

Quickly Generating Billion-Record Synthetic Databases

Jim Gray, Prakash Sundaresan

Digital Equipment Corporation, 455 Market, San Francisco, CA., 94105. {Gray, Prakash}@ SFBay.enet.dec.com

Susanne Englert

Tandem Computers Inc., 19333 Vallco Parkway, Cupertino, CA, 95014. Englert_Susanne @ Tandem.com.

Ken Baclawski

Computer Science, Northeastern University, 360 Huntington Av. Boston, MA. 02115. kenb @ ccs.neu.edu

Peter J. Weinberger

Bell Laboratories, 600 Mountain Ave, Murry Hill, NJ., 07974. pjw @ research.att.com

Abstract: *Evaluating database system performance often requires generating synthetic databases – ones having certain statistical properties but filled with dummy information. When evaluating different database designs, it is often necessary to generate several databases and evaluate each design. As database sizes grow to terabytes, generation often takes longer than evaluation. This paper presents several database generation techniques. In particular it discusses:*

- (1) *Parallelism to get generation speedup and scaleup.*
- (2) *Congruential generators to get dense unique uniform distributions.*
- (3) *Special-case discrete logarithms to generate indices concurrent to the base table generation.*
- (4) *Modification of (2) to get exponential, normal, and self-similar distributions.*

The discussion is in terms of generating billion-record SQL databases using C programs running on a shared-nothing computer system consisting of a hundred processors, with a thousand discs. The ideas apply to smaller databases, but large databases present the more difficult problems.

1. Introduction

Loading or generating large databases may take days or weeks. The goal here is to quickly generate a large database by using parallel algorithms and execution. To make the problem concrete, the goal is to generate a billion record ACCOUNTS table for the TPC-A benchmark [TPC]. Sequential algorithms to generate and load this table would take several days. The goal is to invent algorithms and techniques that generate this hundred-gigabyte table and its indices in an hour.

In outline, the paper first postulates a model of parallel computer hardware and software, so that we can quantify the performance of each algorithm. Then, the paper shows

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

SIGMOD 94- 5/94 Minneapolis, Minnesota, USA
© 1994 ACM 0-89791-639-5/94/0005..\$3.50

how to convert a sequential load to a parallel load by partitioning the job and forking a process-per-partition. Next, the tasks of synthetic data generation are investigated. Parallel algorithms are given for generating dense-unique-pseudo-random sequences, and for generating indices on these sequences. After that, the paper investigates generating non-dense non-uniform distributions with special attention paid to Zipfian and self-similar distributions.

2. The Computation Model

We assume a shared-nothing computer architecture typified by the Tandem and Teradata machines, by workstation clusters from DEC, IBM, HP, Novell, and Sun, and by processor arrays like the Intel Hypercube [Stonebraker, Horst, Teradata, DeWitt 1, DeWitt 3]. In these systems, each processor has a private memory and one or more discs. The processors are connected via a high-speed network and processes communicate via messages. Parallelism and minimal process interaction are major design goals in these systems. Shared-disc systems like IBM's Sysplex and Digital's VMSclusters are gravitating toward this shared-nothing architecture as the number of processors grows.

The ideas presented here apply to a spectrum of execution environments – a *many-small environment* of hundreds of processors typified by Teradata and Gamma, and a *few-big environment* of tens of processors typified by Tandem and VMScluster. In Table 1, the two systems each have 3 BIPS (billion instructions per second) of processing power, 10 GB of RAM, and 1 TB of disc storage (1000 discs). The prices are estimates: 100 \$/MIPS, 30 \$/MB RAM, and 1 \$/MB disc. These approximate 1994 prices.

The ideas here apply to both architectures – but to simplify the discussion, the many-small-processors design is assumed here.

In this presentation, CPUs are numbered 0, 1, ..., 99. Discs are attached to particular CPUs, and are correspondingly named D\$ddd where ddd is the disc number [0...999]. For example, disc 223 is named D\$223 and is the third disc of processor 22.

Table 1: Two large database machines of 1994, each costing 1.6M\$.			
MANY-SMALL PROCESSORS (100)		FEW-BIG PROCESSORS (10)	
CPUS + MEMORY	DISCS	CPUS + MEMORY	DISCS
number: 100 X (30 MIPS+100 MB)	1000 X (1 GB/disk)	10 X (300 MIPS+1 GB)	1000 X (1 GB/disk)
sum	3 BIPS +10 GB	3 BIPS + 10 GB	1 TB
price	100 x (3 k\$ + 3 k\$)	10 x (30 k\$ + 30 k\$)	1000 x 1 k\$
sum price:	600 k\$	600 k\$	1 M\$
total price	1.6M\$	1.6M\$	

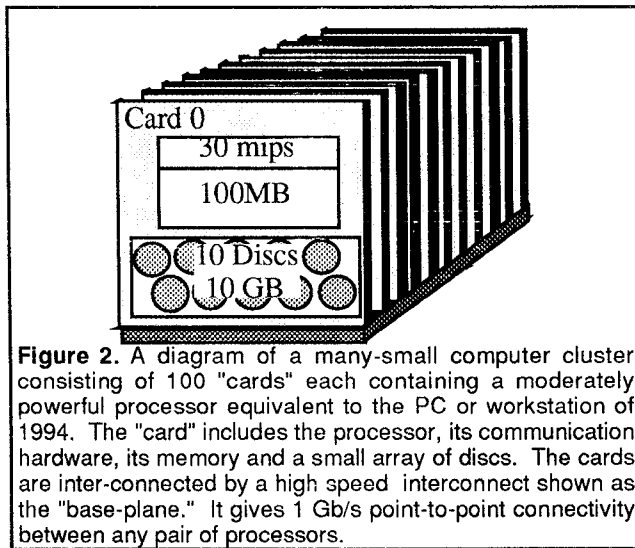


Figure 2. A diagram of a many-small computer cluster consisting of 100 "cards" each containing a moderately powerful processor equivalent to the PC or workstation of 1994. The "card" includes the processor, its communication hardware, its memory and a small array of discs. The cards are inter-connected by a high speed interconnect shown as the "base-plane." It gives 1 Gb/s point-to-point connectivity between any pair of processors.

