

Generating Realistic *Impressions* for File-System Benchmarking

NITIN AGRAWAL, ANDREA C. ARPACI-DUSSEAU, and
REMZI H. ARPACI-DUSSEAU
University of Wisconsin-Madison

The performance of file systems and related software depends on characteristics of the underlying file-system image (i.e., file-system metadata and file contents). Unfortunately, rather than benchmarking with realistic file-system images, most system designers and evaluators rely on *ad hoc* assumptions and (often inaccurate) rules of thumb. Furthermore, the lack of standardization and reproducibility makes file-system benchmarking ineffective. To remedy these problems, we develop Impressions, a framework to generate statistically accurate file-system images with realistic metadata and content. Impressions is flexible, supporting user-specified constraints on various file-system parameters using a number of statistical techniques to generate consistent images. In this article, we present the design, implementation, and evaluation of Impressions and demonstrate its utility using desktop search as a case study. We believe Impressions will prove to be useful to system developers and users alike.

Categories and Subject Descriptors: D.4.3 [**Operating Systems**]: File Systems Management; D.4.8 [**Operating Systems**]: Performance

General Terms: Measurement, Performance

Additional Key Words and Phrases: File and storage system benchmarking

ACM Reference Format:

Agrawal, N., Arpaci-Dusseau, A. C., and Arpaci-Dusseau, R. H. 2009. Generating realistic *Impressions* for file-system benchmarking. *ACM Trans. Storage* 5, 4, Article 16 (December 2009), 30 pages. DOI = 10.1145/1629080.1629086 <http://doi.acm.org/10.1145/1629080.1629086>

An earlier version of this article appeared in the *Proceedings of the 7th USENIX Conference on File and Storage Technologies (FAST'09)*.

This material is based on work supported by the National Science Foundation under the following grants: CCF-0621487, CNS-0509474, as well as by generous donations from Network Appliance and Sun Microsystems. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of NSF or other institutions.

Author's address: email: {nitina,dusseau,remzi}@cs.wisc.edu.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org. © 2009 ACM 1553-3077/2009/12-ART16 \$10.00

DOI 10.1145/1629080.1629086 <http://doi.acm.org/10.1145/1629080.1629086>

ACM Transactions on Storage, Vol. 5, No. 4, Article 16, Publication date: December 2009.

1. INTRODUCTION

File system benchmarking is in a state of disarray. In spite of tremendous advances in file system design, the approaches for benchmarking still lag far behind. The goal of benchmarking is to understand how the system under evaluation will perform under real-world conditions and how it compares to other systems; however, recreating real-world conditions for the purposes of benchmarking file systems has proven challenging. The two main challenges in achieving this goal are generating representative *workloads*, and creating realistic *file-system state*.

While creating representative workloads is not an entirely solved problem, significant steps have been taken towards this goal. Empirical studies of file-system access patterns [Baker et al. 1991; Gribble et al. 1998; Ousterhout et al. 1985] and file-system activity traces [Riedel et al. 2002; SNIA 2007] have led to work on synthetic workload generators [Anderson and Chase 2002; Ebling and Satyanarayanan 1994] and methods for trace replay [Anderson et al. 2004; Mesnier et al. 2007].

The second, and perhaps more difficult, challenge is to re-create the file-system *state* such that it is representative of the target usage scenario. Several factors contribute to file-system state, important among them are the *in-memory* state (contents of the buffer cache), the *on-disk* state (disk layout and fragmentation) and the characteristics of the *file-system image* (files and directories belonging to the namespace and file contents).

One well understood contributor to state is the *in-memory* state of the file system. Previous work has shown that the contents of the cache can have significant impact on the performance results [Dahlin et al. 1994]. Therefore, system initialization during benchmarking typically consists of a cache “warm-up” phase wherein the workload is run for some time prior to the actual measurement phase. Another important factor is the *on-disk* state of the file system, or the degree of *fragmentation*; it is a measure of how the disk blocks belonging to the file system are laid out on disk. Previous work has shown that fragmentation can adversely affect performance of a file system [Smith and Seltzer 1997]. Thus, prior to benchmarking, a file system should undergo *aging* by replaying a workload similar to that experienced by a real file system over a period of time [Smith and Seltzer 1997].

Surprisingly, one key contributor to file-system state has been largely ignored—the characteristics of the *file-system image*. The properties of file-system metadata and the actual content within the files are key contributors to file-system state, and can have a significant impact on the performance of a system. Properties of file-system metadata include information on how directories are organized in the file-system namespace, how files are organized into directories, and the distributions for various file attributes such as size, depth, and extension type. Consider a simple example: the time taken for a find operation to traverse a file system while searching for a file name depends on a number of attributes of the file-system image, including the depth of the file-system tree and the total number of files. Similarly, the time taken for a grep operation to search for a keyword also depends on the type of files (i.e., binary vs. others) and the file content.

File-system benchmarking frequently requires this sort of information on file systems, much of which is available in the form of empirical studies of file-system contents [Agrawal et al. 2007; Douceur and Bolosky 1999; Irlam 1993; Mullender and Tanenbaum 1984; Satyanarayanan 1981; Sienknecht et al. 1994]. These studies focus on measuring and modeling different aspects of file-system metadata by collecting snapshots of file-system images from real machines. The studies range from a few machines to tens of thousands of machines across different operating systems and usage environments. Collecting and analyzing this data provides useful information on how file systems are used in real operating conditions.

In spite of the wealth of information available in file-system studies, system designers and evaluators continue to rely on *ad hoc* assumptions and often inaccurate rules of thumb. Table I presents evidence to confirm this hypothesis; it contains a (partial) list of publications from top-tier systems conferences in the last ten years that required a test file-system image for evaluation. We present both the description of the file-system image provided in the paper and the intended goal of the evaluation.

In the table, there are several examples where a new file system or application design is evaluated on the evaluator's personal file system without describing its properties in sufficient detail for it to be reproduced [Cipar et al. 2007; Hutchinson et al. 1999; Prabhakaran et al. 2005]. In others, the description is limited to coarse-grained measures such as the total file-system size and the number of files, even though other file-system attributes (e.g., tree depth) are relevant to measuring performance or storage space overheads [Cox et al. 2002; Cox and Noble 2003; Gopal and Manber 1999; Muthitacharoen et al. 2001]. File systems are also sometimes generated with parameters chosen randomly [Storer et al. 2008; Zhang and Ghose 2003], or chosen without explanation of the significance of the values [Fu et al. 2002; Padioleau and Ridoux 2003; Sobti et al. 2004]. Occasionally, the parameters are specified in greater detail [Rowstron and Druschel 2001], but not enough to recreate the original file system.

The important lesson to be learned here is that there is no standard technique to systematically include information on file-system images for experimentation. For this reason, we find that more often than not, the choices made are arbitrary, suited for ease-of-use more than accuracy and completeness. Furthermore, the lack of standardization and reproducibility of these choices makes it near-impossible to compare results with other systems.

To address these problems and improve one important aspect of file system benchmarking, we develop *Impressions*, a framework to generate representative and statistically accurate file-system images. *Impressions* gives the user flexibility to specify one or more parameters from a detailed list of file system parameters (file-system size, number of files, distribution of file sizes, etc.). *Impressions* incorporates statistical techniques (automatic curve-fitting, resolving multiple constraints, interpolation and extrapolation, etc.) and uses statistical tests for goodness-of-fit to ensure the accuracy of the image.

We believe *Impressions* will be of great use to system designers, evaluators, and users alike. A casual user looking to create a representative file-system

Table I. Choice of File System Parameters in Prior Research

Paper	Description	Used to Measure
HAC [Gopal and Manber 1999]	File system with 17000 files totaling 150 MB	Time and space needed to create a Glimpse index
IRON [Prabhakaran et al. 2005]	None provided	Checksum and metadata replication overhead; parity block overhead for user files
LBFS [Muthitacharoen et al. 2001]	10702 files from /usr/local, total size 354 MB	Performance of LBFS chunking algorithm
LISFS [Padioleau and Ridoux 2003]	633 MP3 files, 860 program files, 11502 man pages	Disk space overhead; performance of search-like activities: UNIX find and LISFS lookup
PAST [Rowstron and Druschel 2001]	2 million files, mean size 86 KB, median 4 KB, largest file size 2.7 GB, smallest 0 Bytes, total size 166.6 GB	File insertion, global storage utilization in a P2P system
Pastiche [Cox et al. 2002]	File system with 1641 files, 109 dirs, 13.4 MB total size	Performance of backup and restore utilities
Pergamum [Storer et al. 2008]	Randomly generated files of “several” megabytes	Data transfer performance
Samsara [Cox and Noble 2003]	File system with 1676 files and 13 MB total size	Data transfer and querying performance, load during querying
Segank [Sobti et al. 2004]	5-deep directory tree, 5 subdirs, and 10 8 KB files per directory	Performance of Segank: volume update, creation of read-only snapshot, read from new snapshot
SFS read-only [Fu et al. 2002]	1000 files distributed evenly across 10 directories and contain random data	Single client/single server read performance
TFS [Cipar et al. 2007]	Files taken from /usr to get “realistic” mix of file sizes	Performance with varying contribution of space from local file systems
WAFL backup [Hutchinson et al. 1999]	188 GB and 129 GB volumes taken from the Engineering department	Performance of physical and logical backup, and recovery strategies
yFS [Zhang and Ghose 2003]	Avg. file size 16 KB, avg. number of files per directory 64, random file names	Performance under various benchmarks (file creation, deletion)

image without worrying about carefully selecting parameters can simply run Impressions with its default settings; Impressions will use prespecified distributions from file-system studies to create a representative image. A more sophisticated user has the power to individually control the knobs for a comprehensive set of file-system parameters; Impressions will carefully work out the statistical details to produce a consistent and accurate image. In both cases, Impressions ensures complete reproducibility of the image, by reporting the used distributions, parameter values, and seeds for random number generators.

