



Research Retrospective

BOB PAIGE

New York University

The group was exciting in the 1970's, when we were groping for direction and divided by different orientations. I guess it was in this atmosphere that combined purpose with uncertainty where I found my own voice. The common goal was a transformational program development methodology that would improve productivity of designing and maintaining correct software. The emphasis was on algorithmic software.

We differed as to how to achieve this goal, and my approach was out on a limb. Based on a few transformations, the most exciting of which was Jay Earley's iterator inversion combined with high level strength reduction, and also on an overly optimistic faith in the power of science to shed light on this subject, I believed that algorithms and algorithmic software could be designed scientifically from abstract problem specifications by application of a small number of rules, whose selection could be simplified (even automated in some cases) if it could be guided by complexity. Most all others (including the SETL crowd at Courant) disagreed, and accepted the notion that algorithm design was 'inspired', and that the most significant steps in a derivation were unexplainable 'Eureka' steps.

I knew that my goals were ambitious and with little supporting evidence. In fact the case was too flimsy even to begin work on the components of a program development methodology. It seemed better to gather facts first, to uncover compelling examples that might lead to a theory, and to put together this theory only after the pieces were sufficiently understood and developed.

In order to test the viability of a new transformational idea, I tried to demonstrate how it could be used to improve some aspect of a 'conventional' area of CS; e.g., databases, algorithms, programming languages, etc. (but not transformational programming itself). A conservative, neutral test would be determined by an evaluation of the improvement (independent of the transformational idea) by members of the conventional area. A more subjective, but still useful clinical test could be made by me checking whether the new idea could improve the quality of education or facilitate my research. (I used the class room not only as a laboratory for clinical tests, but also as a way for me to learn and retool in new subject areas.) After developing the transformational idea further, it could be evaluated directly by the program transformation community, e.g., at WG2.1 meetings, other conferences, or through the publication process.

Perhaps the most important component of a transformational methodology is a wide spectrum language for expressing all levels of abstraction from problem specifications down to hardware-level implementations. Since I was a student at NYU with Jack Schwartz as my adviser, naturally I went with SETL, and found it convenient to use in my thesis work in finite differencing (from 1977–1979). As it turned out I was teaching full time (4 courses

per year) at Rutgers starting Fall 1977, and tried lots of research ideas out in the class room early on.

Despite serious semantic flaws in SETL (some of which were fixed much later in SETL2), a small subset of the language, augmented with a few standard abstract notations found in the mathematical preliminaries chapters of standard CS texts, provided students with a simple, powerful, and sometimes executable mathematical notation. SETL's small but powerful repertoire of abstract operations proved convenient in being widely reusable in modeling a variety of computer science concepts without the need for language extension.

In algorithms courses the builtin abstractions of SETL allowed data structures and rudimentary strategies of many algorithms to be specified perspicuously and without extraneous implementation detail. Default implementations of SETL's abstract operations allowed such specifications to be studied as executable prototypes, and tested empirically. Such specifications could also be used as the starting points for derivations of more efficient implementations. This made algorithm understanding easier, even (especially) with weak students who were overwhelmed by the detail of Mix-like or Algol-like implementation-level specifications. It seemed that this approach allowed me to cover much more material and convey much greater understanding.

In database courses SETL notations seemed better than the standard query languages. It facilitated teaching different data models, and different levels of description including implementations at the level of file systems and indexes. A Phillips Research TR by Lacroix and Pirotte, which compared over 100 DB queries in 9 different query languages plus English was instructive. SETL variants sketched out by Koenig and me seemed to have greater clarity. A uniform wide spectrum notation facilitated description and analysis of file and database organization, and made it easier to describe mappings from query to physical level and from one data model to another.

The wide spectrum nature of the language made it easy to illustrate transformations at and between various levels of abstraction. The first transformation that I developed was finite differencing (a generalization of conventional strength reduction plus Earley's Iterator Inversion), whose goal was to speedup programs by replacing costly repeated computations by less expensive differential counterparts.

