

Rigorous Design Documentation of Knuth's Solution to the Common Words Problem

Prabhaker Mateti
 Department of Computer Science and Engineering
 Wright State University
 Dayton, Ohio 45435
 pmateti@wright.edu (937) 775 5114

CONTENTS

I	Introduction	1	VIII	Tries	7
I-A	Goals of this Paper	1	VIII-A	The trie	7
I-B	On Design	1	VIII-A1	Nil ID and Root ID	7
I-B1	Vertical Refinement	1	VIII-A2	The Children are Ordered	7
I-B2	Horizontal Decomposition	1	VIII-A3	Parent is Unique	7
I-B3	Rigorous Description of Software Designs	1	VIII-A4	Search for the Word in a Trie	8
I-C	Structure of the Paper	2	VIII-B	Tries with Rings	8
II	The Common Words Problem: Spec	2	VIII-B1	Insert New Word	8
II-A	What is Text and What are Words?	2	VIII-B2	Word from Trie	8
II-B	Occurrences of a Word	2	VIII-C	Design D5 Using Ringed Tries	8
II-C	Alphabetically Sorted	3	VIII-D	Frequency Sorting of the Words	8
II-D	Word-Nat Pairs	3	VIII-D1	Find Frequent Words	8
II-E	k Most Frequent Words	3	VIII-D2	Preparing to Dispense with FOQ	9
II-F	Input File of Text	3	VIII-E	Design D6 Using Ringed Tries but Without FOQ	9
II-G	Output File of Text	3	IX	Hash Tries	9
III	CWP Designs	3	IX-A	Cell-Ids Refined	9
III-A	Design D0 with a Generic Container of Words	3	IX-B	Reading a Word from the Hash Trie	10
III-A1	Module cow	3	IX-C	Initial Hash Trie	10
III-B	Mapping Input Text to Words	4	IX-D	Search of a Word	10
III-B1	nextlexeme: A Spec	4	IX-E	Insertion of a Word	10
III-B2	nextlexeme: A Design	4	IX-F	Sorting the Words by Frequency	11
III-B3	Module lex	4	IX-G	Design D7 Using a Hash Trie	11
III-C	Design Levels	4	X	Related Work	11
IV	Bag of Words	4	X-A	Open Source Projects	11
IV-A	Design D1 Using a Bag of Words	4	X-B	Rigorous Description of Software Designs	11
IV-B	Design D1 is a model of D0	5	XI	Evaluation	12
V	Dictionaries	5	XI-A	Critique of Knuth's Solution	12
V-A	Tables v. Dictionaries	5	XI-B	Why is it this long?	12
V-B	Design D2 Using Tables	5	XI-C	ÔM, Our Design Language	12
V-C	Mapping a Bag of Words to a Dictionary	5	XII	Conclusion	12
VI	Word-Nat Pairs Sorted on Spelling	5			
VI-A	Design D3 Using Word-Nat Pairs	5			
VII	N-ary Trees	6			
VII-A	Search for a Word	6			
VII-B	Design D4 Using N-ary Trees	6			
VII-C	N-ary Tree to Bag of Words	7			

LIST OF FIGURES

1	Example n -ary tree	6
2	Example Plain Trie	7
3	Example Trie with Rings	8
4	Example hash trie [Bentley 86](p 479)	10

[Draft Revised 2018/09/29. It even includes notes to self!]

Abstract—We document the design of Knuth’s solution to the Common Words Problem (CWP) as an example of how the designs of intricate programs might be presented using rigorous justification. The CWP and Knuth’s solution use data structures known as dictionaries, and hash tries, and notions such as lexical structure. These have been the main source of ambiguity. We give precise definitions for all these in a design specification language called ÔM. We explicitly define all our objects and also exhibit the design hierarchy that we were able to reverse engineer from his solution.

I. INTRODUCTION

We have [Knuth 84] suggesting that program development and documentation be treated as ‘literate programming’, and has contributed several examples, both large (T_EX[Knuth 86a], Metafont [Knuth 86b]) and small ([Knuth 84] and [Bentley 86]). The tradition pioneered by [Kernighan and Plauger 76] is continued in such books as [Comer 84], [Tanenbaum 87], [Wirth and Gutknecht 92] and [Fraser and Hanson 95], which include complete listings of source code, along with cogent explanations of why they work. These are impressive accomplishments. But all the above examples neglect to emphasize design descriptions, and concentrate on implementation details.

There is a fundamental difference between Knuth’s examples, and the books by others mentioned above. Knuth’s examples are meant for peers to read, understand and evaluate, whereas the above books are for students to emulate and learn the programming techniques. The designs of Knuth’s examples get buried in a myriad of implementation surface details. The descriptions in the other books are too imprecise.

There are now a large number of large open source projects, with thousands of pages of documentation, but with hardly any design descriptions. [Lenberg et al. 2015] [Salman et al. 2015] [FM 2018]

A. Goals of this Paper

CWP was solved by that grand master Knuth himself as an example of literate programming and as an example usage of WEB. Knuth’s solution had served its purpose. In fact, its reviewer McIlroy found “... Knuth’s program convincing as a demonstration of WEB and fascinating for its data structure.” The complexity of that fascinating data structure, called the hash trie, was a major obstacle in its understandability. On the other hand, thanks to this complexity, Knuth’s program provided enough material for the WEB system in describing “real” and complex systems. An alternate solution to the CWP was given by Hanson [Van Wyck 87].¹

Our goal is not to do yet another literate programming example, but to focus on “design.” What should design descriptions of software contain? How should they be organized? These two questions are implicitly answered, for the CWP problem, by the material of this paper and its organization. We convey our concerns for the precise expression of software designs by reworking the CWP.

This paper deals with the *expression of the design ideas already present* in Knuth’s solution to the CWP. It does not take

sides on the question of whether hash tries are too complex a solution to be used for simple problem such as CWP. If hash tries are indeed appropriate in a situation, and a designer chose to use them, how should the concept of hash tries and its incorporation into the solution be presented is what this paper deals with. Our goal is to explain the design with precision.

B. On Design

A design document captures the end result of having designed something. The decomposition of the given problem into subproblems, and “judging” the effectiveness of the decomposition, the previous approaches to similar problems that the designer considered — all these must be recorded. In the long run, other designers can benefit from such clarification and cataloguing of design principles, and greater understanding of the structure of design descriptions.

A design is a plan of realizing an object from given primitives, materials and environment while meeting its specifications. A software design is a “program” written in an appropriately chosen language. A software specification and design needs to be captured in machine processible form. Since design documents are meant for humans also, it is equally important to include natural language descriptions of the formal objects manipulated in the design.

1) *Vertical Refinement*: At the highest level, a software design can be (ought to be?) state-free and highly declarative. Intermediate level designs will look similar to programs in a procedural language, but manipulate abstract data types such as sets, bags, and tables. At the most detailed level, our designs are “only a step away” from a modern programming language.

The levels are based on *abstraction*. As to when something is abstract or concrete is ultimately a subjective matter, but it is based on the architecture of the target computer and programming language. Also, abstraction is not an either-you-have-it-or-you-don’t item. We call the process of adding detail to an abstract thing to make it more concrete *vertical refinement*. Reification is an alternate term that [VDM] uses.

2) *Horizontal Decomposition*: We call the splitting of the given problem into subproblems, and grouping the internals of each subproblem *horizontal decomposition* giving us modules. A design module describes

- compositions of parts
- functionality of parts
- data structures + their properties
- algorithms + their properties
- reasons for design decisions made
- perceived “good” (possibly “unprovable”) properties

3) *Rigorous Description of Software Designs*: The problem and Knuth’s solution use structures that are known as dictionaries, and hash tries, and notions such as lexical structure. These have been the main source of ambiguity. We give precise definitions for all these. Natural languages, and so-called pseudo-code are unsuitable.

It is important to note that our definitions of these structures are not only precise, and rigorous but also *formal*. Most definitions in combinatorics and algorithms books, including those of Knuth, are *informal*.

¹lecture I watched on video has a good summary.

