

MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

**An Empirical Study of
Insertion and Deletion in Binary Search Trees**

Jeffrey L. Eppinger

**Department of Computer Science
Carnegie-Mellon University
Pittsburgh, Pennsylvania 15213**

December 2, 1982

AD A125210

**DEPARTMENT
of
COMPUTER SCIENCE**

DTIC FILE COPY



**DTIC
ELECTE
MAR 05 1983
E**

Carnegie-Mellon University

This document has been approved
for public release; its sale is

88 03 02 014

**An Empirical Study of
Insertion and Deletion in Binary Search Trees**

Jeffrey L. Eppinger

**Department of Computer Science
Carnegie-Mellon University
Pittsburgh, Pennsylvania 15213**

December 2, 1982



Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By _____	
Distribution/ _____	
Availability Codes	
Dist	Avail and/or Special
A	

Abstract: This paper describes an experiment on the effect of insertions and deletions on the path length of unbalanced binary search trees. Given a random binary tree, repeatedly inserting and deleting nodes yields a tree that is no longer random. The expected internal path length differs when different deletion algorithms are used. Previous empirical studies indicated that expected internal path length tends to *decrease* after repeated insertions and asymmetric deletions. This study shows that performing a larger number of insertions and asymmetric deletions actually *increases* the expected internal path length, and that for sufficiently large trees, the expected internal path length becomes *worse* than that of a random tree. However, with a symmetric deletion algorithm, the results indicate that performing a large number of insertions and deletions *decreases* the expected internal path length, and that the expected internal path length remains better than that of a random tree.

This research was sponsored in part by the Office of Naval Research under contract N00014-76-C-0370.

- a -

1. Introduction

A binary tree created by inserting n randomly chosen keys into an empty tree has an expected internal path length of $I_n \approx 1.386n \lg n$.[†] Randomly deleting k nodes from such a tree yields a tree whose expected internal path length is I_{n-k} . Unfortunately, performing insertions after deletions does not produce binary trees whose internal path length is predicted by this function. A theoretical explanation of the effect of performing deletions and then insertions on binary trees is still lacking. [Knuth 73, Section 6.2.2]

This paper presents an empirical study on the effect of applying random insertions and deletions to random binary search trees and analyzes results of experiments comparing asymmetric and symmetric deletion algorithms. In a previous empirical study, Knott [Knott 75] suggests that the expected internal path length tends to *decrease* after repeated insertions and asymmetric deletions. In this study, the large number of insertions and asymmetric deletions performed suggests that the expected internal path length first decreases but eventually begins to *increase*. For sufficiently large trees, expected internal path length becomes *worse* than that of a random tree. However, experiments using the symmetric deletion algorithm show that performing a large number of insertions and symmetric deletions *decreases* the expected internal path length (making the trees better than random).

Section 2 describes the insertion and deletion algorithms used in this study and provides an overview of some of the previous work in this area. The statistics used in this study are defined in Section 3. Section 3 also mentions a few specifics about how the data was gathered. The observations in Section 4 give an interpretation of the data and the conclusions are summarized in Section 5.

2. Background

Insertion Algorithm: The structure of binary trees naturally leads to one insertion algorithm. To insert a node into a binary tree (known not to contain the node), compare the new and current keys and insert the node into the left or right subtree, whichever maintains the invariant of the data structure. The Pascal code for this algorithm is provided in Figure 1, below. For further

[†] Throughout this paper, $\lg x$ denotes $\log_2 x$.

```

PROCEDURE Insert(VAR root : NodePtr; x : DataType);
BEGIN
  IF root = NIL
  THEN BEGIN
    NEW(root); root^.data := x;
    root^.lChild := NIL; root^.rChild := NIL
  END
  ELSE IF x < root^.data
  THEN Insert(root^.lChild, x)
  ELSE Insert(root^.rChild, x)
END;

```

Figure 1: The insertion procedure.

explanation see [Knuth 73, Section 6.2.2, Algorithm T].

Unlike insertion, there are many reasonable deletion algorithms from which to choose. This paper describes experiments with Knuth's asymmetric deletion algorithm and a trivially modified version of this algorithm to make it symmetric.

Asymmetric Deletion Algorithm: A node's *successor* is defined to be the smallest node in the right subtree. Similarly a node's *predecessor* is defined to be the largest node in the left subtree. To delete a node from a binary tree, replace the node with its successor, i.e., the node that contains the next larger key. The Pascal code for this algorithm is given in Figure 2, below. Figure 4[†] shows examples of the insertion algorithm and this deletion algorithm applied to a particular binary tree; for further explanation see [Knuth 73, Section 6.2.2, Algorithm D].

Symmetric Deletion Algorithm: To delete a node from a binary tree, replace the node with its successor or predecessor. Alternately choose the successor and predecessor (so that half the time the *RightDelete* routine is called and half the time a suitably modified version of this routine, *LeftDelete*, is called).

Consider building a binary tree using n keys chosen randomly from a uniform distribution (i.e., all $n!$ permutations of the keys are equally likely). There are $\binom{2^n}{n}/(n+1)$ possible shapes for this tree [Knuth 68, Section 2.3.4.4], each with some probability of occurring; call the distribution D_n . By this definition, inserting a new node into this binary tree would yield a tree of size $n+1$ whose shape occurs with a probability defined by D_{n+1} . Binary trees whose distribution of shapes

[†] Figures 4-11 are at the end of the paper.

```

PROCEDURE RightDelete (VAR root : NodePtr; x : DataType);
  VAR copy, successor, succPtr : NodePtr;
BEGIN
  IF x < root↑.data
    THEN RightDelete (root↑.lChild, x)
  ELSE IF x > root↑.data
    THEN RightDelete (root↑.rChild, x)
  ELSE BEGIN
    copy := root;
    IF root↑.rChild = NIL
      { Case I: There is no successor. }
      THEN root := root↑.lChild
    ELSE IF root↑.rChild↑.lChild = NIL
      { Case II: The successor is the right child. }
      THEN BEGIN
        root↑.rChild↑.lChild := root↑.lChild;
        root := root↑.rChild
      END
    { Case III: The successor is the leftmost child in the right subtree. }
    ELSE BEGIN
      succPtr := root↑.rChild;
      WHILE succPtr↑.lChild↑.lChild <> NIL DO
        succPtr := succPtr↑.lChild;
      successor := succPtr↑.lChild;
      succPtr↑.lChild := successor↑.rChild;
      successor↑.lChild := root↑.lChild;
      successor↑.rChild := root↑.rChild;
      root := successor
    END;
    DISPOSE (copy)
  END
END;

```

Figure 2: The asymmetric deletion procedure.

is D_n are called *random binary trees*.

Thomas Hibbard [Hibbard 62] proved that deleting a random node (i.e., where each node has an equal probability of being deleted) from a binary tree of size n , with distribution of shapes D_n , yields a tree with a distribution of shapes D_{n-1} .

Strangely, performing random insertion and deletion operations on a random tree does not preserve this distribution of shapes. Consider building a binary tree of size n , as described above. Since the keys are chosen from a uniform distribution, the probability of inserting a new node in any particular interkey gap is $\frac{1}{n+1}$. After one random deletion, the distribution of shapes will be

D_{n-1} , but the probability of inserting a new node where the deleted node used to be will be $\frac{2}{n+1}$ (while all other places are still $\frac{1}{n+1}$). Knuth [Knuth 73, Section 6.2.2] describes this phenomenon as follows:

The shape of the tree is random after deletions, but the relative distribution of values in a given tree shape may change, and it turns out that the first random insertion after a deletion actually *destroys* the randomness property on shapes. This startling fact, first observed by Gary Knott in 1972, must be seen to be believed. Empirical evidence suggests strongly that the path length tends to *decrease* after repeated deletions and insertions, so the departure from randomness seems to be in the right direction; a theoretical explanation for this behavior is still lacking.