Algorithms here avoid non-sequential record reads and writes because sequential operations can run at disk device speed (5 MB/second = 50,000 records/second), while random disk operations run a thousand times slower (50 IO/sec = 50 records/second) due to disk seek and rotational latencies. One-way process-to-process messages have a fixed cost of 3K instructions, and a marginal cost of one instruction per byte [Uren, Kronenberg, Thekkath]. Algorithms here minimize sending messages. If necessary, they send a few large messages rather than many small ones. Table 3 summarizes these computational costs

Table 3: Cost of basic operations:	
record size	100 bytes
sequential SQL read/insert	5,000 instructions + .01 IO
random SQL read/insert	25,000 instructions + 1 IO
disk sequential read/write rate	5 MB/sec = 50,000 rec/s
disk random read/write rate	50 IO/sec = 50 rec/s
cpu cost of x-byte message	3,000 + x instructions

For simplicity, assume that the data of a table or index is range-partitioned among all the 1000 discs in partitions of equal size. Each system has a slightly different syntax for table partitioning. To be concrete, we use Digital's Rdb syntax [Hobbs]. Partitions are called storage areas in Rdb and are often named by the disc on which they reside. Using Rdb syntax, the ACCOUNTS table of the TPC-A Benchmark [TPC] on the many-small system would be defined by:

```
create table ACCOUNTS (
  ID          unsigned integer  not null,
  BALANCE    decimal(12,2)     not null,
  CUSTOMER   unsigned integer  not null,
  FILLER     character(92)     not null,
  primary key (ID));
create storage map PARTITION for ACCOUNTS
store using (ID)
  in D$000 with limit of 0999999,
  in D$001 with limit of 1999999,
  ...
  in D$998 with limit of 998999999,
  otherwise in D$999;
```

When an application selects, inserts, updates, or deletes a record of the ACCOUNTS table, the SQL system locates the record in the correct partition. The application is unaware of the partitioning. Most SQL systems provide a variant of this transparent partitioning.

The computation cost model is simple. A sequential record read or insert costs 5,000 instructions and a fraction of an IO.

3. Sequential Database Generation

The discussion begins by showing how to sequentially generate and populate a table. This algorithm is then generalized to one that generates each partition in parallel. The TPC-A ACCOUNTS table will have one billion hundred-byte records (.1 TB) partitioned among the 1000 discs as described in the data definition Program (1) above. Each of the one thousand discs will store 100 MB of data as a B-tree [Knuth]. Since B-trees are only 69% full at equilibrium, each disc will use 150 MB of storage to hold its B-tree. This is well below the 1 GB capacity of small discs. The remainder of the disk stores data from other tables.

The ACCOUNTS table can be sequentially populated by the simple SQL + C Program (2). Section 4 will show how to fill in the customer number with a unique and uncorrelated customer ID. The customer number should be unique and uncorrelated with the account number. Here it is set to zero in all records. Also, standard SQL does not allow table

```
/*sequentially load records with key in [base,base+count) */ (2)
void sequential_load( char *tablename, /* table to be loaded */
  long base, /* start key of load */
  long count) /* number to load */
{
  exec sql begin declare section; /* declare host vars */
  long key; /* */
  exec sql end declare section; /* load each key */
  for (key = base; key < base + count; key++) /* in range */
    exec sql insert into :tablename /* */
      values (:key, 0, 0, ""); /* */
} /* end sequential_load() */
```

names to be host variables, but this and later programs assume it is allowed. One would have to use dynamic SQL to get this effect in a standard SQL system.

The entire database for a 10,000 tps-A database [TPC] can be loaded by Program (3) which loads a billion accounts, a hundred thousand tellers and ten thousand branches:

```
void sequential_load("accounts", 0, 1000000000);
void sequential_load("tellers", 0, 100000);
void sequential_load("branches", 0, 10000);
```

4. Parallel Database Generation

Programs (2) and (3) are fine for loading small tables (less than a million records), but they use a single processor and so run one hundred times slower than an algorithm that divides the task into a hundred smaller ones each running in parallel on a separate processor. Program (3) runs at 6,000 records/second given the performance assumptions of Table 1 and 3 (5,000 instructions per insert, 30 MIPS, implies 6,000 inserts per second.) At that rate, the billion-record load would take almost two days. The same load could run in twenty minutes if done in parallel on the hundred-processor cluster described by Figure 2.

Parallel algorithms require a way to create processes on specific CPUS. Assume that a process can be created in a cpu by calling:

```
int fork(int cpu); /* create a process */(3)
```

The `fork()` procedure is much like the UNIX `fork()` [Tanenbaum] but has a parameter specifying the `cpu` in which the process is to be created. As with UNIX, the `fork()` returns zero to the child process, but returns the child process identifier to the parent (forking) process. The child is a clone of the parent, executing the same program with the same current state, but a completely separate process environment.

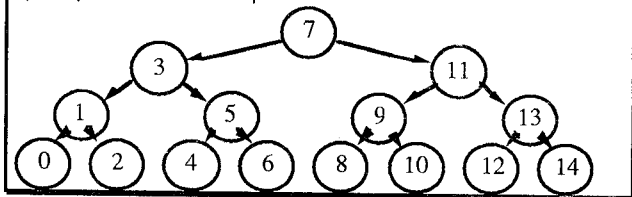
Program (5) shows how to convert Program (2) sequential load to a parallel loader. It starts by spawning a load process in each processor. The spawning time is minimized by forking a binary tree of processes, forking 2^n processes at depth n of the tree. Each node follows the logic:

(1) If I am not a leaf, fork my left child and right child.

(2) Load my partition.

Figure 4 illustrates the idea. If each fork takes about one second (a high estimate), the entire startup of one hundred processes will complete within eight seconds. Assuming the forking logic copies code and data from the forker, there will be no bottlenecks at startup.

Figure 4. The process forking logic for a cluster of 15 processors. `cpu 7` forks loaders into `cpus 3` and `11`, then `cpu 7` proceeds to load partition 7.



Program (5) is easily modified to fork a process-per-disc rather than a process-per-processor if the processors are not the bottleneck. Notice that each generator process uses the same table name to generate data (database systems call this *location transparency*). The table partitioning criterion causes records to go to each loading process's local disc.