We write these definitions in a software design language, named $\hat{\text{OM}}$, which is as “formal” as a typical programming language is these days. We minimally describe the notation and semantics of $\hat{\text{OM}}$ as needed. Full details of $\hat{\text{OM}}$ are given in [Mateti 2018].

C. Structure of the Paper

This paper is organized into three parts. Section 2 clarifies the description of CWP by providing a set of precise definitions, which remove the ambiguities from the original informal specification of the structure of the input. Section 2 covers all that is needed to precisely describe the lexical structure of the input.

Section 3 of the paper presents background “design knowledge” appropriate for the description of Knuth’s solution to CWP. It defines objects such as dictionaries, tries, hash tries, and some function representations used in the solution of CWP.

Sections 4 – 8 are the third part of the paper dealing with the actual design specification. It presents a hierarchy of designs, each described at a different level of abstraction. It also describes the interrelationship between two successive designs in the hierarchy. Also, along with the increasing level of detail, more concrete representation of the abstract objects used in the early design are progressively introduced. Thus, the first design in the hierarchy uses a table as the representation of the dictionary, describes input/output relationships, and assumes the natural number k to be some value. The second design also uses this representation of the dictionary; however, it provides a more more operational view of the designs by explicitly expressing how one obtains the output from the input in a more constructive fashion.

Finally, there is an overall evaluation of the design example.

II. THE COMMON WORDS PROBLEM: SPEC

The June 1986 Programming Pearls column [Bentley 86] posed the following problem:

“Given a text file and an integer k , print the k most common words occurring in the file (and the number of their occurrences) in decreasing frequency.”

We refer to this statement as *the Common Words Problem*(CWP), and the desired program as `cwp`.

As posed, the CWP problem is imprecisely stated. For instance, it is not clear what sequences of characters constitute valid words, whether a mere difference in character cases makes two words to be considered different, how the integer k is to be input, what to do if the input file contains less than k words, more than k words occur all with the same frequency, or if k is negative. Nothing is said about the format of the output. Many of these points have been raised by David Hanson (as reported in [Van Wyck 87]). We make the problem statement more precise by resolving these issues in the next section.

Some would point out that a “precise” specification is not possible to express in any natural language. True as this may be, we quickly understand, roughly perhaps, the details

when explained in a natural language using well-accepted technical terminology, and defining our new terms. It is worth remembering the points made by [Meyer 85].

A. What is Text and What are Words?

A *text* is a sequence of characters. A line is a non-empty sequence of characters such that the last character is a newline, and no other character in it is a newline. A *text file* is an association of a file-name with a file-content. The file-content must be text. A *word* is a non-empty sequence² of letters, upper or lower case.³

```
value upletter: set char := {'A'..'Z'};
value loletter: set char := ['a'..'z']; // English

assert upletter * loletter == {};
value letters := upletter + loletter;
value delimiters := char - letters;

type text is seq char;
type word is (seq letters) - { [] };
```

Call the solution program by the name `cwp`. Via a command shell, it is invoked as in

```
cwp k input-fnm output-fnm
```

The number k is given as an explicit argument, not as part of the input file being scanned, to the program. The k above stands for the radix-10 word of the number k . We do not care to deal with `stdin`, `stdout`, nor with input output redirections. The file-name must be a word. The arguments `input-fnm` and `output-fnm` are file-names.

B. Occurrences of a Word

Uppercase and lowercase letters are considered equal. Two words are *equal* if they are string-equal after ignoring the lower/upper case distinction. For example the two words `ugLy` and `UgLy` are equal. The `lettercode` function gives the ordinal position of a letter in the alphabet.⁴

```
function infix "==" (s, t: word) is
  (lettercode@(s) == lettercode@(t))
```

What does it mean to say that a word w occurs n times in the text t ? For example, how many times does “hi” occur in “hihihoho hi”? It should be 1.

```
function noccur (w: word, t: text) is
  if #t = 0 => 0;
  :: w == t => 1;
  :: else => noccur(w, x) + noccur(w, y) where
    x, y: seq char, d: delimiters such-that
      (t == x + [d] + y)
  fi
```

²We do not wish to say something like “A word is a maximal subsequence of one or more contiguous certain-kind-of-characters of the input file” because we must be able to talk about a word independently of a text file.

³ $\hat{\text{OM}}$: Sets and sequences are built-in data types. If q is a sequence, $\#q$ is the length of q , the indices begin at 1, and $q[i]$ gives the i -th item of q . The $+$ between sequences denotes catenation; $[]$ denotes the empty sequence.

⁴Explain at @, apply-all, symbol.

C. Alphabetically Sorted

```
function infix "<" (s, t: word) is (
  for some z: word (t == s + z)
  or
  for some i: 1 .. min(#s, #t) (
    s[1..i-1] == t[1..i-1],
    lettercode(s[i]) < lettercode(t[i])
  ));
```

```
function sorted(q: seq T) is (
  for i: 1 .. #q-1 (q[i] <= q[i+1]));
```

```
function rsorted(q: seq T) is (
  for i: 1 .. #q-1 (q[i] >= q[i+1]))
```

D. Word-Nat Pairs

```
type wordnat is tuple (w: word, n: nat);
```

Consider wnq , a sequence of k word-number pairs, with the following properties. We can think of it as a table of two columns. Its first “column”⁵ lists words without duplication. Its second column (a sequence of numbers) is sorted in the non-increasing order. Further more, a word w is listed in wnq only if no word, not also listed in wnq , occurs more frequently in the input text than does w .

```
function to-wordnat(itxt: text, k: nat) is
  value wnq : seq wordnat such that (
    k = #wnq = #set(wnq.w),
    sorted(wnq.n),
    for w: word (
      w in wnq.w
      implies
      for x: word (
        x in wnq.w
        or noccur(x, itxt) <= noccur(w, itxt))
    ) )
```

TBD Check⁶

E. k Most Frequent Words

How should we define “ k most frequent words”? Suppose $k = 5$, and no word occurs more than 8 times, and there are 10 distinct words occurring exactly 8 times each. Would any 5 words from the 10 be acceptable?

The above definition of wnq yields a subtle way out for the special case when the input text does not contain k distinct words. It allows, in only this case, wnq to contain words that do *not* occur in the input. Obviously, the frequency of such words must be 0.

The design problem is essentially to construct wnq from the given input text $itxt$ and number k .

F. Input File of Text

Our input file content is $text$. For use in our designs, we map the text to a sequence of words.

G. Output File of Text

The output of the program is a text file. The file-content of the output file must consist of exactly k lines. Each line

⁵OM: $wnq.w$ yields a sequence made out of the w -components of the wnq sequence.

⁶ $itxt$ may not have enough words?

consists of exactly two words. The second word is a radix-10 word of a non-negative number. This number is the count of how many times the first word occurred in the input file content. The first words of the k lines must all be distinct.

If the input file content has less than k different words, the cwp is free to chose arbitrary words with zero as their counts.

Each word in the output is printed along with its frequency count on a separate line. The word is separated from the the frequency count by a fixed number, say one, blanks. No word is to be repeated. The number must not have leading zeroes.⁷

```
function maptotext(q: seq wordnat) is
  if q == [] => [];
  :: else =>
    q[1].w + " " + itoa(q[1].n) + "\n" + ↵
    maptotext(q[2..])
  fi
```

Considering the command issued as an invocation $cwp(k, fin, fout)$, we describe the functionality in terms of the pre- and post-conditions of the cwp .

```
file-content(fout) :=
  maptotext(
    to-wordnat(file-content(fin), k)
  );
```

III. CWP DESIGNS

We will be presenting seven levels of design (see Section III-C), each described at a different level of abstraction.

A. Design D0 with a Generic Container of Words

Our main concern at the highest level of design is functionality. The design D0 is given three arguments: k is the number of common words, fin is the name (as one word) of input file, $fout$ is the name of the output file to be produced.

```
module D0(k: nat, fin: word, fout: word) is (
  import module lex, cow;

  init (
    var itxt := file-content(fin);
    lex.init(itxt);
    var cw := cow.init(lex.nextlexeme(itxt));
    cw.build-all-words();
    file-content(fout) := cw.find-frequent-words(k);
  ) )
```

The file named $fout$ will have as its content the string built from the k most frequent words.