Knuth feels that binary trees tend to improve because "path length tends to decrease." One way to compare binary trees is to measure their internal path lengths. The internal path length of a tree is defined as the sum of the depths of the nodes in the tree,

$$IPL = \sum_{i \in \{\text{nodes}\}} \text{distance}(\text{root}, i).$$

For a random tree containing n nodes, the expected IPL is denoted as I_n and the expected number of comparisons in a successful search is denoted as C_n . Knuth [Knuth 73, Section 6.2.2] gives the expected number of comparisons in a successful search, C_n , as approximately equal to $1.386 \lg n$. Substituting into the relation $I_n = n(C_n - 1)$, one obtains the approximation $I_n \approx 1.386n \lg n$. A distribution of trees is said to be "better than random" when the expected IPL is less than I_n (since the expected number of comparisons is proportional to the IPL).

3. Methodology

If a random sequence of insertions and deletions were applied to a random tree of size n , the resulting tree would probably not have the same number of nodes. The original tree's IPL would therefore not be directly comparable with the IPL of the new tree. In this study, sequences of *insertion/deletion pairs* (I/D pairs) are applied to random trees. Since the resulting tree always has the same size, it is easy to see whether any improvement has been made. (Knott's data was also obtained by using I/D pairs.) The first step of the simulation is therefore to insert n nodes into an empty tree, after which successive pairs of insertions followed by deletions are performed.

Let $IPL_{n,i}$ denote the measured mean IPL of an n -node binary tree after applying i I/D pairs.

Figures 5 through 10 show $IPL_{n,i}/I_n$ plotted as a function of i . This ratio shows the improvement of the resulting tree's expected IPL as a fraction of the random tree's expected IPL.

The deletion algorithm given above generally replaces the node to be deleted with its *successor*, the "left-most node in the right subtree". The left and right subtrees are treated differently and, as observed below, this appears to have a profound affect on the behavior of binary trees. Such a deletion algorithm is called an *asymmetric* deletion algorithm. The *symmetric* deletion algorithm which is examined in this study is a trivially modified version of the asymmetric algorithm. This symmetric algorithm *alternately* replaces the node to be deleted with its successor or its predecessor. The algorithm requires a small amount of state information, but similar results have been obtained by *randomly* replacing the node to be deleted by its successor or predecessor.

To ensure that the results were not an artifact of the random number generator, simulations were performed on both DEC-20s and Perqs. In the DEC-20 simulations the random number generator used the linear congruential method to produce 36-bit pseudorandom numbers [Knuth 69, Section 3.2]. The random number generator for the Perqs is the feedback shift-register pseudorandom number generator as described in [Lewis 73]. The data presented in this paper was generated on the Perqs and took about one month of CPU time, but similar results were obtained for the smaller trees on the DEC-20s.

The outer loop of the simulation program is very simple. First, build a tree with $tsize$ nodes, then gather data before and after each interval of $isize$ I/D pairs.

```
FOR i := 1 TO tsize DO RndInsert;
... gather data ...
FOR i := 1 TO intervals DO BEGIN
  FOR j := 1 TO isize DO BEGIN RndInsert; RndDelete END;
  ... gather data ...
END;
FreeTree;
```

Figure 3: The inner loop of a simulation.

4. Observations

The graphs in Figures 5 and 6 show the expected internal path length of n -node binary trees plotted against the number of insertion and asymmetric deletion pairs. Initially, $IPL_{n,i}$

decreases, as Knott and Knuth observed. After some critical point, though, $IPL_{n,i}$ starts to increase, eventually levelling off after approximately n^2 I/D pairs. Figure 7 is a comparison chart in which $IPL_{n,i}/I_n$ is plotted as a function of i/n^2 for each of the values of n tested. (The latter ratio normalises the x-axis.)

Perhaps the most significant observation is that as n increases so does the asymptotic value for $IPL_{n,i}/I_n$. Since binary trees can be modeled by Markov Chains, and any binary tree may be obtained by applying some combination of I/D pairs to any other binary tree, the $\lim_{i \rightarrow \infty} IPL_{n,i}$ exists [Ross 70, Theorem 4.9]. Figure 7 suggests that

$$\lim_{i \rightarrow \infty} IPL_{n,i} > I_n$$

for sufficiently large values of n (roughly greater than 128). Thus binary trees seem to become "worse than random" after many insertions and deletions.

The comparison chart in Figure 11 shows the asymptotic values of $IPL_{n,i}/I_n$ for both deletion algorithms plotted against n (on a log scale). The data given in Table 1 was obtained by summing all the $IPL_{n,i}$ and $IPL_{n,i}^2$, when $i > n^2$.

n	Samples	$IPL_{n,i > n^2}$	Variance
64	6000	0.97	0.01652
128	6800	1.00	0.01340
256	2300	1.06	0.00985
512	1200	1.16	0.00970
1024	750	1.30	0.01013
2048	5340	1.49	0.00771

Table 1: Data for Asymmetric Deletions.

The asymmetric curve appears to be quadratic. A least-squares multiple regression weighted by the inverse of the variance yields the following approximation:

$$\lim_{i \rightarrow \infty} \frac{IPL_{n,i}}{I_n} \approx 0.0202 \lg^2 n - 0.241 \lg n + 1.69.$$

Substituting $I_n \approx 1.386n \lg n$ we obtain

$$\lim_{i \rightarrow \infty} IPL_{n,i} \approx 0.0280n \lg^3 n - 0.334n \lg^2 n + 2.34n \lg n.$$

The graphs in Figures 8 and 9 show the corresponding plots of the data in Table 2 for the expected internal path length for symmetric deletions.

n	Samples	$IPL_{n,i}$	Variance
64	6000	0.905	0.01654
128	6800	0.890	0.00916
256	2300	0.888	0.00615
512	1200	0.890	0.00347
1024	750	0.881	0.00235
2048	5340	0.883	0.00269

Table 2: Data for Symmetric Deletions.

The $IPL_{n,i}$ decreases initially, as in the case of asymmetric deletions, but the asymptotic value of the expected internal path length seems to remain lower than that of a random tree. The comparison charts in Figures 10 and 11 indicate that

$$1 > \lim_{i \rightarrow \infty} \frac{IPL_{n,i}}{I_n} \approx 0.88$$

or that

$$I_n > \lim_{i \rightarrow \infty} IPL_{n,i} \approx 1.22n \lg n.$$

The comparison chart in Figure 11 shows the asymptotic value of $IPL_{n,i}$ slowly decreasing as n increases. Since a binary tree with n nodes cannot have an internal path length less than that of a perfect tree, we know that

$$\lim_{i \rightarrow \infty} IPL_{n,i} = \Omega(n \log n).$$

5. Conclusions

The expected internal path length of a random binary tree is $I_n = \Theta(n \log n)$. Empirical evidence suggests that performing many insertion and asymmetric deletions yields binary trees with an expected internal path length of $IPL_{n,i} = \Theta(n \log^3 n)$. Thus performing asymmetric deletions causes binary trees to become more unbalanced. Amazingly, the expected path length does not increase by a constant factor, but rather by a factor of $\log^3 n$. However, experiments show

that the symmetric deletion algorithm improves the balance of binary trees leaving the expected internal path length $\Theta(n \log n)$, but with a *smaller* constant coefficient than the expected internal path length of a random binary tree.

Because this is an empirical study, the above conclusions can only be conjectures. No one has provided a theoretical explanation of the behavior of a binary tree's path length after applying deletions and then insertions. There is no proof that the asymptotic value of $IPL_{n,i}$ is less than I_n when performing random insertions and symmetric deletions or that the asymptotic value of $IPL_{n,i}$ is greater than I_n when applying insertions and asymmetric deletions.

In closing, it should be noted that the results of this study will have little impact on the use of binary trees in practice. It takes approximately 1.5 million random insertions and asymmetric deletions to make a 2048-node binary tree worse than a random tree, and 4 million before its expected internal path length reaches the asymptotic value (which is just 50% worse). When so many operations are required, other data structures are probably more appropriate.

6. Acknowledgements

I would like to thank Jon Bentley, James Gosling, Diane Lambert, and Jim Saxe for their help and guidance.