In this article we present the design, implementation and evaluation of the Impressions framework (Section 3), which we have made publicly available.

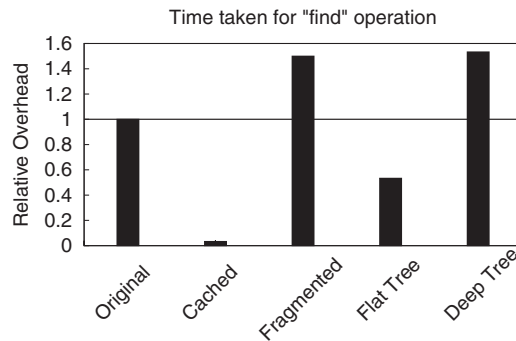


Fig. 1. Impact of directory tree structure. Shows impact of tree depth on time taken by *find*. The file systems are created by Impressions using default distributions (Table II). To exclude effects of the on-disk layout, we ensure a perfect disk layout (layout score 1.0) for all cases except the one with fragmentation (layout score 0.95). The *flat tree* contains all 100 directories at depth 1; the *deep tree* has directories successively nested to create a tree of depth 100.

Impressions is built with the following design goals:

- accuracy* in generating various statistical constructs to ensure a high degree of statistical rigor;
- flexibility* in allowing users to specify a number of file-system distributions and constraints on parameter values, or in choosing default values;
- representativeness* by incorporating known distributions from file-system studies.
- ease of use* by providing a simple, yet powerful, command-line interface.

Using desktop search as a case study, we demonstrate the usefulness and ease of use of Impressions in quantifying application performance, and in finding application policies and bugs (Section 4). To bring the paper to a close, we discuss related work (Section 6), and finally conclude (Section 7).

2. EXTENDED MOTIVATION

We begin this section by asking a basic question: does file-system structure really matter? We then describe the goals for generating realistic file-system images and discuss existing approaches to do so.

2.1 Does File-System Structure Matter?

Structure and organization of file-system metadata matters for workload performance. Let us take a look at the simple example of a frequently used UNIX utility: *find*. Figure 1 shows the relative time taken to run “*find /*” searching for a file name on a test file system as we vary some parameters of file-system state.

The first bar represents the time taken for the run on the original test file system. Subsequent bars are normalized to this time and show performance for a run with the file-system contents in buffer cache, a fragmented version of the same file system, a file system created by flattening the original directory

tree, and finally one by deepening the original directory tree. The graph echoes our understanding of caching and fragmentation, and brings out one aspect that is often overlooked: structure really matters. From this graph we can see that even for a simple workload, the impact of tree depth on performance can be as large as that with fragmentation, and varying tree depths can have significant performance variations (300% between the flat and deep trees in this example).

Assumptions about file-system structure have often trickled into file system design, but no means exist to incorporate the effects of realistic file-system images in a systematic fashion. As a community, we well understand that caching matters, and have begun to pay attention to fragmentation, but when it comes to file-system structure, our approach is surprisingly *laissez faire*.

2.2 Goals for Generating FS Images

We believe that the file-system image used for an evaluation should be *realistic* with respect to the workload; the image should contain a sufficient degree of *detail* to realistically exercise the workload under consideration. An increasing degree of detail will likely require more effort and slow down the process. Thus it is useful to know the degree sufficient for a given evaluation. For example, if the performance of an application simply depends on the size of files in the file system, the chosen file-system image should reflect that. On the other hand, if the performance is also sensitive to the fraction of binary files amongst all files (e.g., to evaluate desktop search indexing), then the file-system image also needs to contain realistic distributions of file extensions.

We walk through some examples that illustrate the different degrees of detail needed in file-system images.

—At one extreme, a system could be completely oblivious to both metadata and content. An example of such a system is a mirroring scheme (RAID-1 [Patterson et al. 1988]) underneath a file system, or a backup utility taking whole-disk backups. The performance of such schemes depends solely on the block traffic.

Alternately, systems could depend on the attributes of the file-system image with different degrees of detail:

- The performance of a system can depend on the amount of file data (number of files and directories, or the size of files and directories, or both) in any given file system (e.g., a backup utility taking whole file-system snapshots).
- Systems can depend on the structure of the file system namespace and how files are organized in it (e.g., a version control system for a source-code repository).
- Finally, many systems also depend on the actual data stored within the files (e.g., a desktop search engine for a file system, or a spell-checker).

Impressions is designed with this goal of flexibility from the outset. The user is given complete control of a number of file-system parameters, and is provided

with an easy to use interface. Transparently, Impressions seamlessly ensures accuracy and representativeness.

2.3 Existing Approaches

One alternate approach to generating realistic file-system images is to randomly select a set of actual images from a corpus, an approach popular in other fields of computer science such as Information Retrieval, Machine Learning, and Natural Language Processing [NIST 2007]. In the case of file systems the corpus would consist of a set of known file-system images. This approach arguably has several limitations which make it difficult and unsuitable for file systems research. First, there are too many parameters required to accurately describe a file-system image that need to be captured in a corpus. Second, without precise control in varying these parameters according to experimental needs, the evaluation can be blind to the actual performance dependencies. Finally, the cost of maintaining and sharing any realistic corpus of file-system images would be prohibitive. The size of the corpus itself would severely restrict its usefulness especially as file systems continue to grow larger.

Unfortunately, these limitations have not deterred researchers from using their personal file systems as a (trivial) substitute for a file-system corpus.

3. THE IMPRESSIONS FRAMEWORK

In this section we describe the design, implementation and evaluation of Impressions: a framework for generating file-system images with realistic and statistically accurate metadata and content. Impressions is flexible enough to create file-system images with varying configurations, guaranteeing the accuracy of images by incorporating a number of statistical tests and techniques.

We first present a summary of the different modes of operation of Impressions, and then describe the individual statistical constructs in greater detail. Wherever applicable, we evaluate their accuracy and performance.

3.1 Modes of Operation

A system evaluator can use Impressions in different modes of operation, with varying degree of user input.

Sometimes, an evaluator just wants to create a representative file-system image without worrying about the need to carefully select parameters. Hence, in the *automated* mode, Impressions is capable of generating a file-system image with minimal input required from the user (e.g., the size of the desired file-system image), relying on default settings of known empirical distributions to generate representative file-system images. We refer to these distributions as *original* distributions.

At other times, users want more control over the images, for example, to analyze the sensitivity of performance to a given file-system parameter, or to describe a completely different file-system usage scenario. Hence, Impressions supports a *user-specified* mode, where a more sophisticated user has the power

Table II. Parameters and Default Values in Impressions
List of distributions and their parameter values used in the default mode.

Parameter	Default Model & Parameters
Directory count w/ depth	Generative model
Directory size (subdirs)	Generative model
File size by count	Lognormal-body ($\alpha_1 = 0.99994, \mu = 9.48, \sigma = 2.46$) Pareto-tail ($k = 0.91, \lambda_m = 512\text{MB}$)
File size by containing bytes	Mixture-of-lognormals ($\alpha_1 = 0.76, \mu_1 = 14.83, \sigma_1 = 2.35$ $\alpha_2 = 0.24, \mu_2 = 20.93, \sigma_2 = 1.48$)
Extension popularity	Percentile values
File count w/ depth	Poisson ($\lambda = 6.49$)
Bytes with depth	Mean file size values
Directory size (files)	Inverse-polynomial (degree = 2, offset = 2.36)
File count w/ depth (w/ special directories)	Conditional probabilities (biases for special dirs)
Degree of Fragmentation	Layout score (1.0) or Pre-specified workload

to individually control the knobs for a comprehensive set of file-system parameters; we refer to these as user-specified distributions. Impressions carefully works out the statistical details to produce a consistent and accurate image.

In both the cases, Impressions ensures complete reproducibility of the file-system image by reporting the used distributions, their parameter values, and seeds for random number generators.

Impressions can use any dataset or set of parameterized curves for the *original* distributions, leveraging a large body of research on analyzing file-system properties [Agrawal et al. 2007; Douceur and Bolosky 1999; Irlam 1993; Mullender and Tanenbaum 1984; Satyanarayanan 1981; Sienknecht et al. 1994]. For illustration, in this article we use a recent static file-system snapshot dataset made publicly available [Agrawal et al. 2007]. The snapshots of file-system metadata were collected over a five-year period representing over 60,000 Windows PC file systems in a large corporation. These snapshots were used to study distributions and temporal changes in file size, file age, file-type frequency, directory size, namespace structure, file-system population, storage capacity, and degree of file modification. The study also proposed a generative model explaining the creation of file-system namespaces.