Differencing is a mathematical idea with an old history. The chain rule mechanism and other features of finite differencing can be described independently of any particular language, but the wide spectrum nature and simplicity of SETL seemed ideal for illustrating the power and broad applicability of the transformation in a separate implementation design. Using SETL finite differencing, one can easily explain how to put together and analyze (using a natural notational form of worst case and amortized complexity) so very many efficient algorithms in a systematic way. This made teaching algorithms much easier.

One of the best examples of how surprising connections can be made, due to finite differencing, took place at the Chamrouse WG2.1 Workshop in France. We were all trying to derive one of those efficient max sequence algorithms, and the use of finite differencing at a major step suggested a similar step (to compute nested collective minimum values) in the derivation of Dijkstra's SSSP algorithm.

During these early years there was also one interesting objective test of finite differencing in the area of databases. In his Turing award address Codd said that integrity control was

an important open problem in relational databases. Successful use of finite differencing to solve view maintenance and integrity control was reported in [13] and by Paige at a workshop in Toulouse, in 1983.

In Compiler classes at Rutgers I also developed a crude form of dominated convergence in order to derive workset algorithms found in Hecht's book [10] for solving global program analysis problems specified by fixed points. Work on this transformation progressed as it was seen to be progressively more useful in deriving an increasingly wider range of algorithms. And I continued to use it in combination with finite differencing in my compiler lectures in order to derive the more difficult algorithms. However, it was not until my collaboration with Cai that a comprehensive investigation and development of this transformation was first completed in [4].

It became apparent from the beginning that I was looking for a transformational program development methodology whose final step would be data structure selection beyond which lay conventional compiler optimization. I also favored an approach limited to composable simple data structures as are found in SETL's method of basings and Wiederhold's file system design methodology instead of an expert system style as proposed by Barstow and Kant. Despite some ingenious ideas, SETL data structure selection was largely ad hoc, not amenable to formal complexity analysis, overrelying on hashing, missing a linked list, and never proven (either analytically or empirically) to be an improvement over naive unoptimized data structures.

Unfortunately SETL data structure selection wasn't good enough to produce high performance code by a complexity driven approach. What was needed were data structures guaranteed to support associative access (e.g., set membership testing) in unit worst case time for search arguments and set elements of arbitrary type. Also needed were principles for selecting such data structures. It has taken many years to develop this transformation, and a first comprehensive treatment will appear in Deepak Goyal's thesis (expected in Fall 1999). Only recently was the transformation equipped with a suitable read method to create the desired data structures. And I still do not understand how to eliminate enough overhead explaining even very simple forms of data structure selection (described in several papers cited earlier) to use this transformation conveniently in a standard algorithms class.

Nevertheless, I believe that it is an essential final stage of and underpinning to a transformational methodology that (1) begins with dominated convergence to compute fixed point specifications, (2) segues into finite differencing to implement repeated expensive high level expressions more efficiently by exploiting the differential maintenance of program invariants, and, finally, (3) uses data structure selection by real-time simulation to implement associative access operations in unit worst case time on a RAM.

Over the years we have uncovered extensive evidence that our methodology will scale up and be capable of improving the productivity of large-scale system implementation. Three kinds of evidence have been considered. Our SETL-based specification languages and the transformations that implement them have been shown to be effective in terms of (1) explaining complex algorithms, (2) discovering new algorithms, and, most importantly, (3) improving the productivity of efficient implementations of algorithms.

Most mainstream research in Program Transformations emphasizes principles for constructing derivations with illustrations drawn from known algorithms. The purpose has

been to uncover common patterns of abstraction in specification and transformation that could form the basis of a useful methodology for making algorithm design and program development easier. Short, elegant transformational proofs (of correctness integrated with analysis) document implementations of complex algorithms, lend greater confidence in the correctness of complex codes, and provide greater assurance in the reliable modification of such codes.