Parallelism often suffers from problems of startup, interference, and skew [Gerber, Smith]. Program (5) minimizes startup problems by parallel forking. Once the load process begins, the underlying system acquires locks on partitions rather than on whole tables – at least that is the way many SQL systems work. So, each partition loader can execute in parallel and in isolation. The load operation is typically not covered by transaction protection, so the recovery log is not a bottleneck—rather it uses the old-master-new-master recovery technique of dumping a copy of the table when the load completes. Consequently, startup, interference, and skew should not be a problem for parallel load. Using the assumptions of Figure 2 and Table 3, the algorithm should generate 600,000 records per second (60 MB/s) and generate 1B records in less than 30 minutes.

```
/* parallel load records into tablename */(5)
/* first fork right and left child processes if right != left */
/* then sequentially load the partitions of this cpu */
#define CPUS 100 /* number of cpus in cluster */
void parallel_load(char * tablename, /* table to be loaded */
                  long records, /* records to load */
                  long left, /* cpu of left subtree */
                  long right, /* cpu of right subtree */
                  integer depth); /* recursion depth */
{ long my_cpu = floor((left+right)/2); /* id of my cpu in cluster */
  long part_size = floor(records/CPUS); /* number of records in my */
  /* table partition */
  long him; /* id of forked process */
  depth++; /* increase recursion depth */
  if ( depth == 1 ) /* return after fork root */
  { if ( fork(my_cpu) ) return; }; /* in center cpu. */
  /* spawn left subtree of processes */
  if ( left < my_cpu ) /* if need left child */
  { him = fork( floor( (left + my_cpu-1) / 2) ); /* fork him */
    if (him == 0) /* child code */
    {parallel_load(tablename, records, left, my_cpu-1, depth); return;}
  }
  /* spawn right subtree of processes */
  if ( my_cpu < right ) /* if need right child */
  { him = fork( floor( (my_cpu + 1 + right) / 2) ); /* fork him */
    if (him == 0) /* child code */
    {parallel_load(tablename, records, my_cpu+1, right, depth); return;}
  }
  /* fill the partition of this process */
  sequential_load(tablename, my_cpu * part_size, part_size); /*
  }; /* end of parallel_load() */
/* invoke parallel loader for a billion account records on 100 cpus */
void = parallel_load("accounts", 1000000000, 0, CPUS - 1, 0); /*
```

5. Generating Dense Unique Random Data

The parallel data generator of the previous section correctly generated the sequential account numbers but it did not generate the customer numbers – it just left them as zeros. This is an example of a general problem: generating synthetic databases often requires a sequence of numbers (i.e., field values) with the following properties:

Dense: All the integers in $[0..N]$ appear in the sequence.

Unique: Each integer appears exactly once.

Random: The sequence appears to be "random" (pseudo-random).

These properties are needed to make the cardinalities of selection expressions and join expressions predictable – for example each customer should have exactly one account. Some applications need synthetic data that is not uniformly distributed. Section 8 gives some ways to transform uniform distributions into Gaussian, exponential, Zipfian, and other distributions.

There are many ways to generate dense-unique-random numbers. The original generator for the Wisconsin Benchmark [Bitton 1] kept an initially zero *bitmap* of length N and used the system random number generator to pick the next free element for the series. The ASAP generator [Bitton 3] replaced *bitmap* with a *shuffle* that built an array of pairs $\langle (i, \text{random}()) / i=1, \dots, N \rangle$. The array was then sorted on the second element to shuffle the first element.

The *bitmap* algorithm uses order N space and order N^2 time. The *shuffle* algorithm takes $N \log N$ time and linear space, so is clearly superior to the *bitmap*. If space is not an issue, *shuffle* is a good way to generate a dense-unique random series. But for large databases a more space-efficient algorithm is needed. The obvious choice is to generalize the *shuffle* to a sort by creating a SQL table of two columns: one containing the dense sequence $1..N$ and the second column containing a random sequence based on a popular random number generator. Then the dense sequence is ordered by the random sequence. Program (6) demonstrates this. It takes $N \log N$ time and $\sim \sqrt{N}$ main memory (sorting) and $\sim N$ disc space (storage of table). For a billion records, the many-little configuration can do the sort in about thirty minutes.

The sort scheme is inconvenient because it constructs a set of files to drive the data generation. It would be more convenient to have a simple subroutine that could generate the next element of the desired sequence in constant time and space (say 25 instructions and 100 bytes of storage). The idea for such an algorithm is to generate the numbers using a generator of the cyclic group of integers under

```

exec sql create table T (          /* create temp table      */ (6)
        sequence integer,        /* column is dense unique 1..N */
        rand integer);          /* column has random values */
for (i = 1; i <= N; i++)          /* fill table with N pairs */
    exec sql insert into T values (:i, :random()); /*
        /* sequence values are dense and unique. */
        /* If ordered by the rand field, they will be random. */
exec sql declare answer cursor for /* define a sql cursor to
select sequence from T order by rand; /*get data in random order */
exec sql open answer;            /* read table via that cursor */
for (i = 0; i < N; i++)          /* set next_value from
    { exec sql fetch answer into :next_value; /* next record
      /* process next value
    }
exec sql close answer;          /*

```

multiplication. In essence, a random number generator is constructed for elements in the desired range. The algorithm is: Pick a prime P larger than N and a generator G for the multiplicative group modulo P . Then the series is:

$$\langle G^i \bmod P \mid i = 1, \dots, P \text{ and } (G^i \bmod P) \leq N \rangle. \quad (7)$$

Program (8) generates the series suggested by Equation (7). This scheme, due to Gray and Englert, was used to generate large Wisconsin Benchmark databases on the Intel Hypercube and is now the standard way to generate Wisconsin databases [DeWitt 2]. To understand how it works, consider the numbers between 1 and 10. The powers of 2 mod 11 form a dense unique sequence of these numbers: 2, 4, 8, 5, 10, 9, 7, 3, 6, 1.

The generator scheme is ideal—it uses linear time and constant space. If N is a prime the multiplicative group consisting of $[1..N-1]$ has many generators: elements whose powers enumerate the group without repetition until they generate 1. But not all generators give a good pseudo-random sequence. There are many tests for "randomness." We used the spectral test recommend by Knuth [Knuth]. In his terminology, all the random number generators in this paper pass the spectral test "with flying colors" in dimensions 2 through 6. We applied the test to primes just larger than powers of ten and recommend the generators of Table 5. Section 7 uses powers of 2 instead of primes.