1) *Module cow*: The cow module provides a container of words. The procedure $build-all-words()$ inserts into cw all the words found in the content of file fin .

```
module cow is (
  import lex.nextlexeme(nat) (nat, nat);
  init (var cw := {}); // empty set
  let old-word(w) == w in cw;
  let incr-count(w) == /*: w was in cw */;
  let add-new-word(w) == /*: w was not in cw */;
```

```
procedure build-all-words() is (
  var m, n: nat;
  var i: nat := 1;
```

⁷result type of $maptotext$?


```

while do
  (m, n) := nextlexeme(i);
  if m > n => break fi;
  let w == itxt[m..n];
  if old-word(w) => incr-count(w);
  :: else => add-new-word(w)
  fi;
  i := n + 1;
od
))

```

In the above, the comments beginning with `/// are “formal” comments expected to be replaced with code in $\hat{\text{OM}}$ later.`

B. Mapping Input Text to Words

Function `nextlexeme` examines `itxt[i..]`, without modifying it, and establishes the borders of the next word. We wish to construct the `cw` incrementally by adding each word delivered by `nextlexeme`.

1) *nextlexeme: A Spec*: The following is a specification, not a design, of function `nextlexeme`.⁸

```

function nextlexeme(i: nat) is value
  (m: nat, n: nat) such-that (
    ( itxt[m..n] in word,
      i <= m <= n <= #itxt,
      n < #itxt implies itxt[n+1] in delimiters,
      set(itxt[i..m-1]) <= delimiters
    )
  )
  or
  (m > n iff set(itxt[i..]) <= delimiters)
)

```

2) *nextlexeme: A Design*: We now present a design of the above that maps the given sequence of characters in the input file into a sequence of words *as and when needed* is described here. This module is not further refined.⁹

3) *Module lex*: Module `lex` supplies function `nextlexeme` that takes one `nat` argument `itxt`, and returns a pair (a, b) : (nat, nat) . The spelling of the next lexeme starts at index a and ends at b .

```

module lex(itxt: text) is (
  assert ( itxt[#] in delimiters,
    itxt[#-1] !in delimiters );

  procedure nextlexeme(i: nat) pre (i < #itxt) is
    var (m: nat, n: nat) such-that (
      m := i;
      while do itxt[m] in delimiters => m := m + ↵
    1 od;
      n := m;
      while do itxt[n] !in delimiters => n := n + ↵
    1 od;
      n := n-1;
    )
  )
)

```

C. Design Levels

In what follows, the module `cow` / data-structure `cw` will be refined several times. Our first design D1 uses a bag of words as `cw`, which is quickly refined into an ordered set of pairs. Each pair consists of a word, and its frequency count in the bag in order make the operations `old-word(w)`, `incr-count(w)`, `add-new-word(w)` efficient. The set

is ordered alphabetically by the spelling of the words it contains.

The representation is progressively refined from a table (designs D2 and D3) to an n -ary tree (design D4), to a trie (designs D5 and D6), and finally to a hash trie (design D7).

In all the designs, we `build-...` first then we `find-...` This suggests that we may choose one representation for `cw` during the `build-...`, and a different one for the `find-...`, transforming the representation once between the two. During the `find-...`, in order to efficiently discover the most frequent word, it would be best if `cw` were a set of word-count pairs ordered not alphabetically but by the frequency counts. Knuth does this by progressively converting, in situ (i.e., without using additional memory), the (hash) trie into a linked list. This is perceived by many readers as tricky and presents us with one more layer of abstraction, from design D5 to D6.

IV. BAG OF WORDS

Even though the specification (Section II) is considering *all* the words when it says `for x: word ...`, in the design we need consider only the words occurring in the input text, `itxt`. On the other hand, we must examine each and every word of `itxt`, otherwise we may miss the most frequent word, or have wrong frequency counts.

It is our goal to generate a/the piece of text that satisfies the output requirement of the specification progressively as the value of the variable `otx`, which stands for `file-content(fout)`.

A. Design D1 Using a Bag of Words

The following design satisfies the CWP. Functions `nextlexeme` is specified in Section III-B1.

```

module bow is
  init (
    var itxt := file-content(fin);
    var bw: bag of word := {};
  )

  let old-word(w) == ();
  let incr-count(w) == ();
  let add-new-word(w) == (bw += {| w |});

  procedure build-all-words() is
    as-in module D0;
  )

```

The empty parens in `old-word(w) == ()` and in `incr-count(w) == ()` need explanation. Maybe `old-word` should be `false`, and `incr-count` should be a `no-op`?

Function `mostfrequent` examines the bag `bw`, and returns a most frequent word and its frequency. Note that we deliberately choose not to uniquely specify which word is to be returned when there are several equally frequent ones in `bw`.

```

function mostfrequent() pre (#bw > 0) is
  value wn: wordnat such-that (
    wn.n = wn.w #in bw,
    for x: bw (x #in bw <= wn.n)
  );

procedure find-frequent-words(k: nat) is

```

⁸ $\hat{\text{OM}}$: In the context of sets, the token `<=` stands for the subset-of relation.

⁹Where did we make sure that `itxt[last]` is a delimiter?

```

var otx: text (
var w: word;
var i, n: nat;
otx := [];
for i: {1..k} (
  if bw == {| |} => break fi;
  (w, n) := mostfrequent();
  bw := bw - {| w ** n |};
  otx := otx + w + [blank] + itoa(n) + [newline]
) )

```

B. Design D1 is a model of D0

The design D1 refines D0. It is essentially the same as D0. Written formally as $D1 \models D0$.

```

module D1(k: nat, fin: word, fout: word) is (
import module lex, bow;
init (
  var itxt := file-content(fin);
  lex.init(itxt);
  var bw := bow.init(lex.nextlexeme);
  bw.build-all-words();
  file-content(fout) := bw.find-frequent-words(k);
) )

```

V. DICTIONARIES

Conceptually, a dictionary is a collection of objects organized in a particular way to ease subsequent search of these objects. Each object in such a collection is attached various attributes of interest. For our purposes here, the only attributes of interest are its spelling and its frequency of a word.

A. Tables v. Dictionaries

A few subtleties aside, the *tables* of $\hat{O}M$ can model the dictionaries nicely. A *table* is a set of like tuples whose first elements are all distinct. A tuple is similar to a sequence but may contain dissimilar items.

```

module dict is (

init ( var dwn: table wordnat := {} );

let old-word(w) == w in dwn.w;
let incr-count(w) == dwn[w].n += 1;
let add-new-word(w) == dwn += { (w, 1) };

procedure build-all-words() is as-in module D1;

procedure find-frequent-words(k: nat)
pre (#dwn >= k) is var otx: text (
  value count-wng: seq wordnat such that (
    bag(count-wng) = bag(dwn),
    rsorted(count-wng.n)
  );
  assert (#count-wng >= k);
  otx := [];
  for i: {1..k} (
    let (w, n) == count-wng[i];
    otx += w + [blank] + itoa(n) + [newline]
  ) ) )

```

The `count-wng` is specified by describing its properties. We never “do a design” for it because in the later refinement of the overall design D2, the `count-wng` disappears.

B. Design D2 Using Tables

D2 looks identical to D1 (Section IV-A) except we now use `dwn` instead of `bow`.

```

module D2(k: nat, fin: word, fout: word) is (
import module lex, dict;
init (
  var itxt := file-content(fin);
  lex.init(itxt);
  var dwn := dict.init(lex.nextlexeme);
  dwn.build-all-words();
  file-content(fout) := ↔
  dwn.find-frequent-words(k);
) )

```

C. Mapping a Bag of Words to a Dictionary

```

function dictionary(bw: bag of word) is
value d: table wordnat such-that (
  for w: set(bw) (d[w].n == w #in bw),
  for w: d.w (d[w].n == w #in bw)
)

```

The above defines a relationship between `d` and `bw`. Clearly, we want all the words in the bag `bw` appear in the first column of the table `d` with the correct count: `for w: set(bw) (d[w].n == w #in bw)`. The second line is requiring that whatever words are in the first column of the table, their occurrences count in the bag be correct. The second line could have been written equivalently as

```
for w: d.1 - set(bw) ( d[w].2 == 0 ).
```

In other words, we are allowing for the possibility of non-`bw` words to appear in the table. This happens to be a *significant and insightful* jump in the design process.