Perhaps the first examples of nontrivial algorithms being derived by finite differencing were presented in [14, 16]. Included among these examples is a SETL specification of Dijkstra's naive Bankers Algorithm and its transformation into Habermann's efficient solution. This derivation was done without knowledge of Habermann's solution, and there were no dead ends. Finite differencing homed right in on a solution matching Habermann's time/space bounds.

It is well known that the construction of optimizing compilers is a costly labor intensive task. Can this labor be reduced by our methods? To answer this question in part, we showed how easy it was to specify dozens of programming language and compiler analysis problems in SQ2+, a functional language composed from the expressions of SETL augmented with least and greatest fixed point expressions. It was also shown how to simplify these specifications, and to transform them by dominated convergence into high level SETL prototypes [4]. In [3] we showed how the constant propagation algorithm of [19] could be expressed as set-theoretic equations in a subset of SQ2+ that could be mapped into RAM code guaranteed to run in linear time in the size of the program dataflow relation.

The use of notation has been regarded as a burden to algorithm designers ever since Knuth came out with Mix [12]. But can notation also help satisfy the needs of the algorithm community – precise algorithmic analysis and succinct exposition? I was thrilled that coauthor Bob Tarjan agreed to explain our new linear time solution to the Single Function Coarsest Partition Problem as being derived from an SQ2+ specification by dominated convergence, finite differencing, and real-time simulation [17] (a special TCS issue of best papers selected from ICALP84). In [18] a much improved explanation of the algorithmic tool called Multiset Discrimination in [6] was obtained by specifying the algorithm in low level SETL and using its type system to formally explain and analyze the low level implementation that would be obtained by real-time simulation. The earlier presentation of this algorithmic device involved so much indirection from pointer-based primitives, that most readers were confused. Some of these earlier readers agreed that the low level SETL presentation clarified their understanding.

The viability of a transformational methodology can be demonstrated by using it to 'explain' or 'prove' well known algorithms. However, if in the course of such formal explication no new deep structure is uncovered that leads to improved solutions, then its impact on algorithmics and programming productivity is likely to be limited. More powerful evidence favoring a transformational methodology is provided if it can be demonstrated to help facilitate the discovery of new algorithms. Our success with algorithm discovery may be attributed to our reliance on complexity in both specification and transformation.

Our first instance of algorithm discovery by transformation was reported in [15], where we used all three transformations to turn an SQ2+ specification of Horn Clause Propositional Satisfiability into a new linear time pointer machine algorithm. Previously, Dowling and

Gallier found a linear time algorithm [8] that relied heavily on array access. In [2] we used dominated convergence and finite differencing to derive a low level SETL executable prototype from an SQ2+ specification of ready simulation. We then showed informally how the low level SETL prototype could be turned into an algorithm that runs 5 orders of magnitude faster than the previous solution in [1]. All three transformations, but especially real-time simulation (where types were shown to be useful in modeling complex data structures), were involved in the discovery of a new improved solution to the classical problem of DFA minimization [11]. Finite differencing was used extensively in [7] to derive a new improved solution to the classical problem of turning regular expressions into DFA's. Perhaps our most convincing paper-and-pencil result was in [9], where Goyal and I used low level SETL specifications and real-time simulation to improve Willard's time bound for query processing from linear expected to linear worst-case time without degrading space. Willard's original algorithm was extremely difficult, and involved over 80 pages of proofs. Our transformational approach yielded much shorter but also more precise constructive proofs that led to an implementation design. Here is a first successful example of scaling up.

Two summers ago, Ph.D. student Deepak Goyal designed and implemented 'practical' algorithms in Java at Microsoft. He believes that his use of low level SETL and the data structure selection transformation as part of a programming methodology increased his productivity and improved the quality of the code that was produced.

Perhaps the most compelling evidence that our transformational methodology will scale up and provide a dramatic improvement in the productivity of large high performance complex systems may be found in the experiments by Cai and Paige [5]. In that paper we developed a simple but conservative model of productivity. Within that model we demonstrated a five-fold improvement in productivity of high performance algorithm implementation in C in comparison to hand-coded C.