7. References

- [Hibbard 62] Hibbard, Thomas N.
Some Combinatorial Properties of Certain Trees
with Applications to Searching and Sorting.
Journal of the Association of Computing Machinery 9(1):13-28, January 1962.
- [Knott 75] Knott, Gary D.
Deletion in Binary Storage Trees.
Ph.D. thesis, Stanford University, May, 1975.
STAN-CS-75-491.
- [Knuth 68] Knuth, Donald E.
The Art of Computer Programming.
Volume I: *Fundamental Algorithms.*
Addison-Wesley, 1968, Section 2.3.4.4.
- [Knuth 69] Knuth, Donald E.
The Art of Computer Programming.
Volume II: *Seminumerical Algorithms.*
Addison-Wesley, 1969, Section 3.2.
- [Knuth 73] Knuth, Donald E.
The Art of Computer Programming.
Volume III: *Searching and Sorting* (Second Printing, March 1975).
Addison-Wesley, 1973, Section 6.2.2.
Note: The Second Printing contains important changes in Section 6.2.2.
- [Lewis 73] Lewis, T. G., and W. H. Payne
Generalized Feedback Shift Register Pseudorandom Number Generator
Journal of the Association of Computing Machinery 20(3):456-468, July 1973.
- [Ross 70] Ross, Sheldon M.
Applied Probability Models with Optimization Applications.
Holden-Day, 1970, Section 4.3.

Figure 4: Examples of Insertion and Asymmetric Deletion.