As can be readily seen, `dictionary({| |}) = {}`. Suppose `dwn = dictionary(bw)`. Then after `add-new-word(w)` and after the `if`-statement in `build-...` we will have the same relationship holding with the updated values for `dwn` and `bw`.

VI. WORD-NAT PAIRS SORTED ON SPELLING

Let us consider a design where we have the dictionary continually sorted for ease of searching for a word. During the building up of this “dictionary” it will be maintained as a sequence of tuples, `var alpha-wng: seq wordnat`, sorted based on the alphabetic order of words.

```

function alphasorted(q: seq wordnat) :=
  (for i: 2..#q (q[i-1].w <= q[i].w));

```

A. Design D3 Using Word-Nat Pairs

Design D3 refines D2 by using module `alpha-sorted-dict` instead of `dict`. D3 does not refine the procedure `find-frequent-words` of D2 further.

```

module alpha-sorted-dict is (
init ( var alpha-wng : seq wordnat := [] );

let old-word(w) == (w #in alpha-wng.w > 0);
let incr-count(w) ==
  (alpha-wng[i].n := 1 + alpha-wng[i].n);
let add-new-word(w: word) ==
  alpha-wng := value uqwn: wordnat such-that (
    set(uqwn) = set(alpha-wng) + { <w, 1 > },

```

```

    alphasorted(ugwn)
  );

  procedure build-all-words() is as-in module D2;
  post set(alpha-wnq) = dwn;

  procedure find-frequent-words(k: nat) is
  as-in module D2;
)

```

Note the post-condition `set(alpha-wnq) == dwn`. The dictionary `dwn` was allowed to contain certain words with zero counts; hence, `alpha-wnq` also will. But neither `dwn` nor `alpha-wnq` have indicated specifically what the characterization of these zero-count words are.

VII. N-ARY TREES

An n -ary tree is a rooted tree, where each node has at most n ordered subtrees; see Figure 1. In CWP, we store in each node a letter, and a count. The path from the root to a node yields a word made up of these letters. The `cnt` field of a node contains the number of times the word represented by the path to this node occurs in the input text `itxt`. The `cnt` fields of some nodes may be zero since not every prefix of a word occurs as a word in `itxt`.

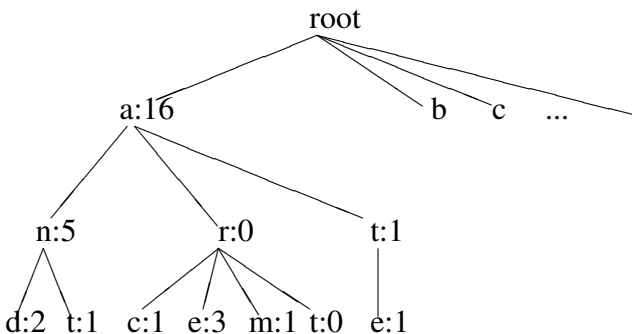


Fig. 1. Example n -ary tree

We also would like to alphabetically order the subtrees of every node based on the letters in their roots.

```

type ntree-ao is
  tuple (
    ltr: letter,
    cnt: nat,
    stq: seq ntree-ao
  ) such that ( for all t: ntree-ao
    (sorted(t.stq.ltr)) )

```

The tree of Figure 1 is the result several operations: `add-new-word` with `a`, `an`, `and`, `ant`, `arc`, `are`, `arm`, `at`, `ate`, and `incr-count` with `a` 15 times, with `an` 4 times, with `and` once, and so on. Note the node with `t:0` shown as a leaf in this figure. Using the operations of `add-new-word` and `incr-count` it is *not* possible to have a *leaf* node with 0 count; it is an internal node and the subtree below that node (e.g., due to an `add-new-word("article")` is not shown. However, the tree of Figure 1 is, as-is, a valid n -ary tree.

As an example insertion of a new word, let us consider `arduous`. After `old-word("arduous")` says false, we make a subtree by invoking `mk-ntree-ao("duous")`,

which must become a subtree of node `(r:0)` between the subtrees at `c:1` and `e:3`.

A. Search for a Word

We search for a given word `w` as follows in a copy of the tree `this`. The search routine does not alter the tree in any way.

```

procedure search(w: word) is (
  var t: ntree-ao := this, p: nat := 0;
  while do
    var i: nat := w[p+1] ## t.stq.letter;
    if i == 0 => break fi;
    p += 1;
    if p == #w + 1 => break fi;
    t := t.stq[i]
  od
)

```

The procedure returns a triplet value `(t, p, i)`. Its first component `t` is an n -ary tree rooted at the last node searched. The second component `p` is a natural number. It is set so that the path from the root yields `w[1..p-1]`, and either `p = #w + 1`, or the letter `w[p]` is not among the children of `t`. Thus,

`p = #w + 1` if the word is already present, and `i` is set so that the root of the i -th subtree contains the letter `w[#]`, `p ≤ #w` and `i = 0` otherwise.

B. Design D4 Using N-ary Trees

`D4` refines `old-word(w) == (w #in alpha-wnq. w > 0)` of the preceding section into `search(w).2 = #w + 1`. The initialization `var nt: ntree-ao := (' ', 0, [])` produces an n -ary tree that has just one node (the root) containing the blank¹⁰, the count zero, and the empty sequence as its `stq`.

```

module D4(k: nat, fin: word, fout: word) is (
  import type ntree-ao, module lex(itxt);
  init (
    var itxt: text := file-content(fin);
    var nt: ntree-ao := (' ', 0, []);
  );
  let (tw, pw, iw) == nt.search(w);
  let sq == tw.stq;
  let old-word(w) == (pw == #w + 1);
  let incr-count(w) == (sq[iw].cnt += 1);
  let add-new-word(w) == (
    let m == min ( {#sq + 1} + // set union
      {j such-that ( 0 < j < #sq + 1,
        w[pw] < sq[j].ltr )});
    sq[ @ m := mk-ntree-ao(w[pw..]) ];
  );

```

```

  procedure build-all-words() is as-in module D3;
  procedure find-frequent-words() is
  as-in module D3;
)

```

The `@` stands for the apply-all, and `q[@ i := e]` stands for an updated sequence `q`, where all its items at indices `i` and above are shifted to one-higher index positions, and the i -th item becomes `e`.

```

procedure mk-ntree-ao(w: word) pre (#w > 0) is
  var t: ntree-ao (
    t := (w[#], 1, []);

```

¹⁰letter?

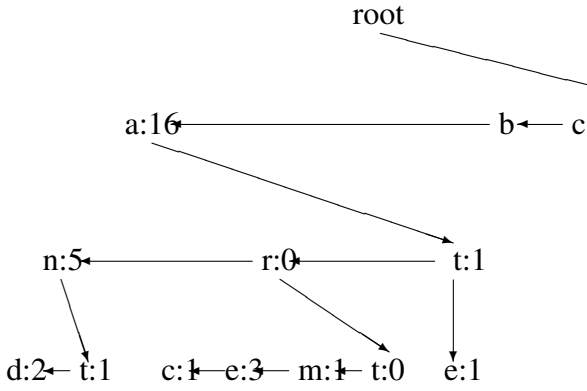


Fig. 2. Example Plain Trie

```

for var i: nat := #w - 1 downto 1 (
  t := (w[i], 0, [t])
)

```

C. N-ary Tree to Bag of Words

For any value of type `ntree-ao` we can find a corresponding value of `alpha-wnq`. The converse, however, is false. This happens due to zero-count words. As a trivial example, the following possible value of `alpha-wnq` is unrepresentable as an `ntree-ao`: `[("are", 1)]` in contrast to `[("a", 0), ("ar", 0), ("are", 1)]`. Recall that neither `dwn` of Section V-C nor `alpha-wnq` of Section VI-A have indicated specifically what the characterization of the zero-count words that they may contain is. For every word `w` with a count greater than 0 occurring in a `ntree-ao`, all the prefixes of `w` will also occur. The counts of these prefixes may or may not be zero depending on whether they have occurred as independent words.