Those experiments tested an approach to producing C programs by writing programs in a statically typed simple variant of low level SETL, and compiling them into high performance C. The high performance of the C code produced by the SETL-to-C translator is based on the translator's ability to simulate associative access (e.g., finite set membership or finite map application) on a RAM in real time.

Measuring productivity improvement depended on two assumptions. The first is that one line of low level SETL takes no more time to compose than one line of C. This assumption is not controversial. Our experience is that one line of low level SETL can be produced faster than a line of C. The generally lower level of discourse in C creates an intellectual gap between the program and the mathematical function it computes. For example, in order to implement SETL's element deletion operation (`s less := x`) efficiently in C would require at least 10 carefully chosen C operations. The need to access data through pointers and cursors in C creates a greater level of indirection (which complicates understanding) than in SETL, where data is accessed directly through values.

The second assumption is that the C code generated automatically from low level SETL has roughly the same number of lines as equivalent hand-coded C. We found that the generated C was between 10% and 30% larger than hand-coded C, and concluded that the errors in the two assumptions would cancel each other out.

Suppose we measure programming productivity in a given programming language as the number of pretty-printed source lines produced per unit time. Suppose also that productivity decreases as the conceptual difficulty in understanding a program grows. Then in our investigation, which is restricted to highly algorithmic programming (as is found in complex language systems and environments), conceptual difficulty is roughly reflected in the size of the dataflow relation, which can be expected to grow nonlinearly with the number of source lines. Thus, programming productivity $P(L)$, as a function of the number of program source lines L , increases as L decreases.

Let L_2 be the size of a C program C_2 compiled from a low level SETL specification S_1 of size L_1 , and let C_3 be a hand-crafted C program equivalent to C_2 . By assumption one, we can use the same productivity function $P(L)$ for programs written in both low level SETL and C. By assumption two, we know that the size of C_3 is roughly the same as the size L_2 of C_2 . Then the improvement in productivity by using our automated program development methodology versus hand-crafted programming is given by the time $L_2/P(L\{2\})$ to manually produce C_2 divided by the time $L_1/P(L\{1\})$ to manually produce S_1 , which is

$$(L_2/L_1)(P(L\{1\})/P(L\{2\})) > L_2/L_1$$

since $P(L\{1\})/P(L\{2\}) > 1$ whenever $L_1 < L_2$.

It should be emphasized that our experiments did not measure productivity directly, which would have required difficult and costly analysis of human factors such as programming expertise and intelligence. Instead we measured productivity improvement by exploiting the two assumptions mentioned above to obtain an objective, inexpensive, and credible framework for conducting comparative productivity experiments that could avoid all human factors. The ratio of lines of C code generated automatically from low level SETL divided by the number of lines of low level SETL being compiled gives a lower bound on productivity improvement. Every algorithm that we tested in [5] yielded ratios that exceeded 5, and that grew as the input size grew.

Although we found that the C code generated automatically from low level SETL had runtime performance comparable to good hand code (whose running time, excluding I/O, was at least 30 times faster than SETL2 running time executed by the standard SETL2 interpreter), only small-scale examples were used. The largest C program was about 10 pages. Low level SETL lacked procedures, and the type system was highly restricted. The SETL-to-C production rate was too slow—about 24 lines of C per minute on a SPARC 2. Finally, high level SETL input had to be translated into low level C input at compilation time.

The hypothesis that productivity improvement will scale up by our methods is based on the methodology and the application domain. Real-time simulation is applied uniformly to each instruction of a program regardless of program size. Algorithms for doing program analysis and transformation have a complex combinatorial nature, which fits our model of productivity.