Impressions provides a comprehensive set of individually controllable file system parameters. Table II lists these parameters along with their default selections. For example, a user may specify the size of the file-system image, the number of files in the file system, and the distribution of file sizes, while selecting default settings for all other distributions. In this case, Impressions will ensure that the resulting file-system image adheres to the default distributions while maintaining the user-specified invariants.

3.2 Basic Techniques

The goal of Impressions is to generate realistic file-system images, giving the user complete flexibility and control to decide the extent of accuracy and detail. To achieve this, Impressions relies on a number of statistical techniques.

In the simplest case, Impressions needs to create statistically accurate file-system images with default distributions. Hence, a basic functionality required by Impressions is to convert the parameterized distributions into real sample values used to create an instance of a file-system image. Impressions uses random sampling to take a number of independent observations from the respective probability distributions. Wherever applicable, such parameterized distributions provide a highly compact and easy-to-reproduce representation of observed distributions. For cases where standard probability distributions are infeasible, a Monte Carlo method is used.

A user may want to use file system datasets other than the default choice. To enable this, Impressions provides automatic curve-fitting of empirical data.

Impressions also provides the user with the flexibility to specify distributions and constraints on parameter values. One challenge thus is to ensure that multiple constraints specified by the user are resolved consistently. This requires statistical techniques to ensure that the generated file-system images are accurate with respect to both the user-specified constraints and the default distributions.

In addition, the user may want to explore values of file system parameters, not captured in any dataset. For this purpose, Impressions provides support for interpolation and extrapolation of new curves from existing datasets.

Finally, to ensure the accuracy of the generated image, Impressions contains a number of built-in statistical tests, for goodness-of-fit (e.g., Kolmogorov-Smirnov, Chi-Square, and Anderson-Darling), and to estimate error (e.g., Confidence Intervals, MDCC, and Standard Error). Where applicable, these tests ensure that all curve-fit approximations and internal statistical transformations adhere to the highest degree of statistical rigor desired.

3.3 Creating Valid Metadata

The simplest use of Impressions is to generate file-system images with realistic metadata. This process is performed in two phases: first, the skeletal file-system namespace is created; and second, the namespace is populated with files conforming to a number of file and directory distributions.

3.3.1 Creating File-System Namespace. The first phase in creating a file system is to create the namespace structure or the *directory tree*. We assume that the user specifies the size of the file-system image. The count of files and directories is then selected based on the file system size (if not specified by the user). Depending on the degree of detail desired by the user, each file or directory attribute is selected step by step until all attributes have been assigned values. We now describe this process assuming the highest degree of detail.

To create directory trees, Impressions uses the generative model proposed by Agrawal et al. [2007] to perform a Monte Carlo simulation. According to

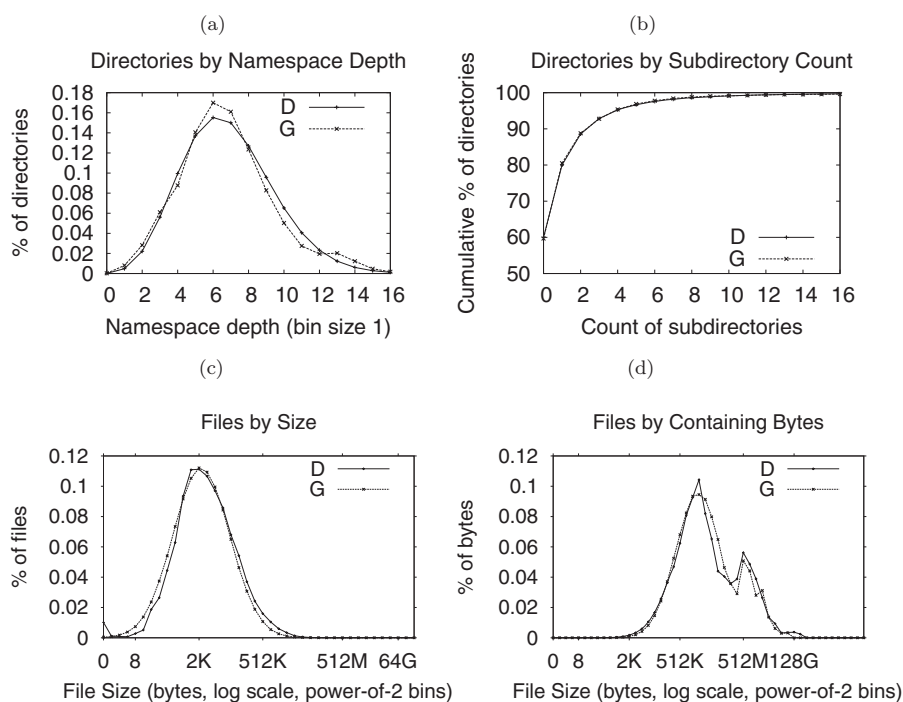


Fig. 2. Accuracy of Impressions in recreating file system properties. Shows the accuracy of the entire set of file system distributions modeled by Impressions. D: the desired distribution; G: the generated distribution. Impressions is quite accurate in creating realistic file system state for all parameters of interest shown here. We include a special abscissa for the zero value on graphs having a logarithmic scale.

this model, new directories are added to a file system one at a time, and the probability of choosing each extant directory as a parent is proportional to $C(d) + 2$, where $C(d)$ is the count of extant subdirectories of directory d . The model explains the creation of the file system namespace, accounting both for the size and count of directories by depth, and the size of parent directories. The input to this model is the total number of directories in the file system. Directory names are generated using a simple iterative counter.

To ensure the accuracy of generated images, we compare the generated distributions (i.e., created using the parameters listed in Table II), with the desired distributions (i.e., ones obtained from the dataset discussed previously in Section 3.1). Figures 2 and 3 shows in detail the accuracy for each step in the namespace and file creation process. For almost all the graphs, the y-axis represents the percentage of files, directories, or bytes belonging to the categories or bins shown on the x-axis, as the case may be.

Figures 2(a) and 2(b) show the distribution of directories by depth, and directories by subdirectory count, respectively. The y-axis in this case is the percentage of directories at each level of depth in the namespace, shown on the x-axis. The two curves representing the generated and the desired distributions match

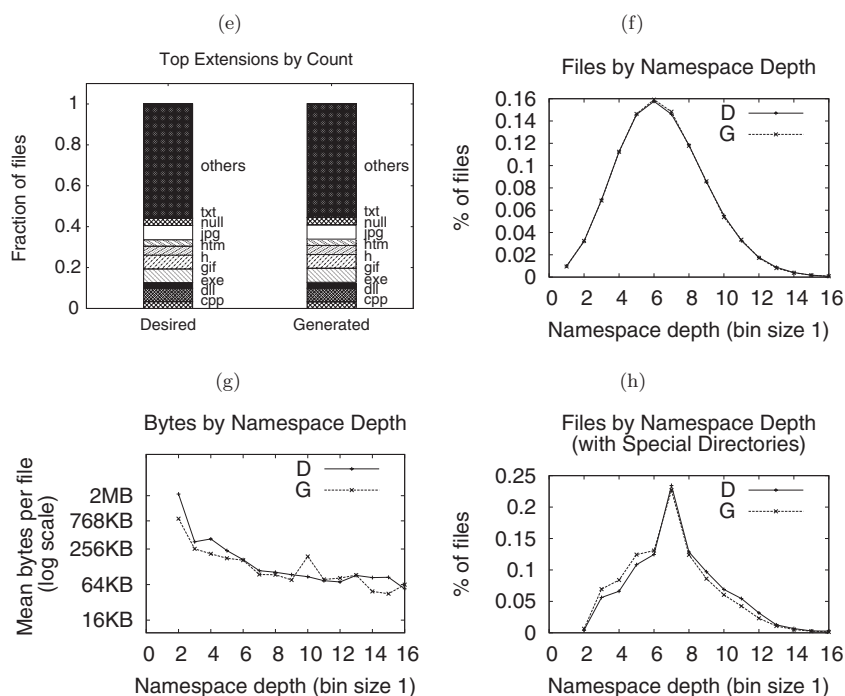


Fig. 3. Accuracy of Impressions in recreating file system properties. Shows the accuracy of the entire set of file system distributions modeled by Impressions. D: the desired distribution; G: the generated distribution. Impressions is quite accurate in creating realistic file system state for all parameters of interest shown here. We include a special abscissa for the zero value on graphs having a logarithmic scale.

quite well, indicating good accuracy and reaffirming prior results [Agrawal et al. 2007].

3.3.2 Creating Files. The next phase is to populate the directory tree with files. Impressions spends most of the total runtime and effort during this phase, as the bulk of its statistical machinery is exercised in creating files. Each file has a number of attributes such as its size, depth in the directory tree, parent directory, and file extension. Similarly, the choice of the parent directory is governed by directory attributes such as the count of contained subdirectories, the count of contained files, and the depth of the parent directory. Analytical approximations for file system distributions proposed previously [Douceur and Bolosky 1999] guided our own models.