VIII. TRIES

A trie is a binary tree that has a certain relationship to our n -ary tree. This relationship is analogous to that discussed in Section 2.3.2 of [Knuth 97]. Figure 2 is a trie made from that of Figure 1. In terms of the figures, we see that the left most edges of each node are retained, but the other edges from a node to its subtrees are replaced by “horizontal” edges from one node to its sibling.

A. The trie

For our discussion here, `cellid` is any arbitrary type that has a “sufficient” number of values. We define four tables whose keys (i.e., the first components of its elements) are values from this type. A trie is a subset of `cellids`, and the four tables that collectively satisfy certain constraints. These constraints amount to requiring that the structure we are defining better be a binary tree.

```

type cellid;
value emt: iletter := 1 + max(#upletter, #loletter);
value hdr: iletter := 0;

value nilid : cellid := new-cellid();
value rootid: cellid := new-cellid();

```

```

type trie := table (
  cid: cellid,
  ltr: iletter,
  cnt: nat,
  nxt: cellid,
  hic: cellid
) such-that for t: trie (
  is-finite(t.cid),
  the-children-are-ordered(t),
  unique-parent(t),
  req-nilid-rootid(t)
);

let cids == t.cid; // a few abbreviations
let cnt(x) == t[x].cnt;
let nxt(x) == t[x].nxt;
let hic(x) == t[x].hic;

```

Knuth [Bentley 86] used our `hdr` value as his `emt` and vice-versa.

1) *Nil ID and Root ID:* We reserve two values, that we name as `nilid` and `rootid`, from the `cellid`. Every value `t` of type `trie` will be such that `nilid` and `rootid` are in `t.cid`.

```

function req-nilid-rootid(t: trie) is (
  t[nilid].ltr = emt;  t[rootid].ltr = emt;
  t[nilid].nxt = nilid; t[rootid].nxt = nilid;
  t[nilid].hic = nilid; t[rootid].hic = nilid
  t[nilid].cnt = 0;   t[rootid].cnt = 0;
)

```

2) *The Children are Ordered:* The children of a node `u` are ordered based on the letters they contain. The list of children starts with `hic(u)`, and the rest is given by `nxti(hic(u))`. The last child has no next. If a `cellid` `u` has no children, `hic(u) = nilid`. Otherwise, the value `v = hic(u)` is the `cellid` of the child `v` of `u` containing the highest letter among the children of `u`. If `y = nxt(x)` is not `nilid`, we require that `x` be a sibling of `y`, that is `parent(x) = parent(y)`, and `ltr(y) < ltr(x)`. Obviously, `nxti(d) = nilid`, for some $0 \leq i \leq \#letter$.

```

function childrenq(u: cids) is
  value q: seq cids such-that
  if hic(u) = nilid => q = [];
  :: else => (
    q[#] = hic(u);
    nxt(q[1]) = nilid;
    for i: 1..#q - 1 (q[i] = nxt(q[i+1]))
  )
fi

```

```

function parent(u: cids) is
  if u == rootid => nilid;
  :: u == nilid => nilid;
  :: else => value p: cids such-that (u in childrenq(p))
fi

```

```

function the-children-are-ordered() is
  for u: cids (
    let q == childrenq(u);
    for i: 1..#q - 1 (
      ltr(q[i]) < ltr(q[i+1])
    )
  )

```

3) *Parent is Unique:* Note that because `cids` is finite, there exists a finite i such that `parenti(u) = rootid` for all `u` except `nilid`.

```

function unique-parent() is (
  for u, v: cellid (

```

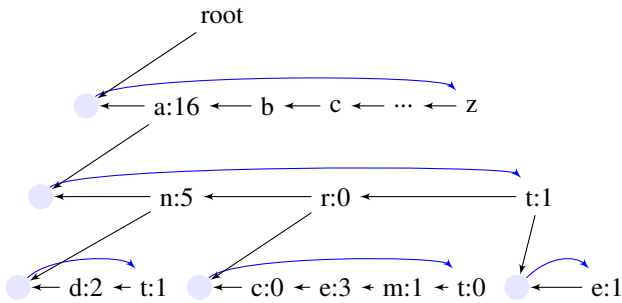



Fig. 3. Example Trie with Rings

```

set(childrenq(u)) * set(childrenq(v))
  != {} iff u == v
))

```

4) *Search for the Word in a Trie:* The procedure `ltr-among-children` is a refinement of `w[p+1] ## t.subs.ltr` of Section VII-A. The procedure is guaranteed to terminate because `ltr(nilid) = hdr < a`, for any `a`. Note that `x` is “behind” `y` in this routine.

```

procedure ltr-among-children(a: letter,
  p: cellid) is (
  var x, y: cellid;
  (x, y) := (p, hic(p));
  while do ltr(y) > a => (x, y) := (y, nxt(y)) od
);

procedure search(w: word) returns
  var (vi: cellid, p: nat, ui: cellid) is (
  vi := t.rootid;
  p := 0;
  while do
  let a == w[p+1];
  (var wi: cellid, ui) := ltr-among-children(a, ↔
  vi);
  if ltr(ui) == a => break fi;
  p += 1;
  assert (ltr(ui) == w[p], ui == nxt(wi));
  if p == #w => break fi;
  vi := ui;
  od
)

```

B. Tries with Rings

We now enhance such tries into those that have circular lists of siblings. See Figure 3. Suppose `u` is a parent whose highest child is `v` and lowest child is `w`. In the unringed tries, `hic(u) = v`, and `nxt(w) = nilid`. In ringed-tries, we introduce an extra cellid `h`, called a header, for each parent with children, so that `hic(u) = h`, `nxt(h) = v`, and the `nxt(v)` etc. remain as they were, except `nxt(w)` which was `nilid` now becomes `h`. The parent of `h` is `u`. The node `h` of course has no children, and there is no meaning yet for either `hic(h)`, or `ltr(h)`. We define `hic(h)` as `u`, and `ltr(h) = dot`. Thus, for all cellids `x`, other than the `rootid`, either `hic(hic(x)) = x` or `hic(x) = nilid`.

```

type ringed-trie is trie except (
  function childrenq(u: cids) is
  value q: seq cids such-that
  if hic(u) = nilid => q = [];
  :: else => (
  q[1] = hic(u);

```

```

  ltr(q[1]) = hdr;
  hic(q[1]) = u;
  nxt(q[1]) = q[#];
  for i: 1..#-1 (q[i] == nxt(q[i+1]))
  fi
)

```

1) *Insert New Word:* As before, `add-new-word(w) == mk-ring-trie(ui, pn+1)`. This adds tuples to the global trie `var t`.

```

procedure mk-ring-trie(ui: cellid, k: nat) is (
  var pi := new-cellid();
  t += {(pi, w[k], 0, nxt(ui), nilid)};
  nxt(ui) := pi;
  for j: nat := k+1 .. #w do
  var hi := new-cellid();
  var ni := new-cellid();
  t += {(hi, hdr, 0, ni, nilid)};
  t += {(ni, w[j], 0, hi, nilid)};
  hic(pi) := hi;
  pi := ni;
  od;
  cnt(pi) := 1;
)

```

2) *Word from Trie:* It can be seen readily that a cellid corresponds to a node of the n -ary tree. The word it represents is quite simple to compute. To get the parent of any node `u`, start from `u`, go to its header `h` by traversing the `nxt` “fields”, and the `hic(h)` yields the parent.

```

procedure header(u: cids) is var u (
  pre (u != rootid, u != nilid);
  while do ltr(u) != hdr => u := nxt(u) od;
)

let parent(x) == hic(header(x));

function word-from-trie(p: cids) :=
  if p == rootid => [];
  :: p == nilid => [];
  :: else =>
  word-from-trie(parent(p))
  + [ltr(p)]
fi