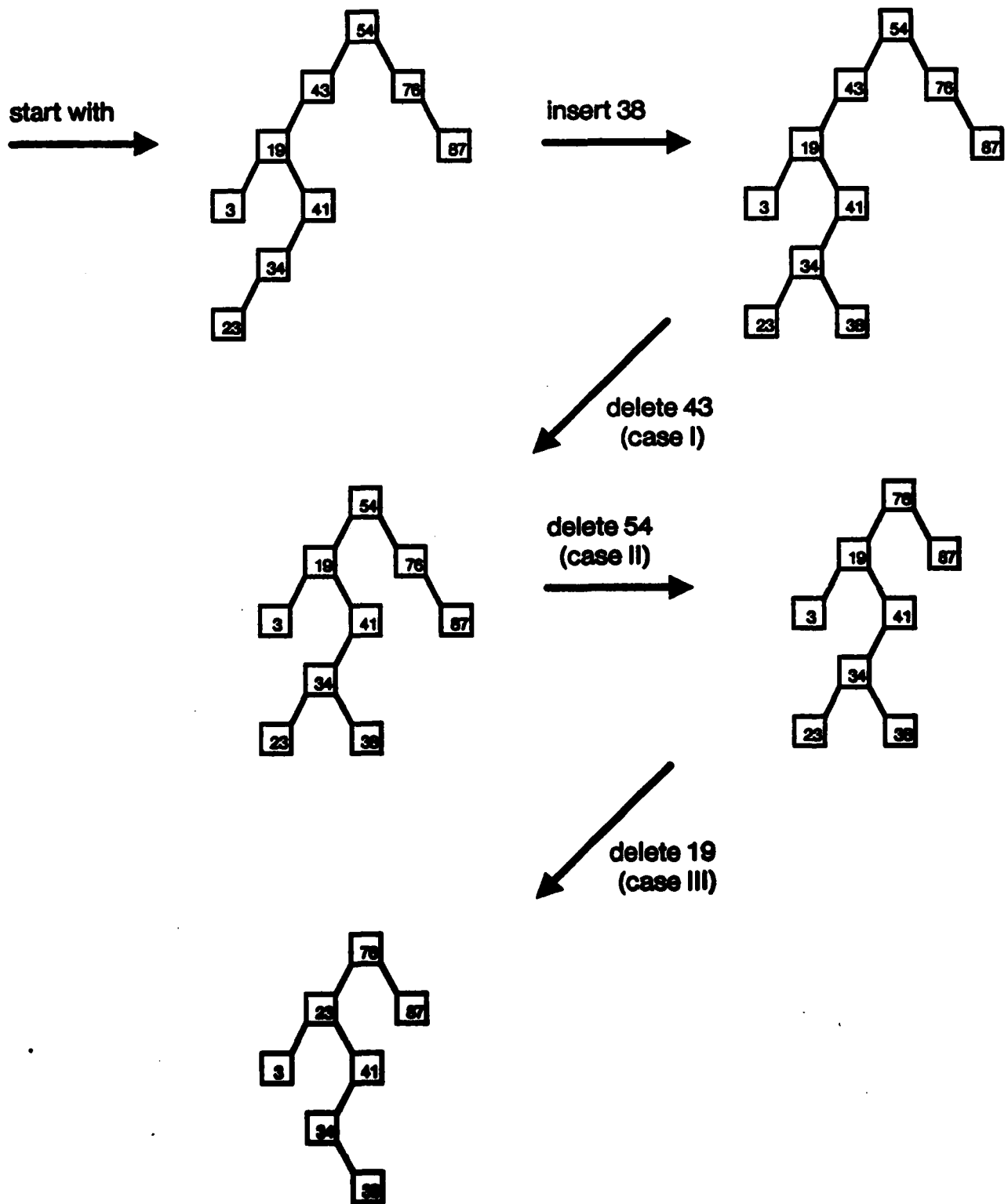


Figure 5

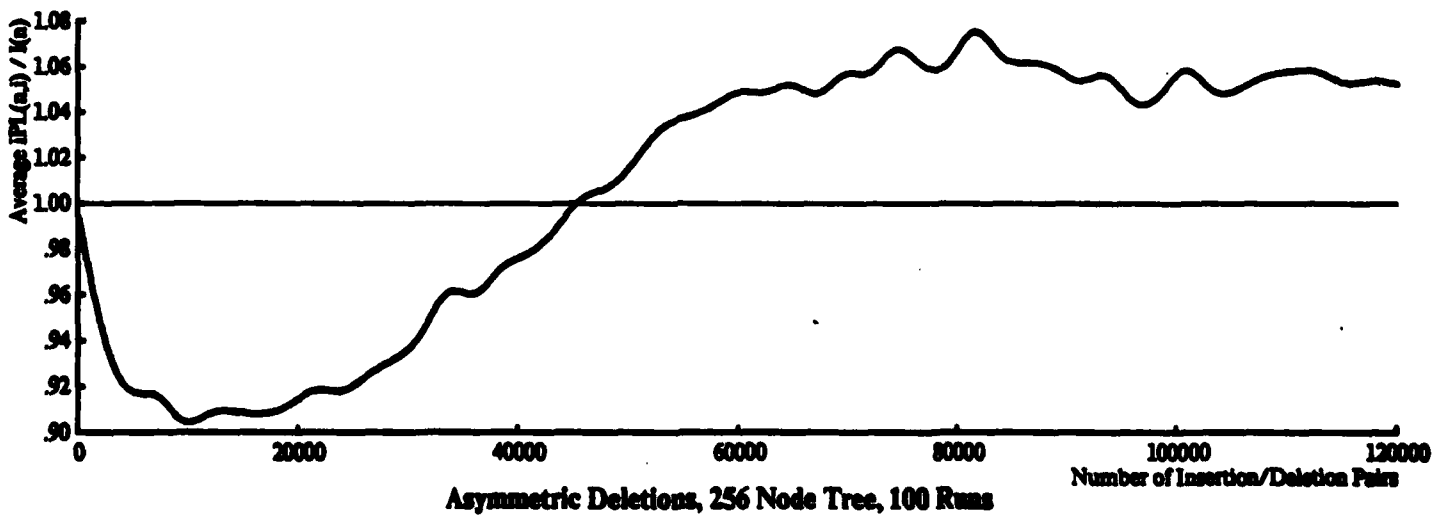
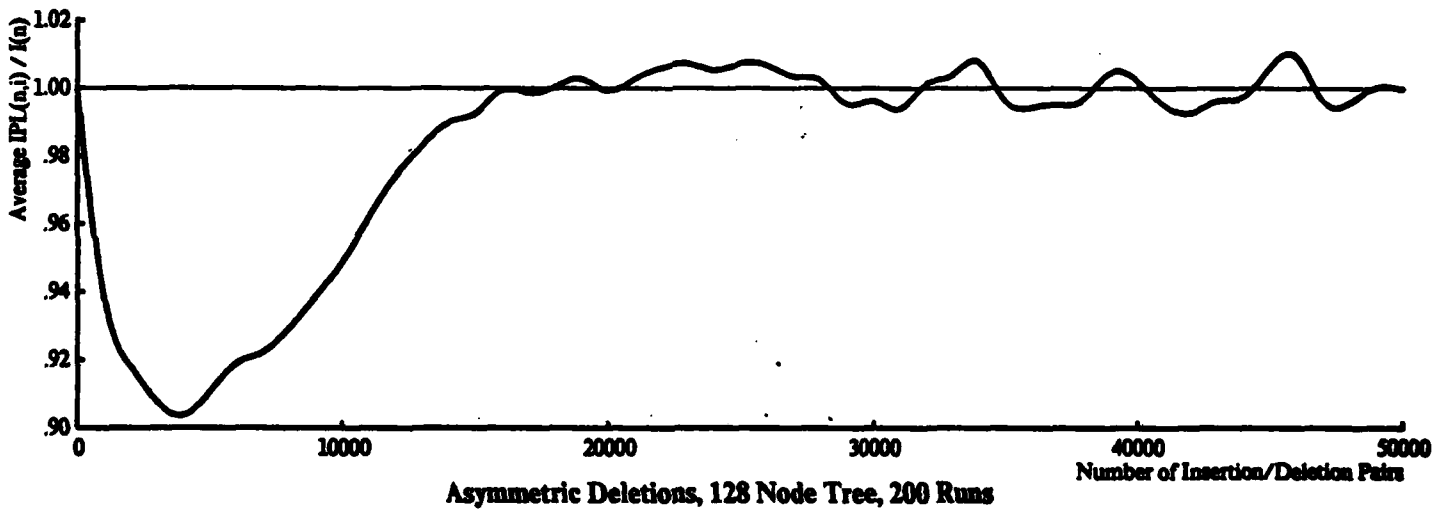
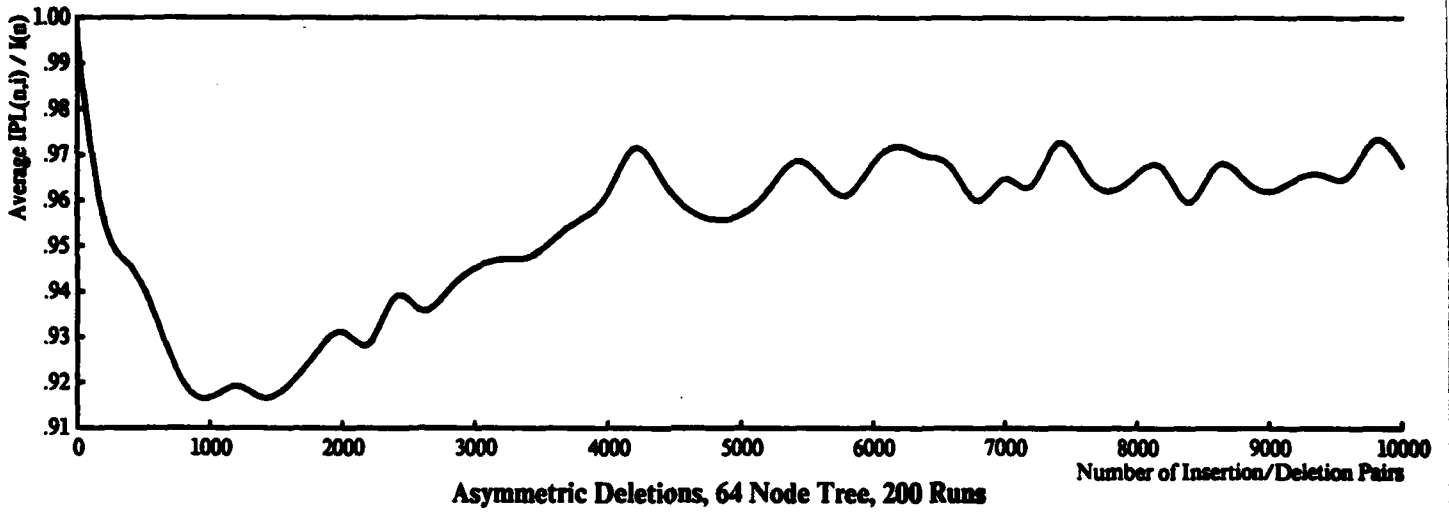


