” Productivity was measured in terms of “lines of code per time unit.” Quality was measured in terms of “defects per KLOC” where “K” was the symbol for 1000 lines of code. The LOC metric was reasonably effective for all three purposes.

As additional higher-level programming languages were created the LOC metric began to encounter problems. LOC metrics were not able to measure non-coding activities such as requirements and design which were becoming increasingly expensive. LOC metrics not able to measure defects in requirements and design.

These problems became so severe that a controlled study in 1994 that used both LOC metrics and function point metrics for 10 versions of the same application coded in 10 languages reached an alarming conclusion: LOC metrics violated the standard assumptions of economic productivity so severely that using LOC metrics for studies involving more than one programming language comprised professional malpractice.

Capers Jones

Email: CJonesiii@cs.com

Copyright © 2008 by Capers Jones & Associates LLC. All rights reserved.

A Short History of Lines of Code (LOC) Metrics

INTRODUCTION

It is interesting to consider the history of lines of code (LOC) metrics and some of the problems with LOC metrics that led IBM to develop function point metrics. Following is a brief history of LOC metrics from 1960 through today, with projections to 2010:

Lines of Code Metrics Circa 1960

When the LOC metric was first introduced there was only one programming language and that was basic assembly language. Programs were small and coding effort comprised about 90% of the total work. Physical lines and logical statements were the same thing for basic assembly language.

In this early environment, LOC metrics were useful for both economic, productivity, and quality analyses. Unfortunately as the software industry changed, the LOC metric did not change and so became less and less useful until by about 1980 it had become extremely harmful without very many people realizing it.

Lines of Code Metrics Circa 1970

By 1970 basic assembly had been supplanted by macro-assembly. The first generation of higher-level programming languages such as COBOL, FORTRAN, and PL/I were starting to be used. Usage of Basic assembly language was beginning to drop out of use as better alternatives became available. This was perhaps the first instance of a long series of programming languages that died out, leaving a train of aging legacy applications that would be difficult to maintain as programmers and compilers stopped being available.

The first known problem with LOC metrics was in 1970 when many IBM publication groups exceeded their budgets for that year. It was discovered (by the author) that technical publication group budgets had been based on 10% of the budget for programming. The publication projects based on assembly language did not overrun their budgets, but manuals for the projects coded in PL/S (a derivative of PL/I) had major overruns. This was because PL/S reduced coding effort by half, but the technical manuals were as big as ever.

The initial solution to this problem was to give a formal mathematical definition to language “levels.” The “level” was defined as the number of statements in basic assembly language needed to equal the functionality of 1 statement in a higher-level language. Thus COBOL was a “level 3” language because it took 3 basic assembly statements to equal 1 COBOL statement. Using the same rule, SMALLTALK is a level 18 language.

For several years before function points were invented, IBM used “equivalent assembly statements” as the basis for estimating non-code work such as user manuals. Thus instead of basing a publication budget on 10% of the effort for writing a program in PL/S, the budget would be based on 10% of the effort if the code were basic assembly language.

The documentation problem was one of the reasons IBM assigned Allan Albrecht and his colleagues to develop function point metrics. Additional programming languages such as APL were starting appear, and IBM wanted both a metric and estimating methods that could deal with non-coding work as well as coding in an accurate fashion.

The use of macro-assembly language had introduced reuse, and this caused measurement problems too. It raised the issue of how to count reused code in software applications or any other reused material. The solution here was to separate productivity into two topics: 1) development productivity; 2) delivery productivity. The former dealt with the code and materials that had to be constructed from scratch. The latter dealt with the final application as delivered, including reused material. For example using macro-assembly language a productivity rate for *development productivity* might be 300 lines of code per month. But to due to reusing code in the form of macro expansions, *delivery productivity* might be as high as 750 lines of code per month.

This is an important business distinction that is not well understood even in 2008. The true goal of software engineering is to improve the rate of delivery productivity and not development productivity. To be successful reused code needs to approach or achieve zero-defect status. It does not matter what the development speed is, if once completed the code can then be used in hundreds of applications. As “service-oriented architecture” (SOA) and “software as a service” (SaaS) approach, their goal is to make dramatic improvements in the ability to deliver software features. Development speed is comparatively unimportant so long as quality approaches zero-defect levels.

Another issue shared between macro-assembly language and other new languages was the difference between physical lines of code and logical statements. Some languages, such as Basic, allowed multiple statements to be placed on a physical line. Other languages, such as COBOL, divided some logical statements into multiple physical lines. The difference between a count of physical lines and a count of logical statements could differ by as much as 500%. For some languages there would be more physical lines than logical statements, but for other languages the reverse was true. This problem was never fully resolved by LOC users and remains troublesome even in 2008.