```

C. Design D5 Using Ringed Tries

```

module D5(k: nat, fin: word, fout: word) is (
  import type ringed-trie;

  init (
  var itxt: text := file-content(fin);
  var t: ringed-trie := {};
  );

  import module lex: init (itxt);

  let (vi, pn, ui) == t.search(w);
  let old-word(w) == pn = #w;
  let incr-count(w) == cnt(nxt(ui)) += 1;
  let add-new-word(w) == t.mk-ring-trie(ui, pn+1);

  procedure build-all-words() := as-in module D4;
  procedure find-frequent-words() is /*: below ↔
  */;
)

```

D. Frequency Sorting of the Words

1) *Find Frequent Words:* We now refine the procedure `find-frequent-words` of D2 further. We introduce (temporarily) an extra field `foq` to our trie that will contain a

“linked-list”, frequency-ordered, of all the words with non-zero counts.

```

type trie is table (
  ...
  foq: cellid;    // new/ temporarily
)
such-that ...
);

let foq(x) == t[x].foq;

procedure insert-into-foq(p: cids) is (
  var q: cellid := rootid;
  while do cnt(p) < cnt(foq(q)) => q := foq(q) od;
  foq(p) := foq(q);
  foq(q) := p;
);

procedure rsort-the-trie() is (
  var q: cellid := rootid;
  foq(q) := nilid;
  for p: (cids - {nilid, rootid})
    if cnt(p) > 0 => insert-into-foq(p) fi
);

procedure find-frequent-words() is (
  pre (#dwn >= k) ;
  var q: cellid := foq(rsort-the-trie());
  otx := [];
  for var i: nat in {1..k} (
    let w == word-from-trie(q);
    let n == cnt(q);
    otx += w + [blank] + itoa(n) + [newline]
    q := foq[q];
  )
)

```

2) *Preparing to Dispense with FOQ:* The function `word-from-trie()` depends on `nxt` fields only because `header()` uses `nxt` fields. If the header can be computed in some other manner, the contents of all `nxt` fields is irrelevant to the `for`-loop of `find-frequent-words()`. This issue will be dealt with in Section IX-B.

The procedure `rsort-the-trie` of Section VIII-D1 uses the loop `for p: (cids - nilid, rootid)`. Note that the order in which different values for `p` are chosen is unspecified. So, we exercise our design freedom and start with the alphabetically last word and end after the first.

```

function last-word-from(q: cellid) returns var
  p := q is (
  while do hic(p) != nilid => p := nxt(hic(p)) od
);

let last-word() == last-word-from(rootid);
let end-of-words() == rootid;

procedure next-word(p: cellid) returns p is (
  p := nxt(p);
  p :=
  if ltr(p) = hdr => hic(p);
  :: else => last-word-from(p);
  fi
);

function rsort-the-trie() returns var q: cellid ↔
  is (
  q := rootid;
  foq(q) := nilid;
  p := last-word-from(rootid);
  while do
    if cnt(p) > 0 => insert-into-foq(p) fi;
    p := next-word(p);

```

```

    if p == end-of-words() => break fi;
  od
);

```

As of now, `next-word(p)` does not use the `foq` fields, and `insert-...` does not use `nxt` fields. Hence, while preserving semantics, we can rewrite the above:

```

function rsort-the-trie() returns var q: cellid (
  q := rootid;
  foq(q) := nilid;
  p := last-word-from(rootid);
  while do
    q := next-word(p);
    if cnt(p) > 0 => insert-into-foq(p) fi;
    if q == end-of-words() => break fi;
    p := q;
  od
);

```

This moved `next-word(p)` to a position above the `insert-into-foq(p)`. It can be seen readily that after the `nxt` field of a node `p` are used (either via `last-word-from` or via `next-word`, it will never be needed again except in `word-from-trie()`. Thus, we can set such `nxt` fields to whatever. We let the `foq` values time-share the “residences” of `nxt` fields: that is:

We declare `foq` as an alias for `nxt` field.

E. Design D6 Using Ringed Tries but Without FOQ

[To be inserted]¹¹

IX. HASH TRIES

In defining regular tries, `cellid` was simply a set of cellids whose details were left unspecified. The search and insertion speed is, on the other hand, influenced by what values these cellids are. We now superimpose these considerations on the tries.

A. Cell-Ids Refined

The cellids now become natural numbers with a certain numerical relationship among the parent and children. Suppose the cellids `u` and `v` are siblings. Let `acn` be a function, yet to be discussed, that maps cellids of the preceding section to cell numbers. We will select the mapping `acn` in such a way that the integer

$$acn(v) - acn(u) = ltr(v) - ltr(u).$$

```

type hash-trie is ringed-trie such that (
  value trie-size := ...;
  value alpha := ((sqrt(5) - 1)/2) * trie_size;
  rootid = 0;
  nilid = trie-size + 1;
  cellid = rootid .. nilid;
  for t: hash-trie (
    for u: t.cids (
      let q == childrenq(u),
      for i: 1..#q-1 (
        q[i+1] - q[i] = ltr(q[i+1]) - ltr(q[i]) )

```

Note that we can now satisfy `hic(hic(rootid)) == rootid` by making `hic(rootid)` also be a 0. Note also that the rings of different parents can be interleaved.

¹¹[To be inserted]

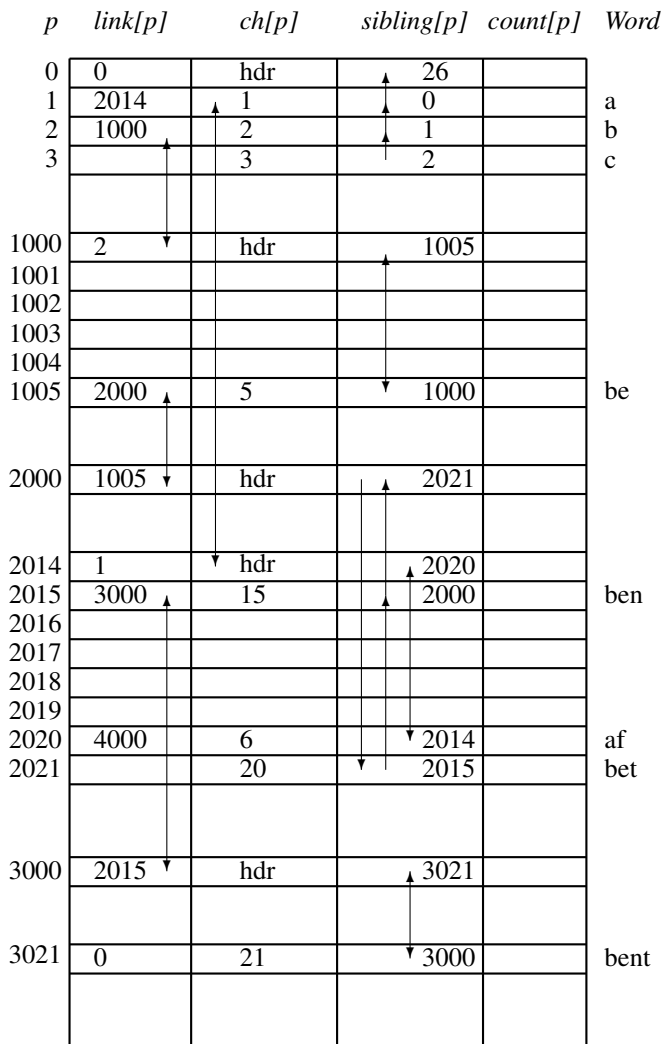


Fig. 4. Example hash trie [Bentley 86](p 479)

B. Reading a Word from the Hash Trie

The header of a node *u* is readily computed, without any traversal.

```
function hash-trie.header(u: cids) is v: cids (
  pre (u != rootid, u != nilid);
  v := u - ltr(u)
```

Reading the word-from-trie is the same as in the preceding section but it now uses the above header().