Figure 6

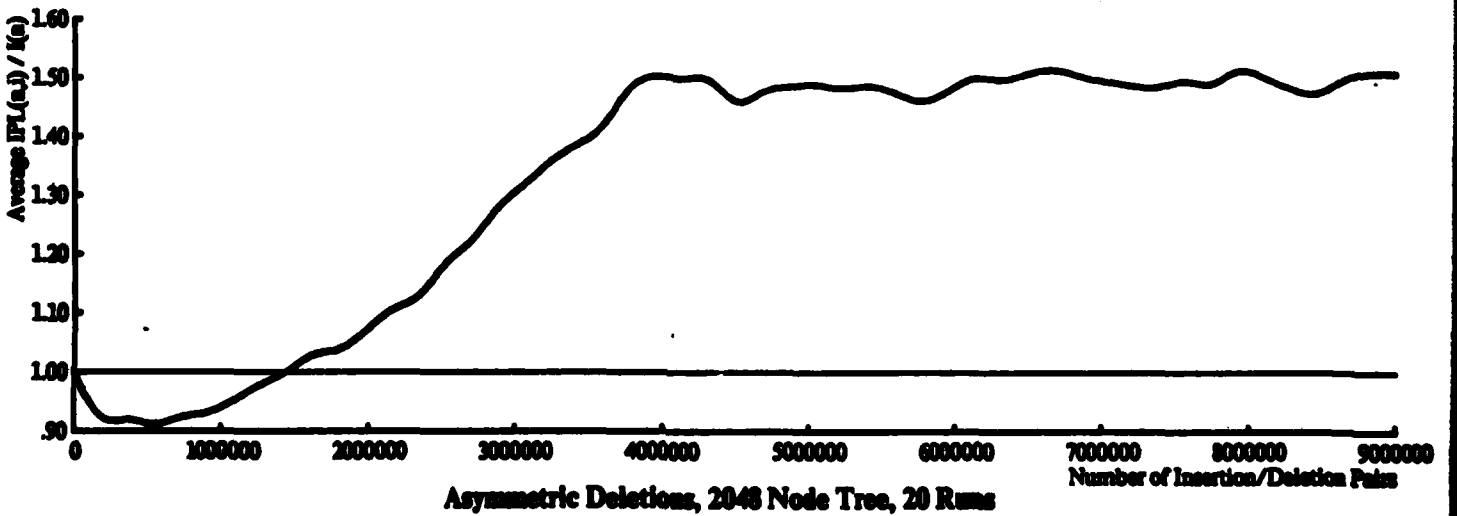
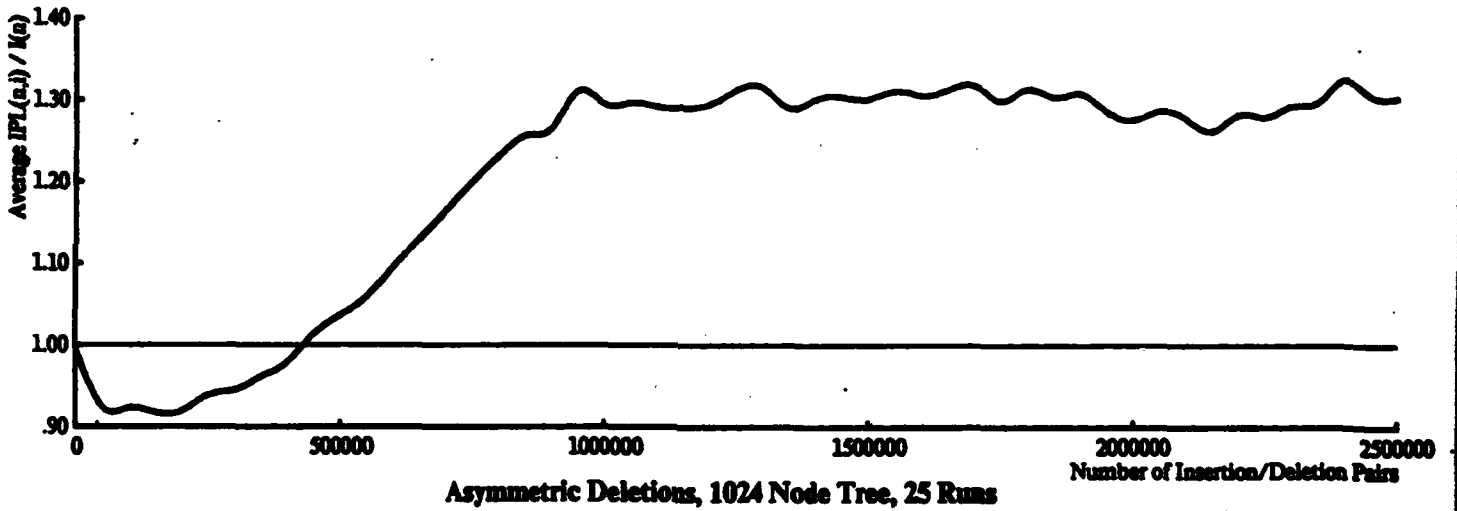
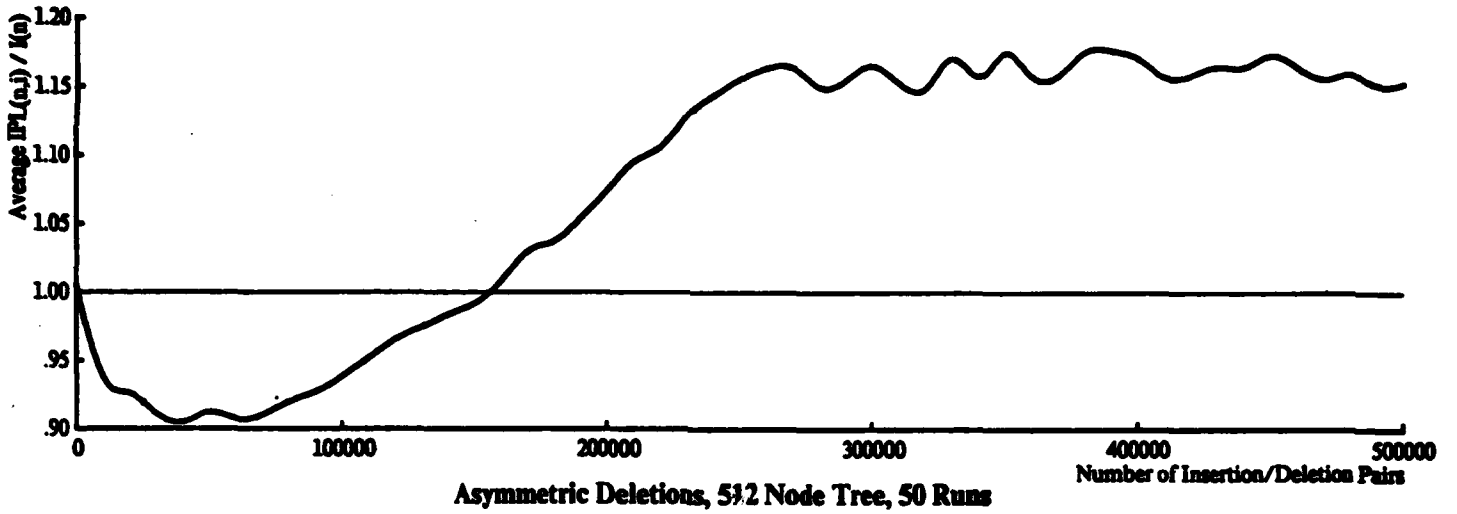