Table 5. A list of recommended primes and generators for each decade from 10 to one billion.

Decade	Prime	Generator
10	11	2
100	101	7
1,000	1,009	26
10,000	10,007	59
100,000	100,003	242
1,000,000	1,000,003	568
10,000,000	10,000,019	1792
100,000,000	100,000,007	5649
1,000,000,000	2,147,483,647	16807

```

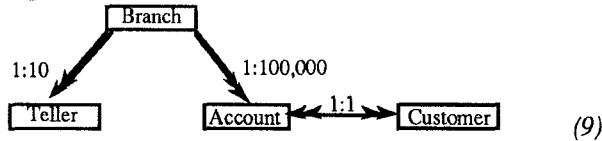
#define P xxx; /* see Table 5 for good values */ (8)
#define G xxx; /* see Table 5 for good values */
static seed = G; /* start the seed at G */
long next_value(long N) /* function to compute next value */
{ do {seed = (G * seed) % P } /* seed = next in series */
  while ( seed > N ); /* mod prime (discard all > N) */
  return seed; /* return new value */
} /* end of generator */

```

When the prime and generator are large compared to the machine's arithmetic, one needs to use the technique due to Schrage [Schrage] to keep the results from overflowing.

6. Generating Correlated Random Data

The ideas of the previous section can now be applied to generate a complete table. Program (3) gave an example of generating several tables with related statistics: the idea there was that the branch:teller:account cardinalities should be in the ratios 1:10:100,000, as in the schema (10). This is a general phenomenon, but the requirements are often more complex. One requirement that was skipped in Program (2) was that the customer id field be unique and be uncorrelated with the account id – rather it was just filled with zeros. The requirement is that each customer have a unique bank account as shown in (9).



Using the ideas of the previous section, Program (2) can be refined to generate each partition of the account table in parallel, including the random-unique-dense customer number as in Program (13).

The generator of the i 'th cpu is invoked as:

```
(void) sequential_load(    "accounts",    (10)
                          records,
                          records/CPUS,
                          i);
```

Each of the one hundred processors will compute the same random series based on the `next_value(records)` procedure. But each generator will only use one hundredth of the series. If the series is $s_0, s_1, s_2, \dots, s_{99999999}$, the i 'th cpu will use the subsequence $s_i, s_{1cpus+i}, s_{2cpus+i}, s_{3cpus+i}, s_{4cpus+i}, \dots$. Each of these subsequences is random and unique, and the union of them is dense. The i 'th generator begins by calling `next_value()` "i-1" times to skip over the values of the other generators. Then after inserting a record, the generator calls the `next_value()` CPUS times, but only uses the last value returned.

This design is inefficient: the series is computed one hundred times and each generator only uses 1% of the values it generates. The premise is that calls to `next_value()` are cheap (~25 instructions) so that 100 calls (2500 instructions) is small compared to the insert cost (~5000 instructions). The generator should produce 1B records in less than 1 hour on the many-small configuration.

Scaling this algorithm to thousands of generators requires a variation that has less wasted work. One might partition the series into 1000 segments and precompute the starting point P_i for each partition. These values could be stored in a global array. Then each partition generator would start at s_x , for some x , and would use the sequence $s_x, s_{x+1}, s_{x+2}, \dots$ and no calls to `next_value()` would be wasted.

A different approach avoids this pre-computation and minimizes wasted computation. Table 5 shows that for large N , P can be chosen just slightly larger than N (within 1% of it). Suppose we generate a database of $P-1$ elements rather than N elements. In that case, no members of the series $\langle G^i \bmod P \mid i = 0, \dots, P \rangle$ would be discarded (compare this to Equation (7)). In turn this means that each partition can compute its next element by multiplying the previous element by $G^{CPUS} \bmod P$. Let $n = \lceil N/CPUS \rceil$ and assign the following series to the j 'th partition for $j = 0, \dots, CPUS - 1$:

$$\langle G^{CPUS \cdot i + j} \bmod P \mid i = 0, \dots, n-1 \rangle \quad (11)$$

This series is very easy for partition j to compute. First it computes the first element $A \equiv G^j \bmod P$ and then $B \equiv G^{CPUS} \bmod P$. These two numbers can be computed in $\ln(CPUS)$ multiplies and divides. Then the j 'th partition uses the series:

$$\langle (A \cdot B^i) \bmod P \mid i = 0, \dots, n-1 \rangle. \quad (12)$$

In this approach, there is no need to precompute a partition table and there are no wasted calls to `next_value()`. If we accept this relaxed definition of partitions (the last partition may be slightly larger than the others and some elements may be a little larger than N), then it will turn out that computing indices is much easier

These techniques generate tables with 1:N relationships. Suppose, as in (9), the BRANCHES table is to have 1,000 records and the TELLERS table is to have 10,000 records. If P is chosen as 1,009 and G is chosen as 229 (as Table 5 recommends), then BRANCHES can be generated as in Program (10). By using the same prime and generator for the BRANCH field of the 10,000 record TELLERS table, the

```

/* sequentially load records into i'th partition of tablename */ (13)
void sequential_load_partition(
    char *tablename,          /* name of table to be loaded */
    long records,            /* number of records in table */
    long count,              /* # records in each partition */
    long part_no)           /* number of partition to load */
{
    long i, j;
    long base = part_no * count; /* loop control variables */
    long parts = records / count; /* partition's start key */
    exec sql begin declare section; /* partitions in file */
    long key, customer;          /* SQL variables */
    exec sql end declare section; /* ids */
    for (i = 1; i < part_no; i++) /* i'th processor skips to */
        customer = next_value(records); /* p(i) in generator sequence */
    for (key = base; key < base + count; key++) /* generate partition */
        {customer = next_value(records) - 1; /* get next customer number */
        exec sql insert into :tablename /* */
            values (id, balance, customer, filler) /* */
            values (:key, 0, :customer, ""); /*generate 1/n'th */
        for (j = 1; j < CPUS; j++) /* skip seed to p(i+CPU) */
            customer = next_value(records); /* i.e. skip p(j) of others */
        } /* end of sequential_load() */
}
  
```

TELLERS table will have exactly 10 tellers for each branch identifier. More generally, each record in one table will match the value of M records in the second table if the second table has M times as many records, and the "join-fields" of the two records use the small-table generator.