First, for each file, the size of the file is sampled from a hybrid distribution describing file sizes. The body of this hybrid curve is approximated by a lognormal distribution, with a Pareto tail distribution ($k = 0.91$, $\lambda_m = 512\text{MB}$) accounting for the heavy tail of files with size greater than 512 MB. The exact parameter values used for these distributions are listed in Table II. These parameters were obtained by fitting the respective curves to file sizes obtained from the file-system dataset previously discussed (Section 3.1). Figure 2(c) shows the

accuracy of generating the distribution of files by size. We initially used a simpler model for file sizes represented solely by a lognormal distribution. While the results were acceptable for files by size (Figure 2(c)), the simpler model failed to account for the distribution of bytes by containing file size; coming up with a model to accurately capture the bimodal distribution of bytes proved harder than we had anticipated. Figure 2(d) shows the accuracy of the hybrid model in Impressions in generating the distribution of bytes. The pronounced double mode observed in the distribution of bytes is a result of the presence of a few large files; an important detail that is otherwise missed if the heavy-tail of file sizes is not accurately accounted for.

Once the file size is selected, we assign the file name and extension. Impressions keeps a list of percentile values for popular file extensions (i.e., top 20 extensions by count, and by bytes). These extensions together account for roughly 50% of files and bytes in a file system ensuring adequate coverage for the important extensions. The remainder of files are given randomly generated three-character extensions. Currently filenames are generated by a simple numeric counter incremented on each file creation. Figure 3(e) shows the accuracy of Impressions in creating files with popular extensions by count.

Next, we assign file depth d , which requires satisfying two criteria: the distribution of files with depth, and the distribution of bytes with depth. The former is modeled by a Poisson distribution, and the latter is represented by the mean file sizes at a given depth. Impressions uses a multiplicative model combining the two criteria, to produce appropriate file depths. Figures 3(f) and 3(g) show the accuracy in generating the distribution of files by depth, and the distribution of bytes by depth, respectively.

The final step is to select a parent directory for the file, located at depth $d - 1$, according to the distribution of directories with file count, modeled using an inverse-polynomial of degree 2. As an added feature, Impressions supports the notion of “Special” directories containing a disproportionate number of files or bytes (e.g., “Program Files” folder in the Windows environment). If required, during the selection of the parent directory, a selection bias is given to these special directories. Figure 3(h) shows the accuracy in supporting special directories with an example of a *typical* Windows file system having files in the web cache at depth 7, in Windows and Program Files folders at depth 2, and System files at depth 3.

Table III shows the average difference between the generated and desired images from Figure 3 for 20 trials. The difference is measured in terms of the MDCC (Maximum Displacement of the Cumulative Curves). For instance, an MDCC value of 0.03 for directories with depth, implies a *maximum* difference of 3% on an average, between the desired and the generated cumulative distributions. Overall, we find that the models created and used by Impressions for representing various file-system parameters produce fairly accurate distributions in all the above cases. While we have demonstrated the accuracy of Impressions for the Windows dataset, there is no fundamental restriction limiting it to this dataset. We believe that with little effort, the same level of accuracy can be achieved for any other dataset.

Table III. Statistical Accuracy of Generated Images
Shows average accuracy of generated file-system images in terms of the MDCC (Maximum Displacement of the Cumulative Curves) representing the maximum difference between cumulative curves of generated and desired distributions. Averages are shown for 20 trials. (*) For bytes with depth, MDCC is not an appropriate metric, we instead report the average difference in mean bytes per file (MB). The numbers correspond to the set of graphs shown in Figure 3 and reflect fairly accurate images.

Parameter	MDCC
Directory count with depth	0.03
Directory size (subdirectories)	0.004
File size by count	0.04
File size by containing bytes	0.02
Extension popularity	0.03
File count with depth	0.05
Bytes with depth	0.12 MB*
File count w/ depth w/ special dirs	0.06

3.4 Resolving Arbitrary Constraints

One of the primary requirements for Impressions is to allow flexibility in specifying file system parameters without compromising accuracy. This means that users are allowed to specify somewhat arbitrary constraints on these parameters, and it is the task of Impressions to resolve them. One example of such a set of constraints would be to specify a large number of files for a small file system, or vice versa, given a file size distribution. Impressions will try to come up with a sample of file sizes that best approximates the desired distribution, while still maintaining the invariants supplied by the user, namely the number of files in the file system and the sum of all file sizes being equal to the file system used space.

Multiple constraints can also be implicit (i.e., arise even in the absence of user-specified distributions). Due to random sampling, different sample sets of the same distribution are not guaranteed to produce exactly the same result, and consequently, the sum of the elements can also differ across samples. Consider the previous example of file sizes again: the sum of all file sizes drawn from a given distribution need not add up to the desired file system size (total used space) each time. More formally, this example is represented by the following set of constraints:

$$\begin{aligned} \mathcal{N} &= \{Constant_1 \vee x : x \in \mathcal{D}_1(x)\} \\ \mathcal{S} &= \{Constant_2 \vee x : x \in \mathcal{D}_2(x)\} \\ \mathcal{F} &= \{x : x \in \mathcal{D}_3(x; \mu, \sigma)\}; \left| \sum_{i=0}^{\mathcal{N}} \mathcal{F}_i - \mathcal{S} \right| \leq \beta * \mathcal{S}, \end{aligned}$$

where \mathcal{N} is the number of files in the file system; \mathcal{S} is the desired file system used space; \mathcal{F} is the set of file sizes; and β is the maximum relative error allowed. The first two constraints specify that \mathcal{N} and \mathcal{S} can be user specified constants

or sampled from their corresponding distributions \mathcal{D}_1 and \mathcal{D}_2 . Similarly, \mathcal{F} is sampled from the file size distribution \mathcal{D}_3 . These attributes are further subject to the constraint that the sum of all file sizes differs from the desired file system size by no more than the allowed error tolerance, specified by the user. To solve this problem, we use the following two techniques.

- If the initial sample does not produce a result satisfying all the constraints, we *oversample* additional values of \mathcal{F} from \mathcal{D}_3 , one at a time, until a solution is found, or the oversampling factor α/\mathcal{N} reaches λ (the maximum oversampling factor). α is the count of extra samples drawn from \mathcal{D}_3 . Upon reaching λ without finding a solution, we discard the current sample set and start over.
- The number of elements in \mathcal{F} during the oversampling stage is $\mathcal{N} + \alpha$. For every oversampling, we need to find if there exists \mathcal{F}_{Sub} , a subset of \mathcal{F} with \mathcal{N} elements, such that the sum of all elements of \mathcal{F}_{Sub} (file sizes) differs from the desired file system size by no more than the allowed error. More formally stated, we find if:

$$\exists \mathcal{F}_{Sub} = \{ \mathcal{X} : \mathcal{X} \subseteq \mathbb{P}(\mathcal{F}), |\mathcal{X}| = \mathcal{N}, |\mathcal{F}| = \mathcal{N} + \alpha, \\ | \sum_{i=0}^{\mathcal{N}} \mathcal{X}_i - \mathcal{S} | \leq \beta * \mathcal{S}, \alpha \in \mathbb{N} \wedge \frac{\alpha}{\mathcal{N}} \leq \lambda \}$$

The problem of resolving multiple constraints as formulated above, is a variant of the more general “Subset Sum Problem” which is NP-complete [Cormen et al. 2001]. Our solution is thus an approximation algorithm based on an existing $O(n \log n)$ solution [Przydatek 2002] for the Subset Sum Problem.

The existing algorithm has two phases. The first phase randomly chooses a solution vector which is valid (the sum of elements is less than the desired sum), and maximal (adding any element not already in the solution vector will cause the sum to exceed the desired sum). The second phase performs *local improvement*: for each element in the solution, it searches for the largest element not in the current solution which, if replaced with the current element, would reduce the difference between the desired and current sums. The solution vector is updated if such an element is found, and the algorithm proceeds with the next element, until all elements are compared.

Our problem definition and the modified algorithm differ from the original in the following ways.

- First, in the original problem, there is no restriction on the number of elements in the solution subset \mathcal{F}_{Sub} . In our case, \mathcal{F}_{Sub} can have exactly \mathcal{N} elements. We modify the first phase of the algorithm to set the initial \mathcal{F}_{Sub} as the first random permutation of \mathcal{N} elements selected from \mathcal{F} such that their sum is less than \mathcal{S} .
- Second, the original algorithm either finds a solution or terminates without success. We use an increasing sample size after each oversampling to reduce the error, and allow the solution to converge.
- Third, it is not sufficient for the elements in \mathcal{F}_{Sub} to have a numerical sum close to the desired sum \mathcal{S} , but the distribution of the elements must also be

close to the original distribution in \mathcal{F} . A goodness-of-fit test at the end of each oversampling step enforces this requirement. For our example, this ensures that the set of file sizes generated after resolving multiple constraints still follow the original distribution of file sizes.

The algorithm terminates successfully when the difference between the sums, and between the distributions, falls below the desired error levels. The success of the algorithm depends on the choice of the desired sum, and the *expected* sum (the sum due to the choice of parameters, e.g., μ and σ); the farther the desired sum is from the expected sum, the lesser are the chances of success.