C. Initial Hash Trie

The following builds an initial hash trie. The *emt* and *hdr* values were defined in Section 7.1.

```
hash-trie.init is (
  var t: hash-trie;

  ltr(0) := hdr;
  cnt(0) := 0;
  hic(0) := #letter;

  for i: {1 .. #letter} (
    ltr(i) := i;
    cnt(i) := 0;
    hic(i) := nilid;
    nxt(i) := i - 1;
```

```
);
  for i: {#letter + 1 .. trie-size} (
    ltr(i) := emt; // empty-slot
  )
```

D. Search of a Word

The procedure search of hash-tries is exactly the same as for tries (see Section VIII-A4). The definitions of *old-word(w)* and *incr-count(w)* remain as before.

E. Insertion of a Word

As we enter the words into the hash trie, the family rings grow larger and larger. As a result, at some point the interleaving of rings does not permit an insertion. To resolve such a collision, we need to relocate one of the rings.

```
procedure hash-trie.relocate-children
  (oh: cellid, nh: cellid) is (
  var r: cellid := oh;
  var d: integer := nh - r;
  while do ltr(r) != emt => (
    nxt(r + d) := nxt(r) + d;
    ltr(r + d) := ltr(r);
    cnt(r + d) := cnt(r);
    hic(r + d) := hic(r);
    if hic(r) != nilid => hic(hic(r)) := r + d fi;
    ltr(r) := emt;
    r := nxt(r);
  )
  od
)
```

As before, *add-new-word(w) == hash-trie.make(...)*. This adds tuples to the global hash trie *t*. The *hash-trie.mk* is more complex than *ring-trie.mk* because we cannot merely choose any new but arbitrary *cellid* for the letters to be inserted.

```
procedure hash-trie.mk(ri, ui: cellid, k: nat) is (
  var pi, hi, ni: cellid;

  oh := hic(ri);
  hi := compute-loc(oh, w[k]);
  if hi != oh => relocate-children(oh, hi) fi;
  pi := hi + w[k];

  ui := hi + ui - oh;
  insert-ltr(w[k], pi, ui);

  for j: nat := k+1 .. #w (
    hi := compute-loc(oh, w[j]);
    ni := hi + w[j];
    insert-ltr(hdr, hi, hi);
    insert-ltr(w[j], ni, hi);
    hic(pi) := hi;
    pi := ni;
  );
  cnt(pi) := 1;
);
```

Function *compute-loc* returns a possibly new location for the header implying that the siblings group needs to be relocated; *nh* equals *oh* if there is no such need.

```
procedure hash-trie.next-hdr-loc(oh: cellid) is
  var nh: cellid
  if oh = last-h => nh := 0;
  :: oh = trie-size - NC => nh := NC + 1;
  :: else => nh := oh + 1
  fi
```

```
procedure hash-trie.compute-loc(oh: cellid, a: ↔
  ltr) is var nh: cellid (
```

```

nh := oh;
if ltr(h + a) /= emt =>
  while do
    nh := next-hdr(hn);
    if will-they-fit(a, oh, nh) => break fi;
  od
fi
)

```

In the function below, a node containing the `a`: `iletter` will become the child of a certain node `p`, whose header is presently at `oh` and we wish to move it to `nh`. The `a` is chronologically the latest child to join the siblings. Function `will-they-fit` is true iff the cells `d` units away from each of the children of `p` are vacant. The distance `d` can be a negative integer.

```

function hash-trie.will-they-fit
(a: iletter, oh, nh: cids) is b: boolean (
  let q == siblings(oh);
  let d == nh - oh;
  pre ((oh, nh) * {rootid, nilid} = {});
  post (b = (ltr(nh + a) = emt,
    for u: q (ltr(u + d) = emt)));
);

```

One letter code is inserted by `insert-ltr()`.

```

procedure hash-trie.insert-ltr
(a: iletter, an: cellid, pn: cellid) is (
  ltr(an) := a;
  cnt(an) := 0;
  hic(an) := nilid;
  nxt(an) := nxt(pn);
  nxt(pn) := an;
)

```

F. Sorting the Words by Frequency

Recall that the sorting by frequency of all the words is begun only after all the input has been read, and the hash trie is fully constructed.

```

procedure hash-trie.link-p-into-sorted(p: cellid) ↔
  is (
    let m == sorted(large-count);
    if f < large-count => insert-into-list(p, f);
    :: cnt(p) >= cnt(m) => insert-into-list(p, ↔
      large-count);
    :: else => insert-into-sib(p, m)
  fi;
);

procedure hash-trie.move-to-last-suffix(var p: ↔
  cellid) is (
  while do chi[p] != nilid => p := nxt[chi[p]] od;
)

procedure hash-trie.rsort() is ( // trie-sort()
  var p, q: cellid;
  for p in 1 .. large-count (sorted[p] := nilid);
  p := nxt[rootid];
  move-to-last-suffix(p);
  while do p != nilid => (
    q := nxt[p];
    if count[p] != 0 => link-p-into-sorted(p) fi;
    if ltr[q] != hdr => move-to-last-suffix(p);
    :: else => p := hic[q]
  fi
) od
)

```

The procedure `find-frequent-words` is the same as before except `fop` is replaced by `nxt`. “After `trie-sort` has done its thing, the linked lists `sorted[largecount]`, ..., `sorted[1]` collectively contain all the words of the input file, in decreasing order of frequency. Words of equal frequency appear in alphabetic order.” [Bentley 86]

G. Design D7 Using a Hash Trie

```

module D7(k: nat, fin: word, fout: word) is (
  import module lex, hash-trie;
  init (
    value itxt := file-content(fin);
    lex.init(itxt);
    var t := hash-trie.init(lex.nextlexeme);
    t.build-all-words();
    file-content(fout) := t.find-frequent-words(k);
  ));

module hash-trie is (
  let (vi, pn, ui) == t.search(w);
  let old-word(w) == pn = #w;
  let incr-count(w) == cnt(nxt(ui)) += 1;
  let add-new-word(w) == t.mk-hash-trie(...);
  procedure build-all-words() is as-in module D6;
);

```

X. RELATED WORK

Knuth originated the idea of ‘literate programming’, and has contributed several examples, both large (T_EX [Knuth 86a], Metafont [Knuth 86b]) and small ([Knuth 84] and [Bentley 86]).

The tradition pioneered by [Kernighan and Plauger 76] is continued in such books as [Comer 84], [Tanenbaum 87], [Wirth and Gutknecht 92] and [Fraser and Hanson 95], which include complete listings of source code, along with cogent explanations of why they work. These are impressive accomplishments. But all the above examples neglect to emphasize design descriptions, and concentrate on implementation details.

There is a fundamental difference between Knuth’s examples, and the books by others mentioned above. Knuth’s examples are meant for peers to read, understand and evaluate, whereas the above books are for students to emulate and learn the programming techniques. The designs of Knuth’s examples get buried in a myriad of implementation surface details. The descriptions in the other books are too imprecise.

There are now a large number of large open source projects, with thousands of pages of documentation, but with hardly any design descriptions. [Lenberg et al. 2015] [Salman et al. 2015] [FM 2018]

A. Open Source Projects

Open source projects rarely have Requirements, Specs, Design and Implementation documents. There is typically a so-called README file describing the software briefly, and explaining how it can be built and installed.

B. Rigorous Description of Software Designs

The problem and Knuth’s solution use structures that are known as dictionaries, and hash tries, and notions such as lexical structure. These have been the main source of ambiguity.

We give precise definitions for all these. Natural languages, and so-called pseudo-code are unsuitable.

It is important to note that our definitions of these structures are not only precise, and rigorous but also *formal*. Most definitions in combinatorics and algorithms books, including those of Knuth, are *informal*.

We write these definitions in a software design language, named $\hat{O}M$, which is as “formal” as a typical programming language is these days (2018). We minimally describe the notation and semantics of $\hat{O}M$ as needed. Full details of $\hat{O}M$ are given in [Mateti 2018].