Due to the increasing power and sophistication of high-level programming languages, the percentage of project effort devoted to coding was dropping from 90% down to about 50%. As coding effort declined, LOC metrics were no longer effective for economic, productivity, or quality studies.

After function point metrics were developed circa 1975 the definition of “language level” was expanded to include the number of logical code statements equivalent to 1 function point. COBOL, for example requires about 105 statements per function point in the procedure and data divisions. This expansion is the mathematical basis for “backfiring” or direct conversion from source code to function points. Of course individual programming styles make backfiring a method with poor accuracy.

There are tables available from several consulting companies such as David Consulting, Gartner Group, and Software Productivity Research (SPR) that provide values for source code statements per function point for hundreds of programming languages.

Lines of Code Metrics Circa 1980

By about 1980 the number of programming languages had topped 50 and object-oriented languages were rapidly evolving. As a result, software reusability was increasing rapidly.

Another issue that surfaced circa 1980 was the fact that many applications were starting to use more than one programming language, such as COBOL and SQL. The trend for using multiple languages in the same application has become the norm rather than the exception. However the difficulty of counting lines of code with accuracy was increased when multiple languages were used.

In the middle of this decade the first commercial software cost estimating tool based on function points had reached the market, SPQR/20 which was released in 1985. This tool supported estimates for 30 common programming languages and also could be used for combinations of more than one programming language. This tool also included sizing and estimating of paper documents such as requirements, design, and user manuals. It also estimated non-coding tasks including testing and project management.

Because LOC metrics were still used, the SPQR/20 tool expressed productivity and quality results using both function points and LOC metrics. Because it was easy to switch from one language to another, it was interesting to compare the results using both metrics when changing from Macro-assembly to Fortran or Ada or PL/I.

Within a few years all other commercial software estimating tools would also support function point metrics, so that CHECKPOINT, COCOMO, PRICE-S, SEER, SLIM SPQR/20 and others could express estimates in terms of both function points and LOC metrics.

By the end of this decade coding effort was below 35% of total project effort, and LOC was no longer valid for either economic or quality studies. LOC metrics could not quantify requirements and design defects, which now outnumbered coding defects. LOC metrics could not be used to measure any of the non-coding activities such as requirements, design, documentation, or project management.

The response of the LOC users to these problems was unfortunate: they merely stopped measuring anything but code production and coding defects. The bulk of all published reports based on LOC metrics cover less than 35% of development effort and less than 25% of defects, with almost no data being published on requirements and design defects, rates of requirements creep, design costs, and other modern problems.

Lines of Code Metrics Circa 1990

By about 1990 not only were there more than 500 programming languages in use, but some applications were written in 12 to 15 different languages. There were no international standards for counting code, and many variations were used sometimes without being defined. In 1991 the first edition of the author's book Applied Software Measurement included a proposed draft standard for counting lines of code based on counting logical statements. One year later Bob Park from the Software Engineering Institute (SEI) also published a proposed draft standard, only based on counting physical lines. A survey of software journals in 1993 found that about one third of published articles used physical lines, one third used logical statements, and the remaining third used LOC metrics without even bothering to say how they were counted. Since there is about a 500% variance between physical LOC and logical statements for many languages, this was not a good situation.

The technical journals that deal with medical practice and engineering often devote as much as 50% of the text to explaining and defining the counting methods used to derive the results. The software engineering journals, on the other hand, often fail to define the counting methods at all. The software journals seldom devote more than a few lines of text to explaining the nature of the

measurements used for the results. This is one of several reasons why the term “software engineering” is something of an oxymoron.

But there was a worse problem than ambiguity in counting lines of code. The arrival of Visual Basic introduced a class of programming languages where counting lines of code was not even possible. This is because a lot of Visual Basic “programming” was not done with procedural code but rather with buttons and pull-down menus. Of the approximate 700 programming languages and dialects in existence, there are only counting rules for about 50. About another 500 are similar to other languages and could share the same counting rules. But for at least 50 languages that use graphics or visual means to augment procedural code, there are no code counting rules at all. Unfortunately some of the languages without code counting rules tend to be most recent languages that are used for web site development.

In the middle of this decade a controlled study was done that used both LOC metrics and function points for 10 versions of the same application written in 10 different programming languages including four object-oriented languages. This study was published in American Programmer in 1994. This study found that LOC metrics violated the basic concepts of economic productivity and penalized high-level and OO languages due to the fixed costs of requirements, design, and other non-coding activities. This was the first published study to state that LOC metrics constituted professional malpractice if used for economic studies where more than one programming language was involved.

By the 1990’s a most consulting studies that collected benchmark and baseline data used function points. There are no large-scale benchmarks based on LOC metrics. The International Software Benchmark Standards Group (ISBSG) was formed in 1997 and only publishes data in function point form. Consulting companies such as SPR and the David Consulting Group also use function point metrics.