Consider an example where a user has specified a desired file system size of 90000 bytes, a lognormal file size distribution ($\mu = 8.16$, $\sigma = 2.46$), and 1000 files. Figure 4(a) shows the convergence of the sum of file sizes in a sample set obtained with this distribution. Each line in the graph represents an independent trial, starting at a y-axis value equal to the sum of its initially sampled file sizes. Note that in this example, the initial sum differs from the desired sum by more than a 100% in several cases. The x-axis represents the number of extra iterations (*oversamples*) performed by the algorithm. For a trial to succeed, the sum of file sizes in the sample must converge to within 5% of the desired file system size. We find that in most cases λ ranges between 0 and 0.1 (i.e., less than 10% oversampling); and in almost all cases, $\lambda \leq 1$.

The distribution of file sizes in \mathcal{F}_{Sub} must be close to the original distribution in \mathcal{F} . Figure 4(b) and 4(c) show the difference between the original and constrained distributions for file sizes (for files by size, and files by bytes), for one successful trial from Figure 4(a). We choose these particular distributions as examples throughout this paper for two reasons. First, file size is an important parameter, so we want to be particularly thorough in its accuracy. Second, getting an accurate shape for the bimodal curve of files by bytes presents a challenge for Impressions; once we get our techniques to work for this curve, we are fairly confident of its accuracy on simpler distributions.

We find that Impressions resolves multiple constraints to satisfy the requirement on the sum, while respecting the original distributions. Table IV gives the summary for the above example of file sizes for different values of the desired file system size. The expected sum of 1000 file sizes, sampled as specified in the table, is close to 60000. Impressions successfully converges the initial sample set to the desired sum with an average oversampling rate α less than 5%. The average difference between the desired and achieved sum β is close to 3%. The constrained distribution passes the two-sample K-S test at the 0.05 significance level, with the difference between the two distributions being fairly small (the D statistic of the K-S test is around 0.03, which represents the maximum difference between two empirical cumulative distributions).

We repeat the previous experiment for two more choices of file system sizes, one lower than the expected mean (30K), and one higher (90K); we find that even when the desired sum is quite different from the expected sum, our algorithm performs well. Only for 2 of the 20 trials in the 90K case, did the algorithm fail to converge. For these extreme cases, we drop the initial sample and start over.

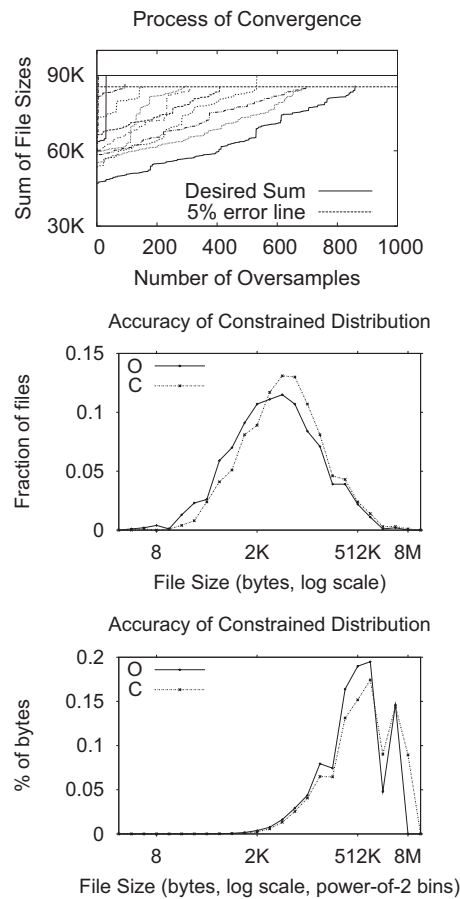


Fig. 4. Resolving multiple constraints. (a) Shows the process of convergence of a set of 1000 file sizes to the desired file system size of 90000 bytes. Each line represents an individual trial. A successful trial is one that converges to the 5% error line in less than 1000 oversamples. (b) Shows the difference between the original distribution of files by size, and the constrained distribution after resolution of multiple constraints in (a). O: Original; C: Constrained. (c) Same as (b), but for distribution of files by bytes instead.

3.5 Interpolation and Extrapolation

Impressions requires knowledge of the distribution of file system parameters necessary to create a valid image. While it is tempting to imagine that Impressions has perfect knowledge about the nature of these distributions for all possible values and combinations of individual parameters, it is often impossible.

First, the empirical data is limited to what is observed in any given dataset and may not cover the entire range of possible values for all parameters. Second, even with an exhaustive dataset, the user may want to explore regions of parameter values for which no data point exists, especially for “what if” style of analysis. Third, from an implementation perspective, it is more efficient to maintain compact representations of distributions for a few sample points, instead of large sets of data. Finally, if the empirical data is statistically

Table IV. Summary of Resolving Multiple Constraints

Shows average rate and accuracy of convergence after resolving multiple constraints for different values of desired file system size generated with a lognormal file size distribution \mathcal{D}_3 ($\mu = 8.16$, $\sigma = 2.46$). β : % error between the desired and generated sum, α : % of oversamples required, D is the test statistic for the K-S test representing the maximum difference between generated and desired empirical cumulative distributions. Averages are for 20 trials. Success is the number of trials having final $\beta \leq 5\%$, and D passing the K-S test.

Num. Files \mathcal{N}	Sizes Sum S (bytes)	Avg. β Initial	Avg. β Final	Avg. α	Avg. D Count	Avg. D Bytes	Success
1000	30000	21.55%	2.04%	5.74%	0.043	0.050	100%
1000	60000	20.01%	3.11%	4.89%	0.032	0.033	100%
1000	90000	34.35%	4.00%	41.2%	0.067	0.084	90%

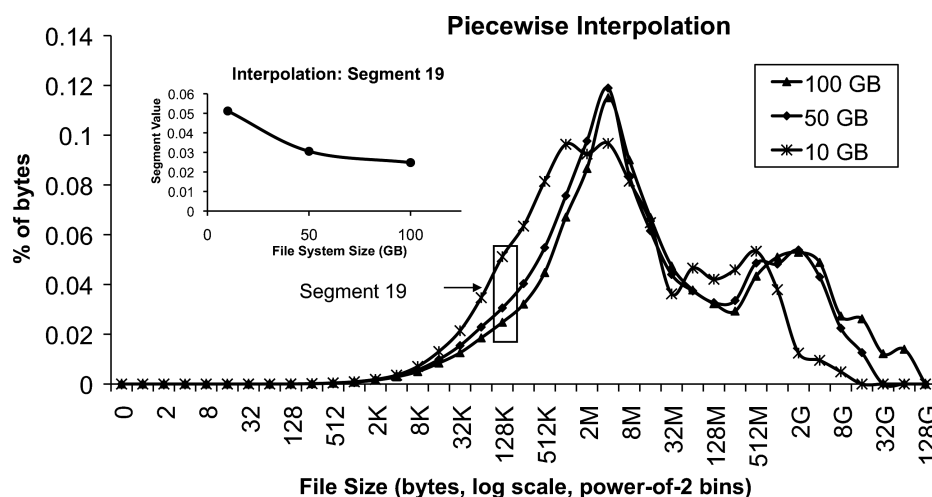


Fig. 5. Piecewise interpolation of file sizes. Piece-wise interpolation for the distribution of files with bytes, using file systems of 10 GB, 50 GB and 100 GB. Each power-of-two bin on the x-axis is treated as an individual *segment* for interpolation (inset). Final curve is the composite of all individual interpolated segments.

insignificant, especially for outlying regions, it may not serve as an accurate representation. Impressions thus provides the capability for interpolation and extrapolation from available data and distributions.

Impressions needs to generate complete new curves from existing ones. To illustrate our procedure, we describe an example of creating an interpolated curve; extensions to extrapolation are straightforward. Figure 5 shows how Impressions uses *piece-wise interpolation* for the distribution of files with containing bytes. In this example, we start with the distribution of file sizes for file systems of size 10 GB, 50 GB, and 100 GB, shown in the figure. Each power-of-two bin on the x-axis is treated as an individual *segment*, and the available data points within each segment are used as input for piece-wise interpolation; the process is repeated for all segments of the curve. Impressions combines the individual interpolated segments to obtain the complete interpolated curve.

To demonstrate the accuracy of our approach, we interpolate and extrapolate file size distributions for file systems of sizes 75 GB and 125 GB, respectively.