Figure 7

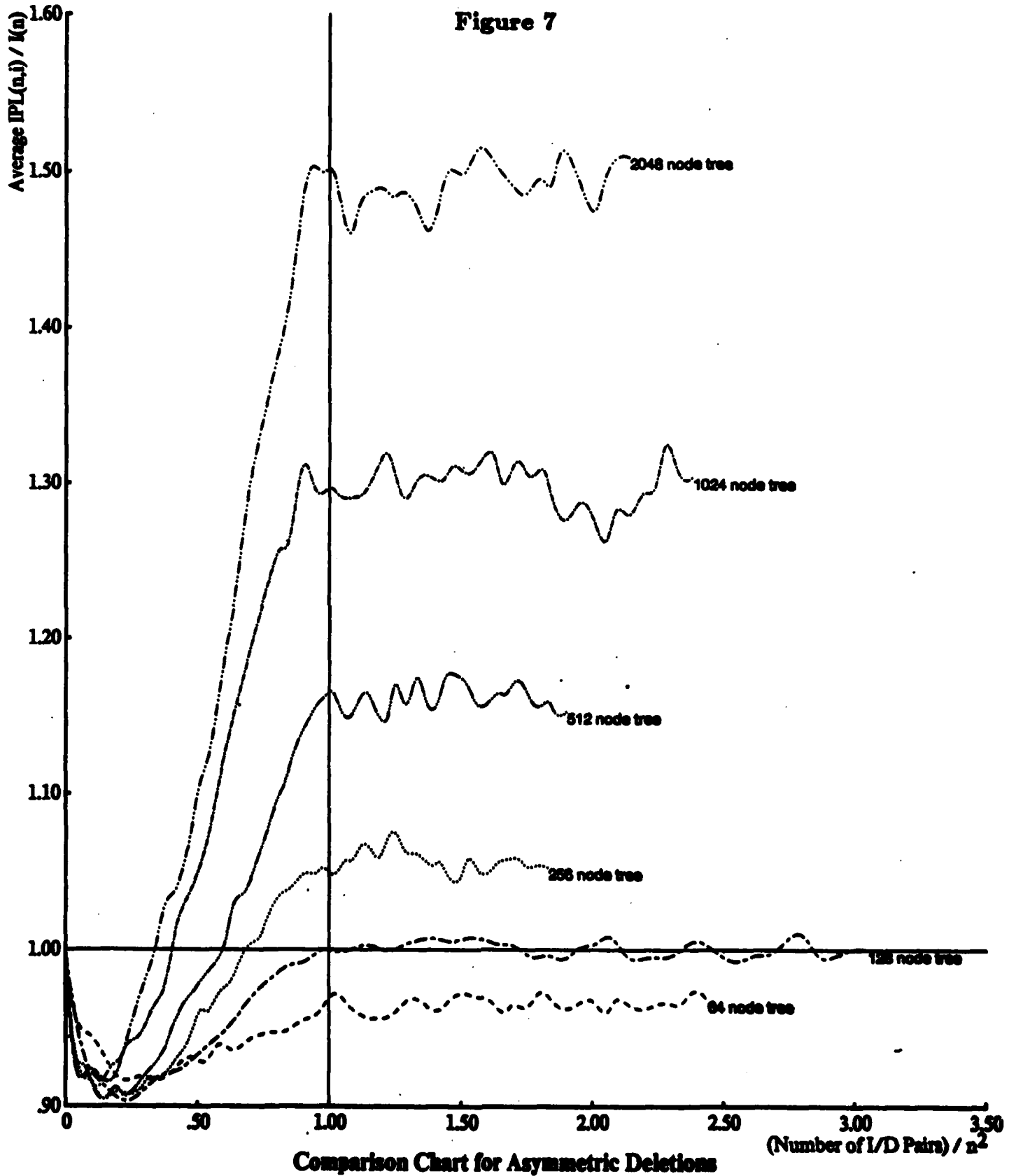


Figure 8

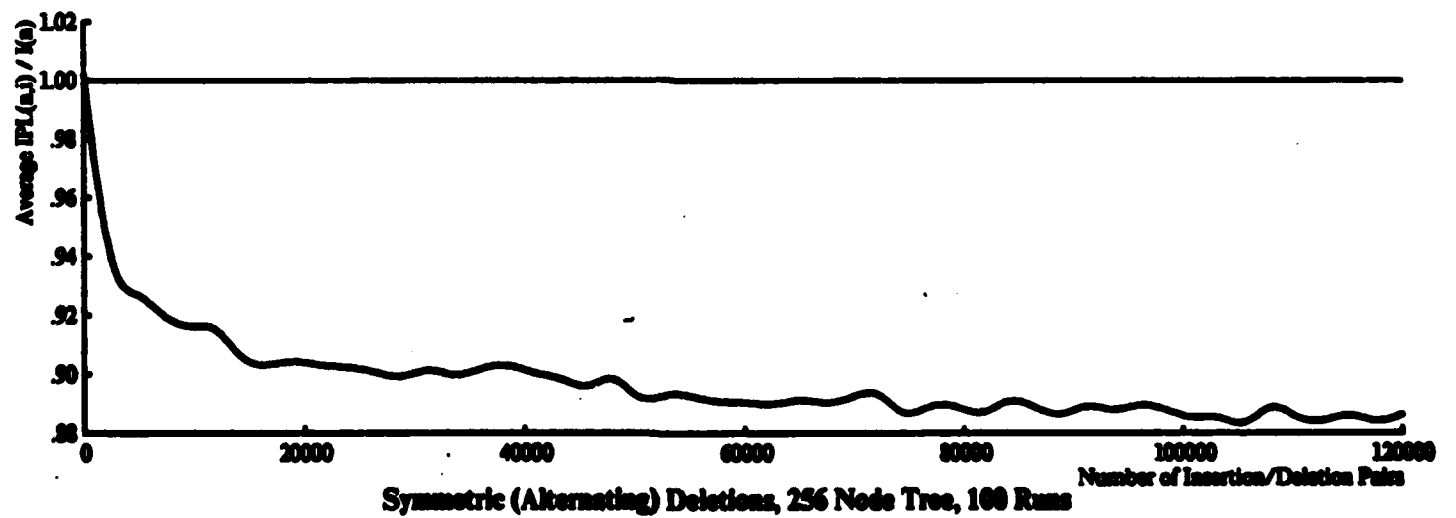
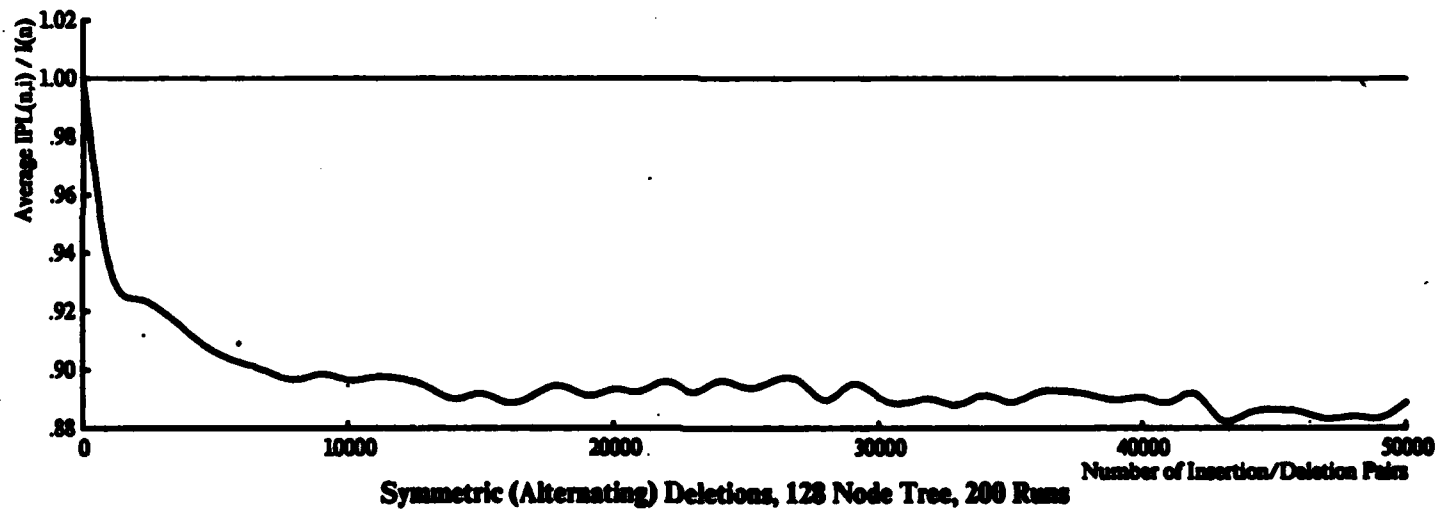
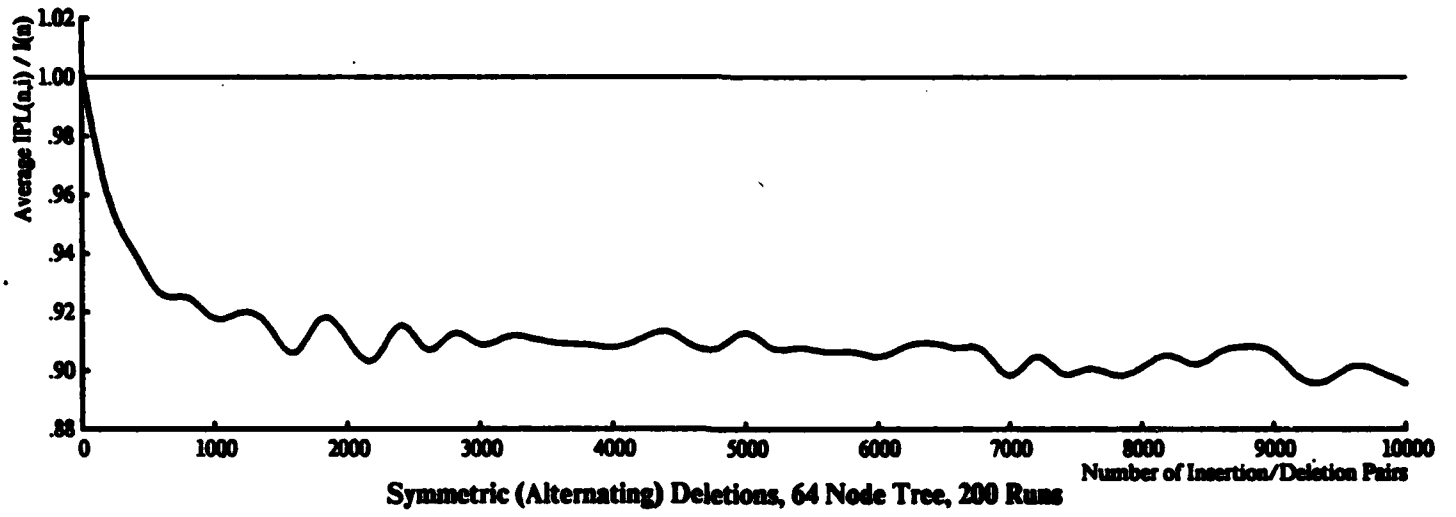