Over the last few years, I've felt that, finally, there was enough evidence and know-how to put together much of our work into a formal program development methodology that could achieve my goals of 20 years ago. Although I was too sick to carry this work out, I'm happy that my student Deepak Goyal has done it in his Ph.D. thesis, which should be

finished some time over the summer. Deepak ought to present it at a future meeting. Annie Liu, who's on Deepak's Committee, can tell you more, or you may look at Deepak's home page, which is in the Ph.D. Student area of our NYU Dept. home page. I think you'll like it.

For anyone interested, Deepak will be working in John Field's group at IBM TJ Watson starting Nov. 1. Until then, he'll be at NYU.

References

1. Bloom, B. Ready simulation, bisimulation, and the semantics of CCS-like languages. Ph.D. Thesis, Massachusetts Institute of Technology, 1989.
2. Bloom, B. and Paige, R. Transformational design and implementation of a new efficient solution to the ready simulation problem. *Science of Computer Programming*, **24**(3) (1995) 189–220.
3. Cai, J. and Paige, R. Binding performance at language design time. In *Proceedings of the Fourteenth Annual ACM Symposium on Principles of Programming Languages*, M.J. O'Donnell (Ed.). ACM Press, München, West Germany, 1987, pp. 85–97.
4. Cai, J. and Paige, R. Program derivation by fixed point computation. *Science of Computer Programming*, **11**(3) (1989) 197–261.
5. Cai, J. and Paige, R. Towards increased productivity of algorithm implementation. In *Proc. ACM SIGSOFT FSE*, 1993, pp. 71–78.
6. Cai, J. and Paige, R. Using multiset discrimination to solve language processing problems without hashing. *Theoretical Computer Science*, **145**(1–2) (1995) 189–228.
7. Chang, C.-H. and Paige, R. From regular expressions to DFA's using compressed NFA's. *Theoretical Computer Science*, **178**(1/2) (1997) 1–36.
8. Dowling, W.F. and Gallier, J. Linear-time algorithms for testing the satisfiability of propositional horn formulae. *Journal of Logic Programming*, **1**(3) (1984) 267–284.
9. Goyal, D. and Paige, R. The formal reconstruction and improvement of the linear time fragment of Willard's relational calculus subset. In *Algorithmic Languages and Calculi*, R. Bird and L. Meertens (Eds.). Chapman & Hall, 1997, pp. 382–414.
10. Hecht, M.S. *Flow Analysis of Computer Programs*, Elsevier North-Holland, 1977.
11. Keller, J.P. and Paige, R. Program derivation with verified transformations—A case study. *CPAM*, **48**(9/10) (1996) 1053–1113.
12. Knuth, D.E. *The Art of Computer Programming, Volume II: Seminumerical Algorithms*. Addison-Wesley, 1969.
13. Koenig, S. and Paige, R. A transformational framework for the automatic control of derived data. In *Proceedings of the 7th International Conference on Very Large Data Bases*. Cannes, France, 1981, pp. 306–318.
14. Paige, R. *Formal Differentiation*. UMI Research Press, 1981.
15. Paige, R. and Henglein, F. Mechanical translation of set theoretic problem specifications into efficient RAM code—A case study. *Journal of Symbolic Computation*, **4**(2) (1987) 207–232.
16. Paige, R. and Koenig, S. Finite differencing of computable expressions. *ACM Transactions on Programming Languages and Systems*, **4**(3) (1982) 402–454.
17. Paige, R., Tarjan, R.E. and Bonic, R. A linear time solution to the single function coarsest partition problem. *Theoretical Computer Science*, **40**(1) (1985) 67–84.
18. Paige, R. and Yang, Z. High Level reading and data structure compilation. In *Proceedings of the Twenty-Fourth Annual ACM Symposium on Principles of Programming Languages*, N.D. Jones (Ed.). ACM Press, Paris, France, 1997, pp. 456–469.
19. Reif, J.H. and Lewis, H.R. Symbolic evaluation and the global value graph. In *Proceedings of the Fourth Annual ACM Symposium on Principles of Programming Languages*, R. Sethi (Ed.). ACM Press, 1977, pp. 104–118.