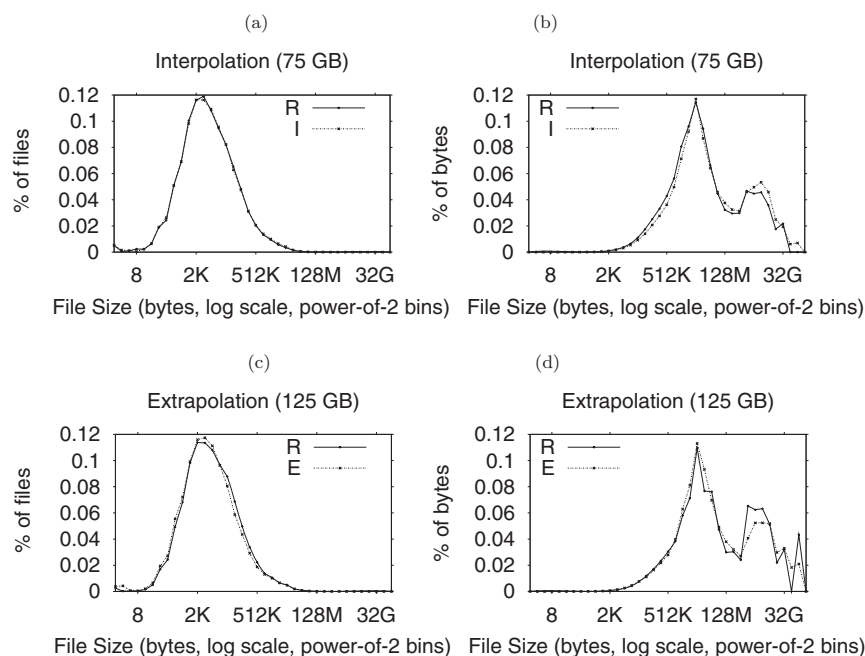


Fig. 6. Accuracy of interpolation and extrapolation. Shows results of applying piecewise interpolation to generate file size distributions (by count and by bytes), for file systems of size 75 GB (a and b, respectively), and 125 GB (c and d, respectively).

Table V. Accuracy of Interpolation and Extrapolation
Impressions produces accurate curves for file systems of size 75 GB and 125 GB, using interpolation (I) and extrapolation (E), respectively.

Distribution	FS Region (I/E)	D Statistic	K-S Test (0.05)
File sizes by count	75 GB (I)	0.054	passed
File sizes by count	125 GB (E)	0.081	passed
File sizes by bytes	75 GB (I)	0.105	passed
File sizes by bytes	125 GB (E)	0.105	passed

Figure 6 shows the results of applying our technique, comparing the generated distributions with actual distributions for the file system sizes (we removed this data from the dataset used for interpolation). We find that the simpler curves such as Figure 6(a) and (c) are interpolated and extrapolated with good accuracy. Even for more challenging curves such as Figure 6(b) and (d), the results are accurate enough to be useful. Table V contains the results of conducting K-S tests to measure the goodness-of-fit of the generated curves. All the generated distributions passed the K-S test at the 0.05 significance level.

3.6 File Content

Actual file content can have substantial impact on the performance of an application. For example, Postmark [Katcher 1997], one of the most popular file

system benchmarks, tries to simulate an email workload, yet it pays scant attention to the organization of the file system, and is completely oblivious of the file data. Postmark fills all the “email” files with the same data, generated using the same random seed. The evaluation results can range from misleading to completely inaccurate, for instance in the case of content-addressable storage (CAS). When evaluating a CAS-based system, the disk-block traffic and the corresponding performance will depend only on the unique content—in this case belonging to the largest file in the file system. Similarly, performance of Desktop Search and Word Processing applications is sensitive to file content.

In order to generate representative file content, Impressions supports a number of options. For human-readable files such as .txt, .html files, it can populate file content with random permutations of symbols and words, or with more sophisticated word-popularity models. Impressions maintains a list of the relative popularity of the most popular words in the English language, and a Monte Carlo simulation generates words for file content according to this model. However, the distribution of word popularity is heavy-tailed; hence, maintaining an exhaustive list of words slows down content generation. To improve performance, we use a word-length frequency model [Sigurd et al. 2004] to generate the long tail of words, and use the word-popularity model for the body alone. According to the word-length frequency model the observed frequencies of word lengths is approximated by a variant of the gamma distribution, and is of the general form: $f_{exp} = a * L^b * c^L$, where f_{exp} is the observed frequency for word-length L , and (a,b,c) are language-specific parameters.

The user has the flexibility to select either one of the models in entirety, or a specific combination of the two. It is also relatively straightforward to add extensions in the future to generate more nuanced file content. An example of such an extension is one that carefully controls the degree of content similarity across files.

In order to generate content for typed files, Impressions either contains enough information to generate valid file headers and footers itself, or calls into a third-party library or software such as Id3v2¹ for mp3; GraphApp² for gif, jpeg and other image files; Mplayer³ for mpeg and other video files; asciidoc for html; and ascii2pdf for PDF files.

3.7 Disk Layout and Fragmentation

To isolate the effects of file system content, Impressions can measure the degree of on-disk fragmentation, and create file systems with user-defined degree of fragmentation. The extent of fragmentation is measured in terms of *layout score* [Smith and Seltzer 1997]. A layout score of 1 means all files in the file system are laid out optimally on disk (i.e., all blocks of any given file are laid out consecutively one after the other), while a layout score of 0 means that no two blocks of any file are adjacent to each other on disk.

¹<http://id3v2.sourceforge.net/>.

²<http://enchantia.com/software/graphapp/>.

³<http://www.mplayerhq.hu/>.

Table VI. Performance of Impressions
 Shows time taken to create file-system images with break down for individual features. *Image₁*: 4.55 GB, 20000 files, 4000 dirs. *Image₂*: 12.0 GB, 52000 files, 4000 dirs. Other parameters are default. The two entries for additional parameters are shown only for *Image₁* and represent times in addition to default times.

FS Distribution (Default)	Time Taken (seconds)	
	<i>Image₁</i>	<i>Image₂</i>
Directory structure	1.18	1.26
File sizes distribution	0.10	0.28
Popular extensions	0.05	0.13
File with depth	0.064	0.29
File and bytes with depth	0.25	0.70
File content (Single-word)	0.53	1.44
On-disk file/dir creation	437.80	1394.84
Total time	473.20 (8 mins)	1826.12 (30 mins)
File content (Hybrid model)	791.20	–
Layout score (0.98)	133.96	–

Impressions achieves the desired degree of fragmentation by issuing pairs of temporary file create and delete operations, during creation of regular files. When experimenting with a file-system image, Impressions gives the user complete control to specify the overall layout score. In order to determine the on-disk layout of files, we rely on the information provided by `debugfs`. Thus currently we support layout measurement only for Ext2 and Ext3. In future work, we will consider several alternatives for retrieving file layout information across a wider range of file systems. On Linux, the `FIBMAP` and `FIEMAP ioctl()`s are available to map a logical block to a physical block.⁴ Other file system-specific methods exist, such as the `XFS_IOC_GETBMAP ioctl` for XFS.

The previous approach however does not account for differences in fragmentation strategies across file systems. Impressions supports an alternate specification for the degree of fragmentation wherein it runs a pre-specified workload and reports the resulting layout score. Thus if a file system employs better strategies to avoid fragmentation, it is reflected in the final layout score after running the fragmentation workload.

There are several alternate techniques for inducing more realistic fragmentation in file systems. Factors such as burstiness of I/O traffic, out-of-order writes and inter-file layout are currently not accounted for; a companion tool to Impressions for carefully creating fragmented file systems will thus be a good candidate for future research.

3.8 Performance

In building Impressions, our primary objective was to generate realistic file-system images, giving top priority to accuracy, instead of performance. Nonetheless, Impressions does perform reasonably well. Table VI shows the breakdown

⁴<http://lwn.net/Articles/260795/>.

App	Parameter & Value	Comment on Validity
GDL	File content < 10 deep	10% of files and 5% of bytes > 10 deep (content in deeper namespace is growing)
GDL	Text file sizes < 200 KB	13% of files and 90% of bytes > 200 KB
Beagle	Text file cutoff < 5 MB	0.13% of files and 71% of bytes > 5 MB
Beagle	Archive files < 10 MB	4% of files and 84% of bytes > 10 MB
Beagle	Shell scripts < 20 KB	20% of files and 89% of bytes > 20 KB

Fig. 7. Debunking application assumptions. Examples of assumptions made by Beagle and GDL, along with details of the amount of file-system content that is not indexed as a consequence.

of time taken to create a default file-system image of 4.55 GB. We also show time taken for some additional features such as using better file content, and creating a fragmented file system. Overall, we find that *Impressions* creates highly accurate file-system images in a reasonable amount of time and thus is useful in practice.

4. CASE STUDY: DESKTOP SEARCH

In this section, we use *Impressions* to evaluate desktop searching applications. Our goals for this case study are twofold. First, we show how simple it is to use *Impressions* to create either representative images or images across which a single parameter is varied. Second, we show how future evaluations should report the settings of *Impressions* so that results can be easily reproduced.

We choose desktop search for our case study because its performance and storage requirements depend not only on the file system size and structure, but also on the type of files and the actual content within the files. We evaluate two desktop search applications: open-source Beagle⁵ and Google’s Desktop for Linux (GDL).⁶ Beagle supports a large number of file types using 52 search-filters; it provides several indexing options, trading performance and index size with the quality and feature-richness of the index. Google Desktop does not provide as many options: a Web interface allows users to select or exclude types of files and folder locations for searching, but does not provide any control over the type and quality of indexing.