7. Generating Indices on Random Data

All the programs so far have carefully (but implicitly) generated data in primary key order. The next generated record is placed right after the previous record in the B-tree or other clustering mechanism. This means that the programs have generated the data in sequential order and so disc IO time has not been an issue. Modern discs can absorb data at 5 MB/s – with the possibility of much higher data rates if striping is used [Kim]. Assuming 100 byte records, this is 50,000 records per second. If each generated record went to a random disk page, data rates would drop by a factor of 2,000 to 25 records per second – each record would cause a seek, rotation, a read transfer and then a rotation and a write transfer. On 1994 discs, each random disk IO consumes about 20ms of disc time and the rate is at most 50 IO/s. So, it is *essential* that records be generated in sequential order unless the entire table fits in main memory.

Program (13) generates the `ACCOUNT` table in `ACCOUNT.ID` order; but it generates the `ACCOUNT.CUSTOMER` field in random order. Applications often need an index on such random fields. For example, an `ACCOUNT.CUSTOMER` index would allow quick lookup of a customer's `ACCOUNT` given the customer number. Such an index is defined in SQL by:

```
create unique index account_customer
on account (customer);
```

The index is range partitioned among the thousand disks, but that partitioning is not shown. The index is actually a table with the schema:

```
create table account_customer (
customer unsigned integer not null,
id unsigned integer not null,
primary key (customer),
foreign key (id) references account(id)
);
```

If (k, i) is a record in the `account_customer` index then: k is the i 'th value returned by `next_value(records)` or using G and P as defined in Table (5) and program (8):

$$k \equiv G^i \pmod{P}. \quad (15)$$

More formally, i is the discrete logarithm of k with respect to $G \pmod{P}$ [Coppersmith].

How can this index be generated in parallel with linear speedup and scaleup? One could compute the discrete logarithm, but [Coppersmith] indicates that each computation would be millions of instructions. Three schemes can be used:

Scan-and-sort: Read the base table in parallel, projecting out the two desired fields, and parallel sort the result into the target index. In SQL this would be expressed as:

```
insert into account_customer (16)
select customer, id
from account
order by customer;
```

Generate-and-sort: In each processor, generate the index data to be stored by that processor's disks, sort it, and then insert it into the local index partitions.

Compute: Compute the discrete logarithm quickly and generate the index in the same way one generates the base table.

Our index is one billion records of eight bytes each, 8 GB in all. So, each of the 100 processors must deal with an 80 MB partition of the index – just fitting in each processor's 100 MB memory. The scan-and-sort approach, lets processors generate index data in parallel to data generation, sending index records to the appropriate partitions (cpus) as the base table is generated locally. The receiving processors can sort the indices locally in their memory as the data arrives or is generated. This is a credible and scaleable technique, needing about 10 MB/s network bandwidth to move the 8 GB from source to destination for a one-hour job. But the technique is memory intensive, just barely fitting in the processor memories. If the table keys were larger of if there were more indices, then the scan-and-sort technique would require a disk-based sort.

Generate-and-sort is more cpu-intensive – but uses no network messages. Each processor generates the entire base-table sequence, and extracts the index subsequence that applies to the local processor. In particular, if each partition has R index records and if the whole sequence is k_1, k_2, k_3, \dots then the n 'th processor uses the subsequence:

$$\langle \langle k_i, i \rangle \mid nR \leq k_i < (n+1)R \rangle \quad (17)$$

These are the index entries for the n 'th partition. They are then sorted on the k_i attribute and inserted into the local index partition in sequential order. Program (18) shows the enumerate-and-sort algorithm.

The generation step of Program (18) should take about 30

```
/* global variables: seed , generator, CPUS, my_cpu */ (18)
/* R is records per partition */
void index_load(long records) /*
{ long R = records / CPUS; /* R records per partition */
long my_first_record = my_cpu * R; /* base of index partition */
long i, j = 0, customer; /* working variables */
exec sql begin declare section; /*
struct { long account; /* in-memory index */
long customer; /* to be sorted on the */
} sorted [R]; /* "random" customer id */
exec sql end declare section; /*
/* fill in the array with the unsorted values */
for (i = 0 ; i < records; i++) /* for each account number */
{ customer = next_value(records); /* get next customer number */
if ( (my_first_record <= customer) & /* if cust no in part */
(customer < my_first_record + R) /* range use the entry */
{sorted.account[j] = i; sorted.customer[j]=customer; j++;} /*
}; /* now sorted[i] = G inverse */
/* sort the array on the second attribute (customer) */
sort(sorted) on sorted.customer; /* sort array by customer */
for (j = 0; j < R; j++) /* copy array to the index */
exec sql insert into account_index /* scan sorted array insert */
values (:sorted.customer[j] , :sorted.account[j]); /* element */
}; /* end of index_load() */
```

```

#define POWER 32 /* P from Table 6 */
#define G 37117 /* G for P from Table 6 */
static long P = 1 << (POWER-2); /* P/4 (the field size) */
static long P_MASK = P - 1; /* mask to avoid mod P division */
/* do mod P multiplication, in the *4 + 1 space. The equation is:
/* ( a * 4 + 1 ) * ( b * 4 + 1 ) = ( a + b + 4 * a * b ) * 4 + 1
long mul(long a, long b) /*
{ return ((a + b + 4 * a * b) & P_MASK); }/*
/* discrete_log(k) returns the discrete log of k with respect to G
/* At entry (G^(ans) = 1 + 4 * k ) mod P
/* the invariant of the loop is: (G^(x+ans) = 1 + 4 * up) mod P
/* where the last i bits of up are zero
long discrete_log(long k) /*
{
long up = k; /* up is j with G^2^i removed
long i; /* index on radix bits of k
long x = 0; /* the answer we are building
long radix = 1; /* radix = 2^i in the loop below
long Gpow = (G - 1) / 4; /* Gpow = ((G-1)/4)^2^i
for ( i = 0; up; i ++ ) /* for each 2^i bit in up
{ if (up & radix) /* if 2^i divides k
{ x = x + radix; /* equation (20) says add 2^i
up = mul( up, Gpow ); } /* preserve loop invariant
radix = radix << 1; /* advance radix = 2^i
Gpow = mul( Gpow, Gpow ); /* advance Gpow = G^2i
} /* end of loop
/* now up = 0 so G^(x+ans) = 1 mod P due to the invariant.
/* by Fermat's theorem, x+ans = P so discrete log of k is P-x
return (P - x); /* return it
} /* end of discrete_log()

```

(21) Table 6 is a catalog of values for POWER and G that have passed the spectral test. One problem is that the series G, G^{2^1}, \dots are all congruent to 1 mod 4 - all their binary representations end in "01". To get a dense series, one must divide by 4. Program (21) encapsulates this *2+1 logic in the mul() routine.