Figure 9

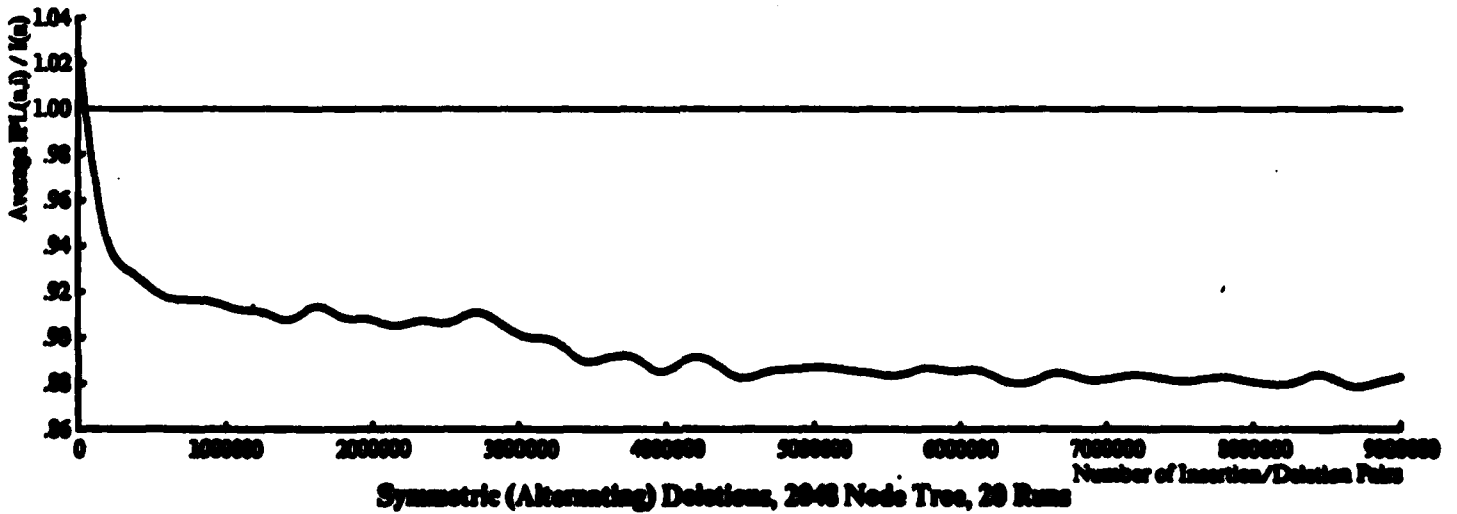
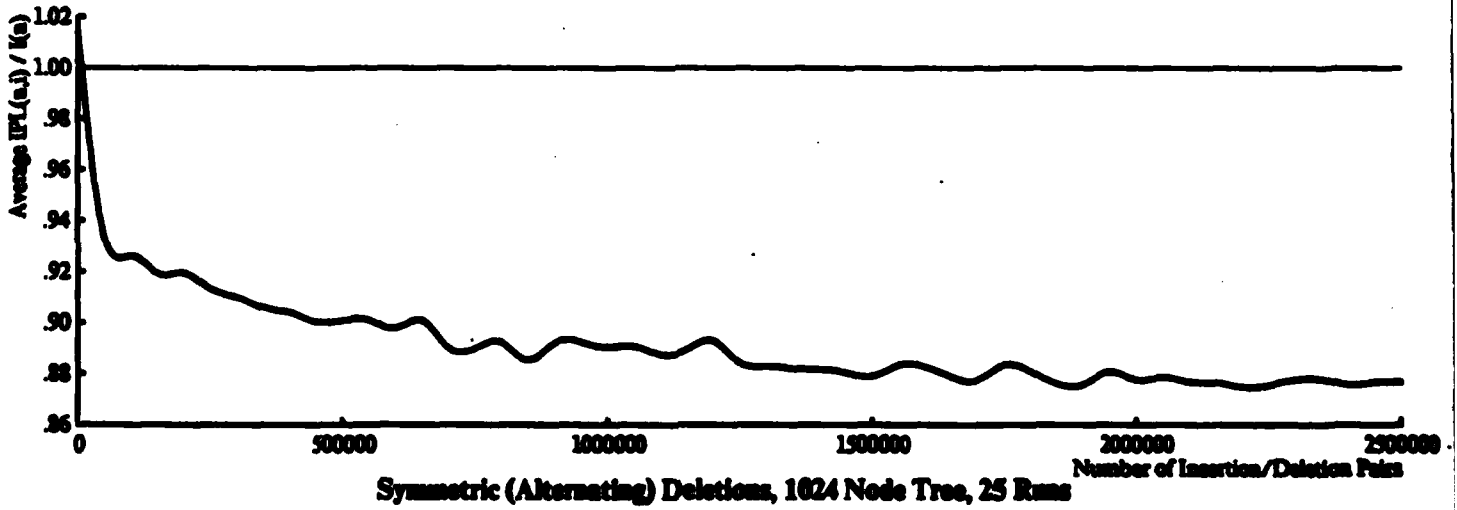
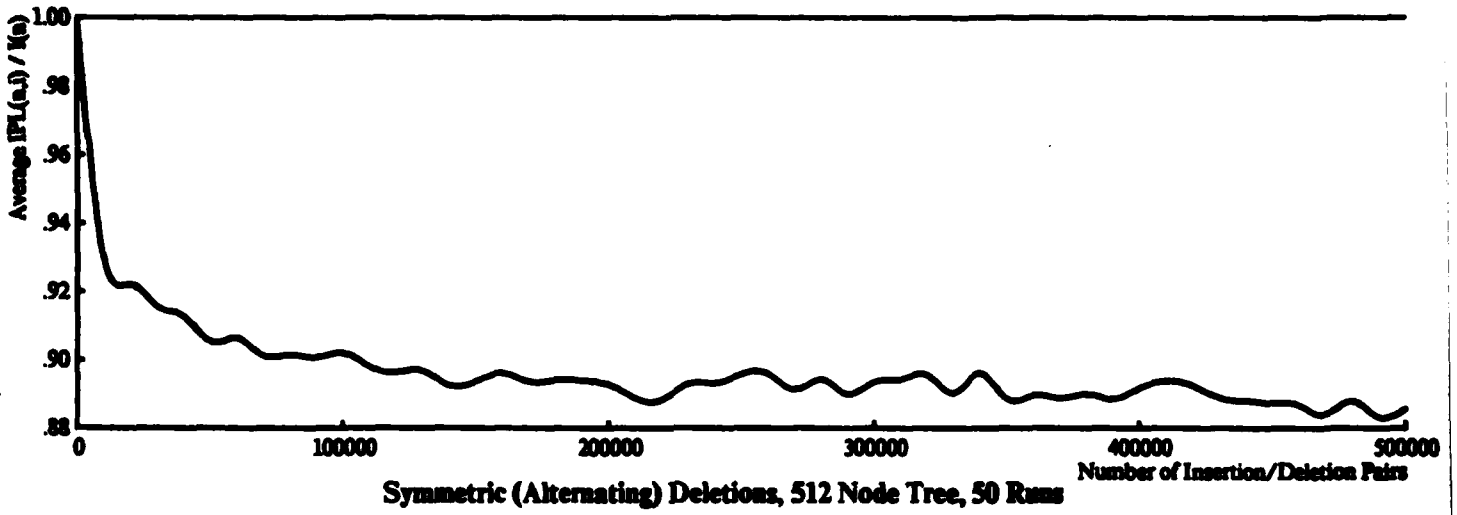