4.1 Representative Images

Developers of data-intensive applications frequently need to make assumptions about the properties of file-system images. For example, file systems and applications can often be optimized if they know properties such as the relative proportion of metadata to data in representative file systems. Previously, developers could infer these numbers from published papers [Agrawal et al. 2007; Douceur and Bolosky 1999; Satyanarayanan 1981; Sienknecht et al. 1994], but only with considerable effort. With *Impressions*, developers can simply create a sample of representative images and directly measure the properties of interest.

Figure 7 lists assumptions we found in GDL and Beagle limiting the search indexing to partial regions of the file system. However, for the representative file

⁵<http://www.beagle-project.org/>.

⁶<http://desktop.google.com/linux/index.html>.

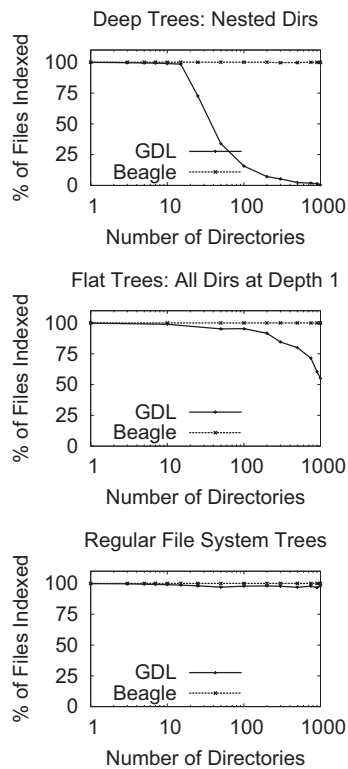


Fig. 8. Tree depth and completeness of index. Shows the percentage of files indexed by Beagle and GDL with varying directory tree depths in a given file-system image.

systems in our data set, these assumptions omit large portions of the file system. For example, GDL limits its index to only those files less than ten directories deep; our analysis of typical file systems indicates that this restriction causes 10% of all files to be missed.

Figure 8 shows one such example: it compares the percentage of files indexed by Beagle and GDL for a set of file-system images. The topmost graph shows the results for *deep* file-system trees created by successively nesting a new directory in the parent directory; a file system with D directories will thus have a maximum depth of D . The y-axis shows the % of files indexed, and the x-axis shows the number of directories in the file system. We find that GDL stops indexing content after depth 10, while Beagle indexes 100% of the files. The middle graph repeats the experiment on flat trees, with all directories at depth 1. This time, GDL's percentage completeness drops off once the number of directories exceeds 10. For regular file system trees, shown in the lowermost graph, we find that both Beagle and GDL achieve near 100% completeness. Since the percentage of user-generated content deeper in the namespace is growing over the years, it might be useful to design search indexing schemes which are better suited for deeper name spaces.

This strange behavior further motivates the need for a tool like *Impressions* to be a part of any application designer’s toolkit. We believe that instead of arbitrarily specifying hard values, application designers should experiment with *Impressions* to find acceptable choices for representative images.

We note that *Impressions* is useful for discovering these application assumptions and for isolating performance anomalies that depend on the file-system image. Isolating the impact of different file system features is easy using *Impressions*: evaluators can use *Impressions* to create file-system images in which only a single parameter is varied, while all other characteristics are carefully controlled.

This type of discovery is clearly useful when one is using closed-source code, such as GDL. For example, we discovered the GDL limitations by constructing file-system images across which a single parameter is varied (e.g., file depth and file size), measuring the percentage of indexed files, and noticing precipitous drops in this percentage. This type of controlled experimentation is also useful for finding non-obvious performance interactions in open-source code. For instance, Beagle uses the *inotify* mechanism⁷ to track each directory for change; since the default Linux kernel provides 8192 watches, Beagle resorts to manually crawling the directories once their count exceeds 8192. This deterioration in performance can be easily found by creating file-system images with varying numbers of directories.

4.2 Reproducible Images

The time spent by desktop search applications to crawl a file-system image is significant (i.e., hours to days); therefore, it is likely that different developers will innovate in this area. In order for developers to be able to compare their results, they must be able to ensure they are using the same file-system images. *Impressions* allows one to precisely control the image and report the parameters so that the exact same image can be reproduced.

For desktop search, the type of files (i.e., their extensions) and the content of files has a significant impact on the time to build the index and its size. We imagine a scenario in which the Beagle and GDL developers wish to compare index sizes. To make a meaningful comparison, the developers must clearly specify the file-system image used; this can be done easily with *Impressions* by reporting the size of the image, the distributions listed in Table II, the word model, disk layout, and the random seed. We anticipate that most benchmarking will be done using mostly default values, reducing the number of *Impressions* parameters that must be specified.

An example of the reporting needed for reproducible results is shown in Figure 9. In these experiments, all distributions of the file system are kept constant, but only either text files (containing either a single word or with the default word model) or binary files are created. These experiments illustrate the point that file content significantly affects the index size; if two systems are compared using different file content, obviously the results are meaningless. Specifically, different file types change even the relative ordering of index size

⁷<http://www.linuxjournal.com/article/8478>.

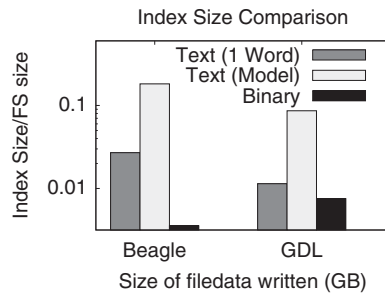


Fig. 9. Impact of file content. Compares Beagle and GDL index time and space for wordmodels and binary files. Google has a smaller index for wordmodels, but larger for binary. Uses Impressions default settings, with FS size 4.55 GB, 20000 files, 4000 dirs.

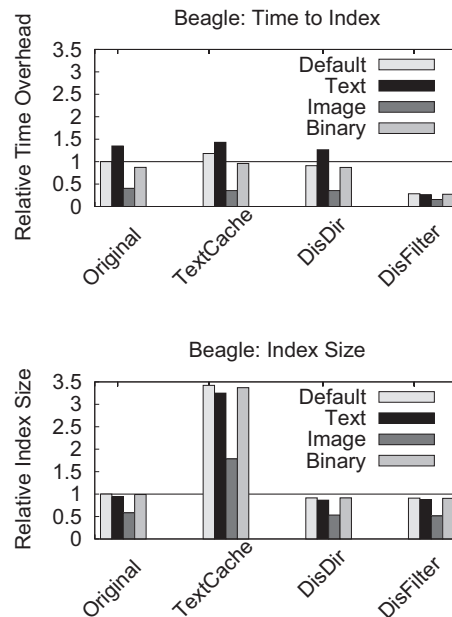


Fig. 10. Reproducible images: impact of content. Using Impressions to make results reproducible for benchmarking search. Vertical bars represent file systems created with file content as labeled. The *Default* file system is created using Impressions default settings, and file system size 4.55 GB, 20000 files, 4000 dirs. Index options: *Original*—default Beagle index. *TextCache*—build text-cache of documents used for snippets. *DisDir*—don't add directories to the index. *DisFilter*—disable all filtering of files, only index attributes.

between Beagle and GDL: given text files, Beagle creates a larger index; given binary files, GDL creates a larger index.

Figure 10 gives an additional example of reporting Impressions parameters to make results reproducible. In these experiments, we discuss a scenario in which different developers have optimized Beagle and wish to meaningfully compare their results. In this scenario, the original Beagle developers reported results for four different images: the default, one with only text files, one with

only image files, and one with only binary files. Other developers later create variants of Beagle: *TextCache* to display a small portion of every file alongside a search hit, *DisDir* to disable directory indexing, and *DisFilter* to index only attributes. Given the reported Impressions parameters, the variants of Beagle can be meaningfully compared to one another.

In summary, Impressions makes it extremely easy to create both controlled and representative file-system images. Through this brief case study evaluating desktop search applications, we have shown some of the advantages of using Impressions. First, Impressions enables developers to tune their systems to the file system characteristics likely to be found in their target user populations. Second, it enables developers to easily create images where one parameter is varied and all others are carefully controlled; this allows one to assess the impact of a single parameter. Finally, Impressions enables different developers to ensure they are all comparing the same image; by reporting Impressions parameters, one can ensure that benchmarking results are reproducible.

5. OTHER APPLICATIONS

Besides its use in conducting representative and reproducible benchmarking, Impressions can also be handy in other experimental scenarios. In this section we present two examples, the usefulness of Impressions in generating realistic rules of thumb, and in testing soundness of hypothesis.

5.1 Generating Realistic Rules of Thumb

In spite of the availability of Impressions, designers of file systems and related software will continue to rely on rules of thumb to make design decisions. Instead of relying on old wisdom, one can use Impressions to generate realistic rules of thumb. One example of such a rule of thumb is to calculate the overhead of file-system metadata—a piece of information often needed to compute the cost of different replication, parity or check summing schemes for data reliability. Figure 11 shows the percentage of space taken by metadata in a file system, as we vary the distribution of file sizes. We find that the overhead can vary between 2 and 14% across the file size distributions in this example. Similarly, Impressions can be used to compute other rules of thumb for different metadata properties.