The use of numbers near a billion strains the word size of 32-bit computers. In particular, if P is bigger than 2¹⁶ or so, multiplication modulo P cannot be done without some programming trick unless you have a computer with 64-bit registers. Schrage's technique, can be used to fit such arithmetic into small words [Schrage]. The mul() routine of Program (21) is easily modified to use that technique.

instructions per iteration. There are 1 billion iterations, so it will take 1000 seconds on a 30 MIPS computer. The sort deals with ten million records. It will need about 300 seconds. Once that is complete, the write of the data in bulk to the index (at 5000 instructions/insert) should take 300 seconds. This adds up to about thirty minutes. In summary, indices for large tables can be built in parallel while the base tables are being built. The generator can run in parallel with the base table generation if sufficient processors and memory are available.

The third index technique involves quickly computing discrete logarithms by carefully choosing G and P. That is, given alternate-key value k, quickly compute primary key value i such that:

$$k = G^i \text{ mod } P. \tag{19}$$

Solving this problem for arbitrary k when P is a large prime is believed to be quite difficult. Indeed, this is what makes some cryptographic protocols seem secure. Even for smaller primes, around a billion, each discrete logarithm calculation takes about \sqrt{P} time and space. The sorting algorithm above would be faster.

Picking P as a power of 2, and a generator G of the form $G = 4u+1$ where u is odd, allows computing discrete logs in n steps (log P time) and constant space. The following equation is the key to computing the discrete log of k when G and P satisfy these properties (it is proved by a simple induction argument):

$$G^{2^i} \equiv 1 + 2^{i+1} \pmod{2^{i+3}} \tag{20}$$

The technique is to use Equation (20) to compute x such that $G^{i+x} \equiv 1 \pmod{P}$. Then, using Fermat's Theorem, $x=P-i$. Program (21) implements this idea, it computes the discrete log of k (the 2-adic log) with respect to P.

In summary, indices on synthetically generated data can be built in one of three ways. Scan-and-sort, generate-and-sort, or compute. The computational method has some restrictions on the size of the table and on

Table 6. Powers of 2 and generators allowing fast computation of discrete logs (as in Program (21)).

Field	(max value)	POWER	G
1...2 ⁸ -1	255	10	29
1...2 ⁹ -1	511	11	29
1...2 ¹⁰ -1	1,023	12	53
1...2 ¹¹ -1	2,047	13	53
1...2 ¹² -1	4,095	14	117
1...2 ¹³ -1	8,191	15	125
1...2 ¹⁴ -1	16,383	16	229
1...2 ¹⁵ -1	32,767	17	221
1...2 ¹⁶ -1	65,535	18	469
1...2 ¹⁷ -1	131,071	19	517
1...2 ¹⁸ -1	262,143	20	589
1...2 ¹⁹ -1	524,287	21	861
1...2 ²⁰ -1	1,048,575	22	1,189
1...2 ²¹ -1	2,097,151	23	1,653
1...2 ²² -1	4,194,303	24	2,333
1...2 ²³ -1	8,388,607	25	3,381
1...2 ²⁴ -1	16,777,215	26	4,629
1...2 ²⁵ -1	33,554,431	27	6,565
1...2 ²⁶ -1	67,108,863	28	9,293
1...2 ²⁷ -1	134,217,727	29	13,093
1...2 ²⁸ -1	268,435,455	30	18,509
1...2 ²⁹ -1	536,870,911	31	26,253
1...2 ³⁰ -1	1,073,741,823	32	37,117
1...2 ³¹ -1	2,147,483,647	33	52,317
1...2 ³² -1	4,294,967,295	34	74,101
1...2 ³³ -1	8,589,934,591	35	104,581
1...2 ³⁴ -1	17,179,869,183	36	147,973

the generator, but is the most efficient approach for large tables. The computational approach is nicely suited to parallel algorithms.

8. Generating Non Uniform Data

Having explained how to generate unique and pseudo-random data, now consider generating non-uniform data distributions. The examples above all required dense-unique values. Often, the database needs values obeying some common distribution. The size of cities, lengths of words, and frequency of words is known to follow a Zipfian distribution. Measurement errors often obey a Gaussian distribution, and the inter-arrival intervals of events often follow a Poisson or negative exponential distribution. Such domains are easily generated by skewing a uniform distribution. This section catalogs the standard distributions, and adds a little to the generation of self-similar and Zipfian distributions.

Program (3) demonstrated the simplest case, repeating some value a constant number of times in another field. It generates ten accounts per branch – repeating each branch number ten times in the ACCOUNT.BRANCH domain. Suppose we wanted the number of child records to follow some more complex distribution. Then the code of program (21) might be appropriate.

Program (21) encapsulates the distribution of child cardinalities. It only remains to describe ways to generate popular distributions. The next page present such code for the popular districutons (uniform distributions have already been covered). See [Ripley, Jain, or Press] for more details. Assume `randf()` returns values distributed uniformly in $[0..1]$ and `random()` returns values distributed uniformly in $[0..∞]$ or some approximation thereof (e.g. $[0..2^{64}]$).

The only "new" idea on the next page is the self-similar distribution often used for situations following the 80/20 rule or some other skewed self-similar distribution. Self-similar distributions have the property that within any

region of the distribution, the skew is the same as in any other region. So, for example, all subranges of the 80/20 self similar distribution follow the 80/20 rule ($h=.20$).