XI. EVALUATION

This section is intended to give an objective evaluation of the experience we have had in doing the exercise reported in this paper. We are concerned with three factors in this evaluation: the complexity of the example, a retrospect of $\hat{O}M$ ¹², and our stylistic use of $\hat{O}M$.

A. Critique of Knuth’s Solution

As we have mentioned earlier, the complexity of this paper as an example of a design document is due to the intricate structure of the hash trie. The hash trie has some nice properties such as storage efficiency and alphabetical orderedness of its content. Unfortunately, these properties do not come for free. McIlroy comments that “Knuth has shown us here how to program intelligently, but not wisely. ... He has fashioned a sort of industrial strength Fabergé egg — intricate, wonderfully worked, refined beyond all ordinary desires, a museum piece from the start.”

We ought to describe this complex structure and the operations to access it in such a way as to exhibit the reasons why it has these properties and how they are preserved. In addition to that, we are very much concerned with precision throughout the descriptions.

However, the complexity inherent to Knuth’s solution motivated us in choosing it as the appropriate material to test the suitability of $\hat{O}M$ to the design of real software systems.

B. Why is it this long?

An alternate design solution by David Hanson [Van Wyck 87] is less than 5 pages. Knuth’s paper [Bentley 86] is much shorter than the present paper. Why? It is longer mostly because of the explanation of the design language details.

C. $\hat{O}M$, Our Design Language

The following criteria are important in evaluating a software design specification language: availability of versatile high level data structures, expressibility of algorithmic descriptions and design decisions, degree of precision realizable in design specifications, executability of designs, and support to design methodologies and principles. These data structures together should allow a “good” designer to compose an abstract representation of any data object necessary to express designs.

In addition, a design language must offer an adequate and efficient syntax if it is going to be used by human designers. The syntax of the language must allow the designer to express designs in a natural manner, without having to cope with overwhelming syntactic details. Short hand builtin notations help avoid unnecessarily lengthy description. With respect to these criteria $\hat{O}M$ appears to be fairly satisfactory.

Due to its nature, $\hat{O}M$ can be used in various styles [?]: functional style, imperative style, sometimes logic style, and very often a mixture of the above. However, the functional and logic styles are more appropriate for high level design specifications whereas the imperative style is used for low level design specifications.

The $\hat{O}M$ usage in the current papaper is a limited subset of the full language Constructs to support expression of design decisions such as those concerning resource usage and choices among alternative designs are still under investigation. In addition, as related to natural design expressibility, the language will have some mechanism for importing notations from the problem domains of the designs. For instance, if we were to design a satellite control program, we would like to be able to express its design in terms of the vector notation that the typical physicist would use in this case. Finally, because $\hat{O}M$ is primarily directed towards an engineering contribution to software design, it does not yet provide any particular documentation support beyond the usual commentaries.

XII. CONCLUSION

Our example is the result of a reverse engineering process. The specifications and designs were written after studying the implementation. The material presented in this paper turned out to be longer (in terms of text size) than Knuth’s implementation for instance. However, one must note that formal specifications and designs, as presented above, convey more information than does their corresponding implementations. Through the example presented above, we have demonstrated some features of $\hat{O}M$ and have shown that indeed the language is powerful enough to support the design complex software systems. With its constructs, we have been able to precisely “specify a design solution” to the CWP common words problem.

REFERENCES

- [Bentley 86] J. Bentley, D. E. Knuth and M. D. McIlroy, “Programming Pearls,” a column in *Communications of the ACM*, Vol. 29, No. 6, 471-483. **i, 1, 2, 7, 10, 11, 12**
- [Cohen et al. 86] Bernard Cohen, W.T. Harwood, and M.I. Jackson, The specification of complex systems, Addison-Wesley, 1986, ISBN: 0-201-14400-X.
- [Comer 84] Douglas Comer, Operating System Design, Prentice-Hall, c1984-c1987, ISBN: 0-13-637539-1 (v. 1), ISBN: 0-13-637414-X (v. 2). **1, 11**
- [Diby et al. 18] Kouakou Diby, and Others, Examples in $\hat{O}M$ of Software Specs and Designs, 1988 - 2018.
- [Diby 90] Kouakou Diby, *Foundations of Hierarchical Design Methods for Software*, Ph. D. Dissertation, Wright State U, June 90.

¹²drop this

- [FM 2018] Klaus Havelund, Jan Peleska, Bill Roscoe, Erik de Vink (Eds), 22nd International Symposium on Formal Methods, FM 2018, Oxford, UK, July 15-17, 2018, [1](#), [11](#)
- [Fraser and Hanson 95] Christopher W. Fraser and David R. Hanson, *A retargetable C compiler : design and implementation*, Benjamin/Cummings Pub. Co. ISBN: 0-8053-1670-1. [1](#), [11](#)
- [Hayes 93] Ian Hayes (Editor), *Specification case studies*, Prentice Hall, 1993, 2nd ed. ISBN: 0-13-832544-8.
- [Kernighan and Plauger 76] Brian Kernighan, and Plauger, *Software Tools*, Prentice-Hall. [1](#), [11](#)
- [Knuth 97] Donald Knuth, *The Art of Computer Programming Vol. 1: Fundamental Algorithms*, Addison-Wesley, 1997, 3rd ed, ISBN: 0-201-89683-4 (v. 1). [7](#)
- [Knuth 73] Donald Knuth, *The Art of Computer Programming, Vol. 3: Sorting and Searching*, Addison-Wesley, Reading, Mass., Section 6.3.
- [Knuth 84] Donald E. Knuth, "Literate Programming," *Computer Journal*, Vol. 27, No. 2, 97–111. Reprinted in *Literate Programming* (book), Center for the Study of Language and Information, c1992, ISBN: 0-937073-80-6. [1](#), [11](#)
- [Knuth 86a] Donald E. Knuth, *TeX: The Program*, Addison-Wesley, Reading, Mass. [1](#), [11](#)
- [Knuth 86b] Donald E. Knuth, *MetaFont: The Program*, Addison-Wesley, Reading, Mass. ISBN: 0-201-13438-1. [1](#), [11](#)
- [Knuth and Levy 94] Donald E. Knuth, and Silvio Levy, *The CWEB System of Structured Documentation: version 3.0*, Addison-Wesley, c1994, ISBN: 0-201-57569-8.
- [Lenberg et al. 2015] Per Lenberg, Robert Feldt, Lars Göran Wallgren, Behavioral Software Engineering: A Definition and Systematic Literature Review, *Journal of Systems and Software*, Volume 107, 2015, Pages 15-37, ISSN 0164-1212, <https://doi.org/10.1016/j.jss.2015.04.084>. [1](#), [11](#)
- [Mateti 2018] Prabhaker Mateti, "ÔM: A Design Specification Language", Unpublished Manuscript, 45+ pp., 2018, Department of Computer Science and Engineering, Wright State University, Dayton, OH 45435. [2](#), [12](#)
- [Meyer 85] Bertrand Meyer, "On Formalism in Specifications," *IEEE Software*, Jan 1985, 6–26. [2](#)
- [Salman et al. 2015] Iftaah Salman ; Ayse Tosun Misirli ; Natalia Juristo, Are Students Representatives of Professionals in Software Engineering Experiments? 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering [1](#), [11](#)
- [Tanenbaum 87] Andrew Tanenbaum, *Operating Systems: Design and Implementation*, Prentice-Hall, ISBN: 0-13-637406-9. [1](#), [11](#)
- [VDM] Vienna Development Method, <http://www.ifad.dk/vdm/vdm.html> [1](#)
- [Van Wyck 87] Christopher J. Van Wyck, "Literate Programming," a column in *Communications of the ACM*, Vol. 30, No. 7, pp. 594-599. [1](#), [2](#), [12](#)
- [Wirth and Gutknecht 92] Jürg Gutknecht, and Niklaus Wirth, Project Oberon. The Design of an Operating System and Compiler, ACM Press, 1992, Addison-Wesley Publishing. [1](#), [11](#)

[Draft PDF generated October 5, 2018.]