By the end of the decade, some projects were spending less than 20% of the total effort on coding, so LOC metrics could not be used for the 80% of effort outside the coding domain. The LOC users remained blindly indifferent to these problems, and continued to measure only coding, while ignoring the overall economics of complete development cycles that include requirements, analysis, design, user documentation, project management, and many other non-coding tasks. By the end of the decade non-coding defects in requirements and design outnumbered coding defects almost 2 to 1. But since non-code defects could not be measured with LOC metrics the LOC literature simply ignores them.

Lines of Code Metrics Circa 2000

By the end of the century the number of programming languages had topped 700 and continues to grow at more than 1 new programming language per month. Web applications are mushrooming, and all of these are based on very high-level programming languages and substantial reuse. The Agile methods are also mushrooming, and also tend to use high-level programming languages. Software reuse in some applications now tops 80%. LOC metrics cannot be used for most web applications and are certainly not useful for measuring Scrum sessions and other non-coding activities that are part of Agile projects.

Function point metrics have become the dominant metric for serious economic and quality studies. But two new problems have appeared that have kept function point metrics from actually becoming the industry standard for both economic and quality studies.

The first problem is the fact that some software applications are now so large (>300,000 function points) that normal function point analysis is too slow and too expensive to be used. There are gaps at both ends of normal function point analysis. Above 15,000 function points the costs and schedule for counting function point metrics become so high that large projects are almost never counted. (Function point analysis operates between 400 and 600 function points per day per counter. The approximate cost is about \$6.00 per function point counted.)

At the low end of the scale, the counting rules for function points do not operate below a size of about 15 function points. Thus small changes and bug repairs cannot be counted. Individually such changes may be as small as 1/50th of a function point and are rarely larger than 10 function points. But large companies can make 30,000 or more changes per year, with a total size that can top 100,000 function points.

The second problem is that the success of function points has triggered an explosion of function point “clones.” As of 2008 there are at least 24 function point variations. This makes benchmark and baseline studies difficult, because there are very few conversion rules from one variation to another. In addition to standard IFPUG function points there are also Mark II function points, COSMIC function points, Finnish function points, Netherlands function points Australian function points, web-object points, and many others.

Although LOC metrics continue to be used, they continue to have such major errors that they constitute professional malpractice for economic and quality studies where more than one language is involved, or where non-coding issues are significant.

Lines of Code Metrics Circa 2010

It would be nice to predict an optimistic future, but if current trends continue within a few more years the software industry will have more than 800 programming languages of which about 750 will be obsolete or becoming dead languages, more than 20 variations for counting lines of code, more than 50 variations for counting function points, and probably another 20 unreliable metrics such as “cost per defect” or percentages of unknown numbers. (The software industry loves to make claims such as “improve productivity by 10 to 1” without defining either the starting or the ending point.)

Future generations of sociologists will no doubt be interested in why the software industry spends so much energy on creating variations of things, and so little energy on fundamental issues. No doubt large projects will still be cancelled, litigation for failures will still be common, software quality will still be bad, software productivity will remain low, security flaws will be alarming, and the software literature will continue to offer unsupported claims without actually presenting quantified data.

What the software industry needs is actually fairly straightforward: 1) measures of defect potentials from all sources expressed in terms of function points; 2) measures of defect removal efficiency levels for all forms of inspection and testing; 3) activity-based productivity benchmarks from requirements through delivery and then for maintenance and customer support from delivery to retirement using function points; 4) certified sources of reusable material near the zero-defect level; 5) much improved security methods to guard against viruses, spyware, and hacking; 6) licenses and

board-certification for software engineering specialties. But until measurement becomes both accurate and cost-effective, none of these are likely to occur. An occupation that will not measure its own performance with accuracy is not a true profession.

SUMMARY AND CONCLUSIONS

The history of lines of code metrics is a cautionary tale for all people who work in software. The LOC metric started out well and was fairly effective when there was only one programming language and coding was so difficult it constituted 90% of the total effort for putting software on a computer.

But the software industry began to develop hundreds of programming languages. Applications started to use multiple programming languages and that remains the norm today. Applications grew from less than 1,000 lines of code up to more than 10,000,000 lines of code. Coding is the major task for small applications, but for large systems the work shifts to defect removal and production of paper documents in the forms of requirements, specifications, user manuals, test plans, and many others.

The LOC metric was not able to keep pace with either change. It does not work well when there is ambiguity in counting code, which always occurs with high-level languages and multiple languages in the same application. It does not work well for large systems where coding is only a small fraction of the total effort.

As a result LOC metrics became less and less useful until sometime around 1985 they started to become actually harmful. Given the errors and misunderstandings LOC metrics bring to economic, productivity, and quality studies it is fair to say that in many situations usage of LOC metrics can be viewed as professional malpractice if more than one programming language is part of the study.