A set of values of the form $1..k$ is called a "hot spot" if any of the values in this set has more weight (and hence is "hotter") than any value outside the set. Self-similar distributions are characterized by hot spots that have distributions similar to the entire distribution. Zipf distributions are characterized by a straight line frequency distribution when plotted on a log-log graph. It is commonly thought that self-similarity and the Zipf distribution are the same, or at least close. This misconception apparently stems from a misleading approximation made by Knuth (volume 3, page 398). Knuth's approximation is adequate for the statistic being computed there but should not be construed as asserting that the self-similar and Zipf distributions are close in the usual probabilistic sense. In particular, other statistics can yield very different results for the self-similar and Zipf distributions. However, Knuth's calculation can be modified to produce a reasonable approximation as follows. The function `zeta()` returns the sum

$$(1/1)^{\theta} + (1/2)^{\theta} + \dots + (1/n)^{\theta}. \quad (22)$$

The approximation uses Knuth's technique, but corrects the weight assigned to the first two values to improve the approximation. The log-log graph of Zipf's distribution with parameter 0.5 is shown on the next page. It shows the largest weight on 1, the second largest on 2, and so on.

If the hot spot is supposed to be randomly spread throughout a range of values, then it is necessary to permute the values randomly. In principle, generating a random permutation and generating a distribution are independent problems. However, hot spot distributions like the self-similar and Zipf distributions can be permuted more easily than general distributions. The trick is to assume that the "cold" values are uniformly distributed. This greatly simplifies the generation of the permutation since it now just involves choosing the relatively small number of "hot" values. This is especially important when

```

exec sql create table parent (master integer not null,          (21)
                             rest      char(96),
                             primary key (master) );
exec sql create table child (master integer not null,
                             detail    integer not null,
                             rest      char(92),
                             primary key (master,detail),
                             foreign key (master) references parent );
/*sequentially load records if key in [base,base+count) */
void sequential_load(char *parent, /* name of parent table */
                    char *child, /* name of child table */
                    long base, /* start key of load */
                    long limit) /* first key after load */
{
    exec sql begin declare section; /*
    long master, detail; /* master-detail key values */
    exec sql end declare section; /*
    for (master = base; master < limit; master++) /*for each parent */
    {exec sql insert into :parent values(:master, "");/*add it */
    count = distribution(); /* create count() child recs*/
    for ( detail = 0; detail < count; detail++ ) /*
    exec sql insert into :child values(:master, :detail, "");/**/
    }
}
/* end master loop */
/* endsequential_load() */

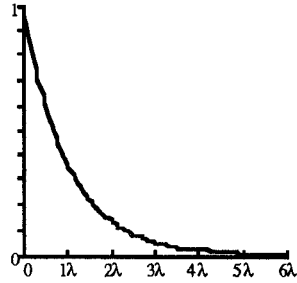
```

it is necessary to perform the computation in parallel. The "hot" values can be chosen at one node and broadcast as a table to the rest of the nodes. The algorithms above are used to compute an index into the table of "hot" values. If the index is too large, then a value is chosen uniformly at random from all possible values (ignoring whether it is already in the table of hot values). The index computations can be done independently at each node so long as the seeds for the pseudo-random number generator are chosen independently.

Negative exponential

equation: $f(x) = \frac{1}{\lambda} e^{-\lambda x} \quad : x \in [0, \infty]$

```
code:
double neg_exp(double lambda)
{ return ln(random()) / lambda; }
```

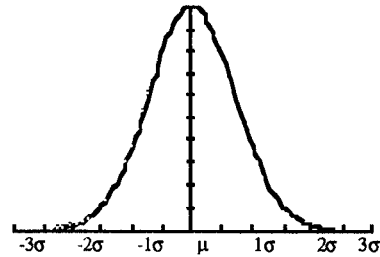


Gauss = Normal Distribution:

equation: $\phi(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad : x \in [0, \infty]$

deviation is ~.5% see [Ripley, pp. 54]

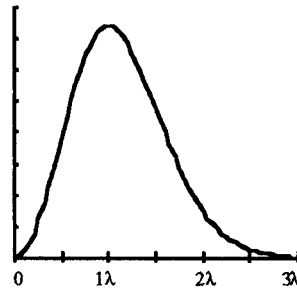
```
code: double gauss(double mu, sigma)
{ int i; double ans = 0.0;
  for (i = 0; i<12; i++)
    {ans = ans + (randf()) -0.5;}
  return (mu + ans/6);
}
```



Poisson

equation: $f(k) = e^{-\lambda} \frac{\lambda^k}{k!} \quad : k = 0, 1, \dots$

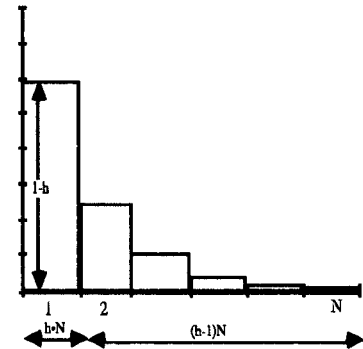
```
code: long poisson(long lambda)
{ long n = 0;
  double c = pow(e, -lambda);
  double p = 1.0;
  while (p >= c)
    {p = p * randf(); n++;}
  return (n-1);
};
```



Self-Similar (80-20 rule)

Integers between 1...N.
 The first h*N integers get 1-h of the distribution.
 For example: if N = 25 and h = .10, then
 80% of the weight goes to the first 5 integers.
 and 64% of the weight goes to the first integer.

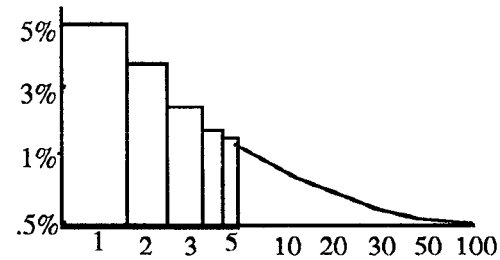
```
code: long selfsimilar(long n, double h)
{ return (1 + (long)
  (N * pow((randf(), log(h)/log(1.0-h))))
); };
```



Zipf's "Law"

Integers between 1...N.
 Integer k gets weight proportional to $(\frac{1}{k})^\theta$ theta
 where $0 < \theta < 1$ is the skew

```
code:
zeta(double n, theta)
{ int i, long ans = 0;
  for (i=1; i<=n; i++) ans += pow(1/n, theta);
  return(ans); }
long zipf(long n, double theta)
{
  double alpha = 1 / (1 - theta);
  double zetan = zeta(n, theta);
  double eta =
    (1 - pow(2.0 / n, 1 - theta)) /
    (1 - zeta(theta, 2) / zetan);
  double u = randf();
  double uz = u * zetan;
  if (uz < 1) return 1;
  if (uz < 1 + pow(0.5, theta)) return 2;
  return 1 +
    (long) (n * pow(eta*u - eta + 1, alpha));
};
```



The number of hot values that require special treatment depends on how accurately one needs to represent the distribution in question. One should bear in mind that the self-similar and Zipf distributions are themselves only approximations to what is observed in actual systems. As an example, consider a self-similar distribution following the 95/5 rule. A table of just 313 out of a billion possible values would account for over 77% of the weight of this distribution.