Figure 10

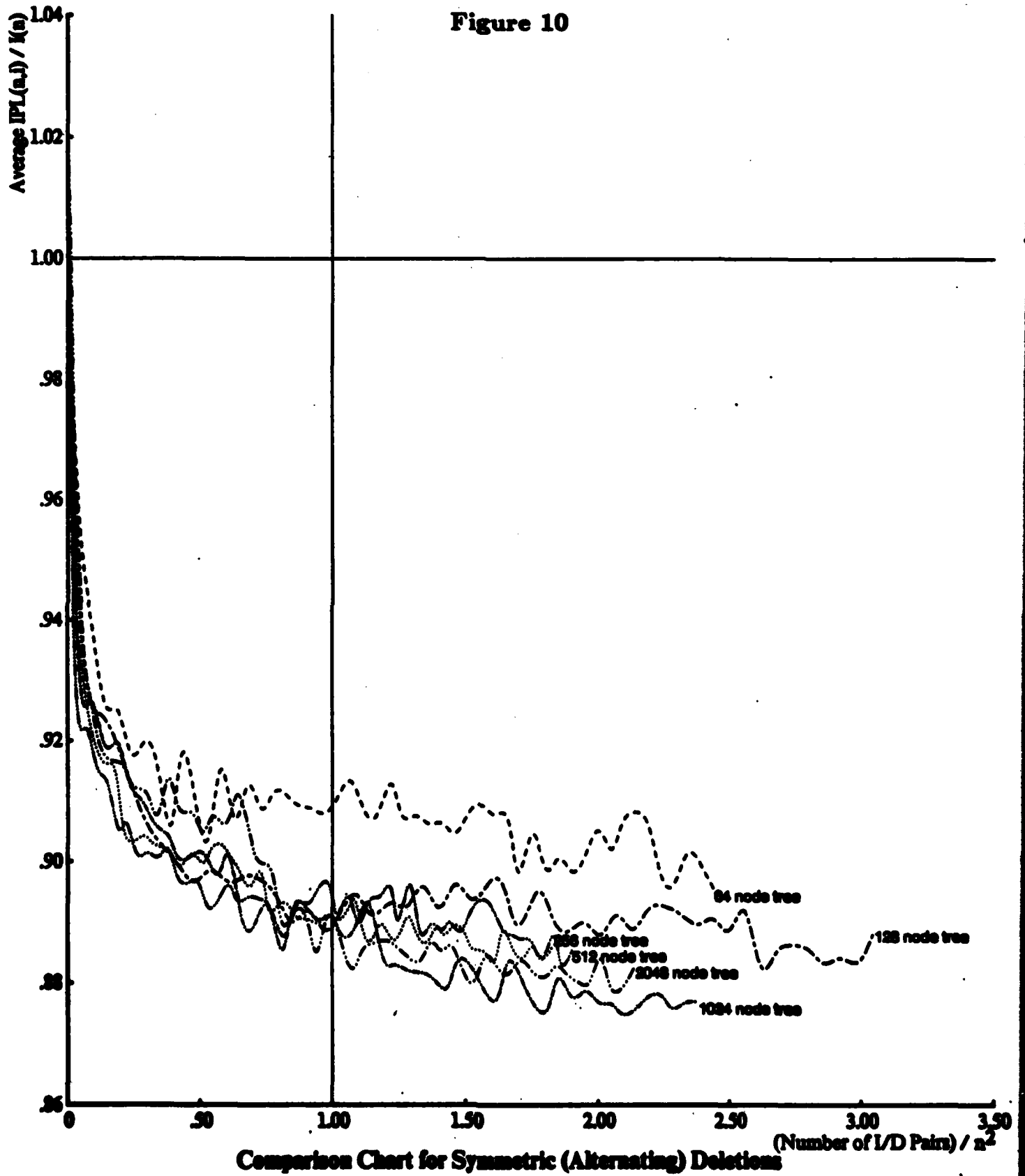
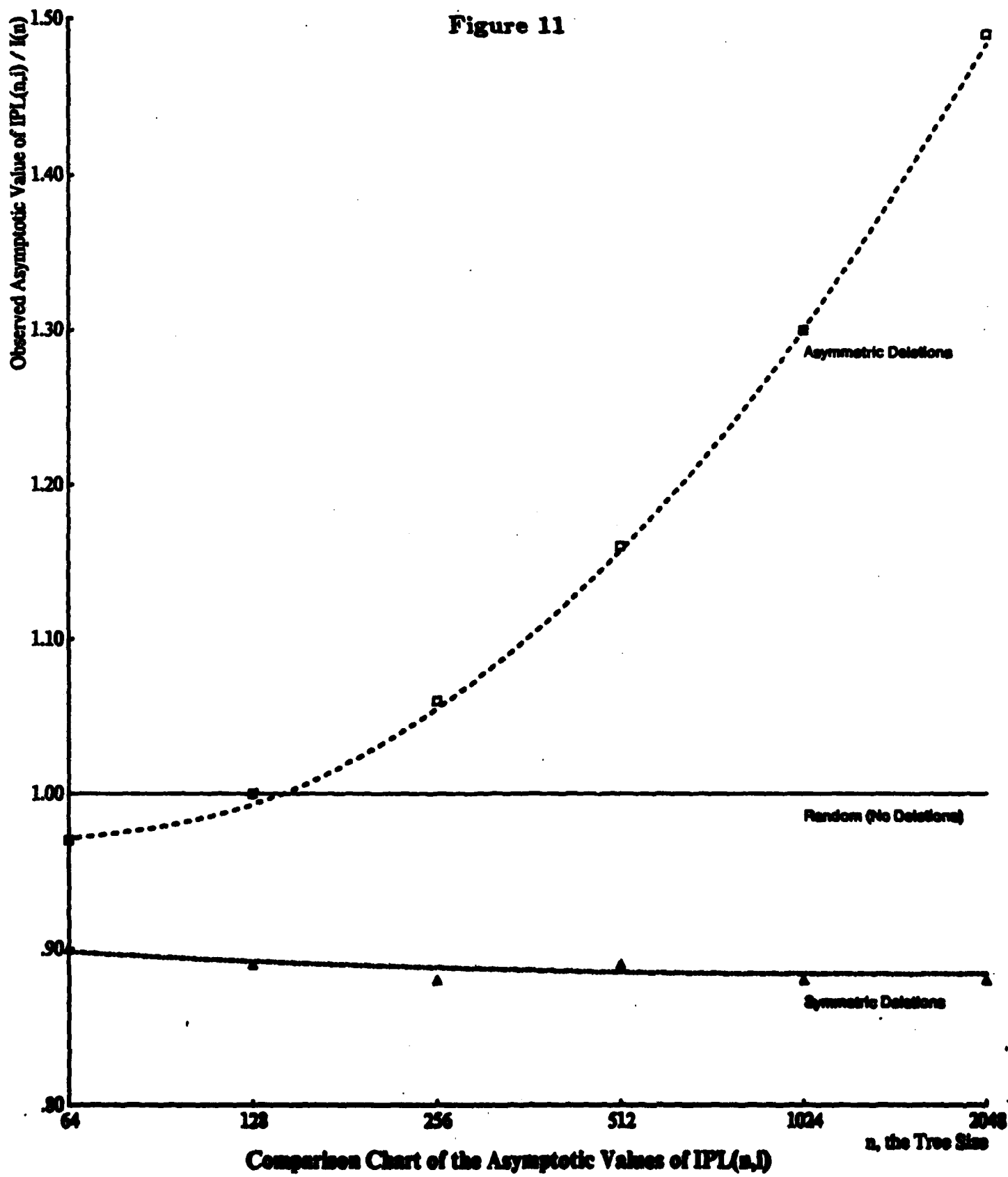


Figure 11



REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER CMU-CS-82-146	2. GOVT ACCESSION NO. A125210	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) AN EMPIRICAL STUDY OF INSERTION AND DELETIO IN BINARY SEARCH TREES	5. TYPE OF REPORT & PERIOD COVERED Interim	
	6. PERFORMING ORG. REPORT NUMBER	
7. AUTHOR(s) Jeffrey L. Eppinger	8. CONTRACT OR GRANT NUMBER(s) N00014-76-C-0370	
9. PERFORMING ORGANIZATION NAME AND ADDRESS Carnegie-Mellon University Computer Science Department Pittsburgh, PA. 15213	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS	
11. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research Arlington, VA 22217	12. REPORT DATE December 2, 1982	
	13. NUMBER OF PAGES 19	
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)	15. SECURITY CLASS. (of this report) UNCLASSIFIED	
	16. DECLASSIFICATION/DOWNGRADING SCHEDULE	
17. DISTRIBUTION STATEMENT (of this Report) APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report) Approved for public release; distribution unlimited		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number)		