READINGS AND REFERENCES ON METRICS AND FUNCTION POINT ANALYSIS

- Boehm, Barry Dr.; Software Engineering Economics; Prentice Hall, Englewood Cliffs, NJ; 1981; 900 pages.
- Brooks, Fred: The Mythical Man-Month, Addison-Wesley, Reading, Mass., 1974, rev. 1995.
- DeMarco, Tom; Why Does Software Cost So Much?; Dorset House, New York, NY; ISBN 0-9932633-34-X; 1995; 237 pages.
- Fleming, Quentin W. & Koppelman, Joel M.; Earned Value Project Management; 2nd edition; Project Management Institute, NY; ISBN 10 1880410273; 2000; 212 pages.
- Galorath, Daniel D. & Evans, Michael W.; Software Sizing, Estimation, and Risk Management: When Performance is Measured Performance Improves; Auerbach, Philadelphia, AP; ISBN 10-0849335930; 2006; 576 pages.
- Garmus, David & Herron, David; Function Point Analysis; Addison Wesley, Boston, MA; ISBN 0-201069944-3; 363 pages; 2001.
- Garmus, David & Herron, David; Measuring the Software Process: A Practical Guide to Functional Measurement; Prentice Hall, Englewood Cliffs, NJ; 1995.
- Harris, Michael D., Herron, David, and Iwanicki, Stasia; The Business Value of IT: Managing Risks, Optimizing Performance, and Measuring Results; CRC Press, Boca Raton, FL; ISBN 978-14200-6474-2; 2008; 266 pages.
- Jones, Capers; Program Quality and Programmer Productivity; IBM Technical Report TR 02.764, IBM San Jose, CA; January 1977.
- Jones, Capers; Sizing Up Software; Scientific American Magazine; New York NY; Dec. 1998, Vol. 279 No. 6; December 1998; pp 104-109.
- Jones, Capers; Applied Software Measurement; McGraw Hill, 3rd edition 2008; ISBN 978-0-07-150244-3; 575 pages; 3rd edition (March 2008).
- Jones, Capers; Software Assessments, Benchmarks, and Best Practices; Addison Wesley Longman, Boston, MA, 2000; 659 pages.
- Jones, Capers; Conflict and Litigation Between Software Clients and Developers; Version 6; Software Productivity Research, Burlington, MA; June 2006; 54 pages.
- Jones, Capers; Estimating Software Costs; McGraw Hill, New York; 2nd edition, 2007; 644 pages; ISBN13: 978- 0-07-148300-1.
- Kan, Stephen H.; Metrics and Models in Software Quality Engineering, 2nd edition; Addison Wesley Longman, Boston, MA; ISBN 0-201-72915-6; 2003; 528 pages.

- Kaplan, Robert S & Norton, David B.; The Balanced Scorecard; Harvard University Press, Boston, MA; ISBN 1591391342; 2004.
- Love, Tom; Object Lessons – Lessons Learned in Object-Oriented Development Projects; SIG Books Inc., New York NY; ISBN 0-9627477-3-4; 1993; 266 pages.
- McConnell, Steve; Software Estimation – Demystifying the Black Art; Microsoft Press, Redmond, Wa; ISBN 10: 0-7356-0535-1; 2006.
- Park, Robert E.: SEI-92-TR-20: Software Size Measurement: A Framework for Counting Software Source Statements; Software Engineering Institute, Pittsburgh, PA; 1992; 220 pages.
- Parthasarathy, M.A.; Practical Software Estimation – Function Point Methods for Insourced and Outsourced Projects; Addison Wesley, Boston, MA; ISBN 0-321-43910-4; 2007; 388 pages.
- Putnam, Lawrence H.; Measures for Excellence – Reliable Software On-Time Within Budget; Yourdon Press, Prentice Hall, Englewood Cliffs, NJ; ISBN 0-13-567694-0; 1992; 336 pages.
- Putnam, Lawrence & Myers, Ware; Industrial Strength Software – Effective Management Using Measurement; IEEE Press, Los Alamitos CA; ISBN 0-8186-7532-2; 1997; 320 pages.
- Strassmann, Paul; The Squandered Computer; Information Economics Press, Stamford, CT; 1997.
- Stutzke, Richard D.; Estimating Software-Intensive Systems – Projects, Products, and Processes; Addison Wesley, Boston, MA; ISBN 0-301-70312-2; 2005; 917 pages.
- Yourdon, Ed; Outsource – Competing in the Global Productivity Race; Prentice Hall PTR, Upper Saddle River, NJ; ISBN 0-13-147571-1; 2004; 251 pages.
- Yourdon, Ed; Death March—The Complete Software Developer’s Guide to Surviving “Mission Impossible” Projects, Prentice Hall PTR, Upper Saddle River, N.J., ISBN 0-13-748310-4, 1997.