9. Summary

This paper showed how to convert a simple sequential load into a parallel load – turning a two-day task into a one-hour task. It then explored the ways to generate synthetic data. At first it focused on generating the primary keys of records and values uncorrelated to these keys: dense-unique-pseudo-random sequences. Then, attention turned to building indices on these synthetic tables – either by sorting, or by using discrete logarithms. By careful selection of generators, the discrete log problem is tractable and indices can be quickly generated within the 1-hour limit we set for the billion-record load.

The paper then looked at skewed distributions. It presented the standard ways to generate uniform, exponential, normal, and Poisson distributions. It went into more detail on the new topic of self-similar and Zipfian distributions.

Using these techniques, one can generate billion-record databases in an hour, and two terabyte databases per day.

10. Acknowledgments

This paper has been in-progress since 1987. Many people have contributed ideas to it. Some inspiration for the work came from Dina Bitton and Jeff Millman who showed us the data generators they used in their DBstar product. Dave DeWitt was an early user of these algorithms and has been a source of ideas and encouragement. Betty Salzberg explained the concept of discrete logs to Gray and Englert and put us in contact with that community. Geotz Graefe and Gayn Winters helped us improve the presentation.

11. References

- [Bitton 1] Bitton, D., DeWitt, D., Turbyfill, C., Source code for Wisconsin Database Generator distributed on the "Wisconsin Benchmark Tape", Computer Science, U. Wisconsin, Madison, WI, 1984
- [Bitton 2] Bitton, D., et al., "Benchmarking Database Systems: A Systematic Approach", 9th VLDB, Nov. 1983.
- [Coppersmith] Coppersmith, D., et. al., "Discrete Logarithms in GF(p)", *Algorithmica*, 1(1), 1986, pp. 1-15.
- [DeWitt 1] DeWitt, D. et. al., "Gamma - A High Performance Dataflow Database Machine", Proc. 12th VLDB, Sept. 1986, pp. 228-236.
- [DeWitt 2] DeWitt, D., et. al., "The Gamma Database Machine Project", IEEE Trans. on Knowledge and Data Engineering, March 1990.
- [DeWitt 3] DeWitt, D.J., Naughton, J.F., Schneider, D.A. "Parallel Sorting on a Shared-Nothing Architecture Using Probabilistic Splitting", Proc. First Int. Conf. on Parallel and Distributed Info. Systems, IEEE Press, Jan. 1992, pp. 280-291.
- [Englert] Englert, S., et. al. "A Benchmark of NonStop SQL Release 2 Demonstrating Near-Linear Speedup and Scaleup on Large Databases", Proc. of 1990 ACM SIGMETRICS Conference. May 1990.
- [Gerber] Gerber, R.H., "Dataflow Query Processing Using Multiprocessor Hash-Partitioned Algorithms", Ph.D. Thesis, Comp. Sci. TR 672, U. Wisconsin, Madison. Oct. 1986.
- [Iobbs] Hobbs, L., England, K., *Rdb/VS : A Comprehensive Guide*, Digital Press, Bedford, MA, 1991.
- [Horst] Horst, R., Chou, T., "The Hardware Architecture and Linear Expansion of Tandem NonStop Systems", Proc. 12th Int. Conf. Computer Architecture, June 1985.
- [Jain] Jain, R., *The Art of Computer Systems Performance Analysis*. John Wiley & Sons, New York, 1991
- [Kim] Kim, M.Y., "Synchronized Disk Interleaving", IEEE TOC, 3(11), Nov. 1986, pp. 978-988.
- [Knuth] Knuth, D., *The Art of Computer Programming*, V. 2, Chapter 3, 2nd ed., Addison Wesley, 1981.
- [Kronenberg] Kronenberg, N., H. Levey, W. Strecker and R. McCrewood. "The VAXcluster Concept; An Overview of a Distributed System." *Digital Technical Journal*. 1(3): Jan. 1987, pp. 7-21.
- [Nyberg] Nyberg, C., Barclay, T., Gray, J., Lomet, D., "AlphaSort - A High-Speed Sort for RISC Machines" Proc 1994 ACM SIGMOD, 1994.
- [Press] Press, W.H., Teukolsky, S.A., Vetterling, W.T., Flannery, B.P., *Numerical Recipes in C - 2'nd Edition*, Cambridge Univ. Press. Cambridge, 1992.
- [Ripley] Ripley, B.D., *Stochastic Simulation*, John Wiley, 1987.
- [Schrage] Schrage, L.E., "A More Portable FORTRAN Random Number Generator", *ACM TOMS*, V. 5(2), May 1979, pp 132-138.
- [Smith] Smith, M., et al., "An Experiment on Response Time Scalability", Proc. 6th Int. Workshop on Database Machines, June 1989
- [Stonebraker] Stonebraker, M., "The Case for Shared-Nothing", *Database Engineering*, V. 9(1), Jan. 1986.
- [Tanenbaum] Tanenbaum, A.S., *Operating Systems: Design and Implementation*, Prentice Hall, 1986.
- [Teradata] "The Genesis of a Database Computer: A Conversation with Jack Shemer and Phil Neches of Teradata Corporation", *IEEE Computer*, Nov. 1984. or *DBC/1012 Database Computer System Manual, Release 1.3*, C10-0001-01, Teradata Corp., Los Angeles, Feb. 1985.
- [Thekkath] Thekkath, C.A., Levey, H.M., Limits to Low-Latency Communication on High-Speed Networks, *ACM TOCS*, 11(2), May 1993, pp. 179-203.
- [TPC] "Transaction Processing Performance Council Benchmark A", Chapter 3 of *Performance Handbook for Database and Transaction Processing Systems*, Morgan Kaufmann, San Mateo, 1993.
- [Uren] Uren, S., "Message System Performance Tests", *Tandem Systems Review*, V3.4, pp. 27-32, Dec. 1986.