5.2 Testing Hypothesis

In our experience, we found Impressions convenient and simple to use for testing hypothesis regarding application and file system behavior, hiding away the statistical complexity of the experiment from the end-user. To illustrate this, we describe our experience with a *failed* experiment.

It was our hypothesis that the distribution of bytes and files by namespace depth would affect the time taken to build the search index: indexing file content in deeper namespace would be slower. To test our hypothesis, all we had to do was use Impressions to create file-system images, and measure the time taken by Beagle to build the index, varying only a single parameter in the configuration file for each trial: the λ value governing the Poisson distribution for file

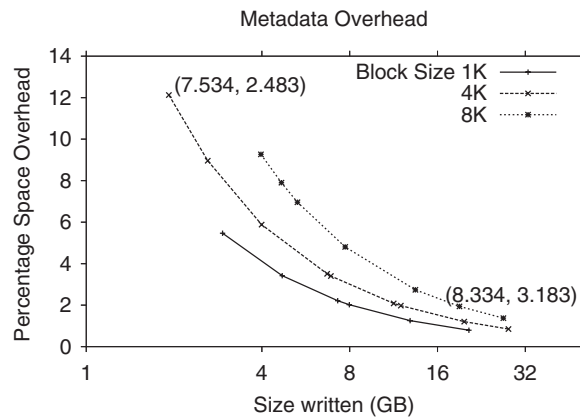


Fig. 11. Metadata overhead. Shows the relative space overhead of file-system metadata with varying file-size distribution, modeled by (μ, σ) parameters of a lognormal distribution (shown in parentheses for the two extremes).

depth. Although our hypothesis was not validated by the results (i.e., we didn't find significant variation in indexing time with depth), we found Impressions to be suitable and easy to use for such experimentation.

6. RELATED WORK

We discuss previous research in four related areas. First, we discuss previous studies on file-system metadata; second, we discuss existing tools for generating file-system images; third, we present prior research on improving file system benchmarking; finally, we discuss existing models for explaining file system metadata properties.

6.1 File System Measurement Studies

Impressions enables file system measurement studies to be put into practice. Besides the metadata studies on Windows workstations [Agrawal et al. 2007; Douceur and Bolosky 1999], previous work in non-Windows environments includes Satyanarayanan's study of a Digital PDP-10 [Satyanarayanan 1981], Irlam's and Mullender's studies of Unix systems [Irlam 1993; Mullender and Tanenbaum 1984], and the study of HP-UX systems at Hewlett-Packard [Sienknecht et al. 1994]. These studies provide valuable data for designers of file systems and related software, and can be incorporated in Impressions.

6.2 Tools for Generating File-System Images

We are not aware of any existing system that generates file-system images with the level of detail that Impressions delivers; here we discuss some tools that we believe provide some subset of features supported by Impressions.

FileBench, a file system workload framework for measuring and comparing file system performance [McDougall] is perhaps the closest to Impressions in terms of flexibility and attention to detail. FileBench generates test file system

images with support for different directory hierarchies with namespace depth and file sizes according to statistical distributions. We believe Impressions includes all the features provided by FileBench and provides additional capabilities; in particular, Impressions allows one to contribute newer datasets and makes it easier to plug in distributions. FileBench also does not provide support for allowing user-specified constraints.

The SynRGen file reference generator by Ebling and Satyanarayanan [1994] generates synthetic equivalents for real file system users. The *volumes* or images in their work make use of simplistic assumptions about the file system distributions as their focus is on user access patterns.

File system and application developers in the open-source community also require file-system images to test and benchmark their systems, tools for which are developed in-house, often customized to the specific needs of the system being developed.

Genbackupdata is one such tool that generates test data sets for performance testing of backup software [Wirzenius 2009]. Like Impressions, but in a much simplified fashion, it creates a directory tree with files of different sizes. Since the tool is specifically designed for backup applications, the total file system size and the minimum and maximum limits for file sizes are configurable, but not the file size distribution or other aspects of the file system. The program can also modify an existing directory tree by creating new files, and deleting, renaming, or modifying existing files, inducing fragmentation on disk.

Another benchmarking system that generates test file systems matching a specific profile is Fstress [Anderson and Chase 2002]. However, it does contain many of the features found standard in Impressions, such as popularity of file extensions and file content generation according to file types, supporting user-specified distributions for file system parameters and allowing arbitrary constraints to be specified on those parameters.

6.3 Tools and Techniques for Improving Benchmarking

A number of tools and techniques have been proposed to improve the state of the art of file and storage system benchmarking. Chen and Patterson proposed a “self-scaling” benchmark that scales with the I/O system being evaluated, to stress the system in meaningful ways [Chen and Patterson 1993]. Although useful for disk and I/O systems, the self-scaling benchmarks are not directly applicable for file systems.

TBBT is a NFS trace replay tool that derives the file-system image underlying a trace [Zhu et al. 2005]. It extracts the file system hierarchy from a given trace in depth-first order and uses that during initialization for a subsequent trace replay. While this ensures a consistent file-system image for replay, it does not solve the more general problem of creating accurately controlled images for all types of file system benchmarking.

The Auto-Pilot tool [Wright et al. 2005] provides an infrastructure for running tests and analysis tools to automate the benchmarking process. Auto-Pilot can help run benchmarks with relative ease by automating the repetitive tasks of running, measuring, and analyzing a program through test scripts.

6.4 Models for File-System Metadata

Several models have been proposed to explain observed file-system phenomena. Mitzenmacher [2002] proposed a generative model, called the Recursive Forest File model to explain the behavior of file size distributions. The model is dynamic as it allows for the creation of new files and deletion of old files. The model accounts for the hybrid distribution of file sizes with a lognormal body and Pareto tail.

Downey's Multiplicative File Size model [Downey 2001] is based on the assumption that new files are created by using older files as templates for example, by copying, editing or filtering an old file. The size of the new file in this model is given by the size of the old file multiplied by an independent factor.

The HOT (Highly Optimized Tolerance) model provides an alternate generative model for file size distributions. These models provide an intuitive understanding of the underlying phenomena, and are also easier for computer simulation. In future, Impressions can be enhanced by incorporating more such models.

7. CONCLUSION

File system benchmarking is in a state of disarray. One key aspect of this problem is generating realistic file-system state, with due emphasis given to file-system metadata and file content. To address this problem, we have developed Impressions, a statistical framework to generate realistic and configurable file-system images. Impressions provides the user flexibility in selecting a comprehensive set of file system parameters, while seamlessly ensuring accuracy of the underlying images, serving as a useful platform for benchmarking.

In our experience, we find Impressions easy to use and well suited for a number of tasks. It enables application developers to evaluate and tune their systems for realistic file system characteristics, representative of target usage scenarios. Impressions also makes it feasible to compare the performance of systems by standardizing and reporting all used parameters, a requirement necessary for benchmarking. We believe Impressions will prove to be a valuable tool for system developers and users alike; we have made it publicly available for download. Please visit the URL <http://www.cs.wisc.edu/adsl/Software/Impressions/> to obtain a copy.

ACKNOWLEDGMENTS

We are grateful to Bill Bolosky for providing us with a copy of the five-year metadata dataset from Microsoft. Lakshmi Bairavasundaram provided many useful discussions and gave valuable comments on earlier drafts of this paper. Finally, we would like to thank Valerie Aurora Henson and the anonymous FAST reviewers for their excellent feedback and comments.

REFERENCES

- AGRAWAL, N., BOLOSKY, W. J., DOUCEUR, J. R., AND LORCH, J. R. 2007. A five-year study of file-system metadata. In *Proceedings of the 5th USENIX Symposium on File and Storage Technologies (FAST'07)*.

- ANDERSON, D. AND CHASE, J. 2002. Fstress: A flexible network file service benchmark. In Tech rep. Duke University.
- ANDERSON, E., KALLAHALLA, M., UYSAL, M., AND SWAMINATHAN, R. 2004. Buttress: A toolkit for flexible and high fidelity I/O benchmarking. In *Proceedings of the 3rd USENIX Symposium on File and Storage Technologies (FAST'04)*.
- BAKER, M., HARTMAN, J., KUPFER, M., SHIRRIFF, K., AND OUSTERHOUT, J. 1991. Measurements of a distributed file system. In *Proceedings of the 13th ACM Symposium on Operating Systems Principles (SOSP'91)*. 198–212.
- CHEN, P. M. AND PATTERSON, D. A. 1993. A new approach to I/O performance evaluation—self-scaling I/O benchmarks, predicted I/O performance. In *Proceedings of the ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems (SIGMETRICS'93)*. 1–12.
- CIPAR, J., CORNER, M. D., AND BERGER, E. D. 2007. Tfs: A transparent file system for contributory storage. In *Proceedings of the USENIX Conference on File and Storage Technologies (FAST'07)*. USENIX Association, Berkeley, CA. 28–28.
- CORMEN, T. H., LEISERSON, C. E., RIVEST, R. L., AND STEIN, C. 2001. *Introduction to Algorithms*, 2nd Ed. MIT Press and McGraw-Hill.
- COX, L. P., MURRAY, C. D., AND NOBLE, B. D. 2002. Pastiche: Making backup cheap and easy. *SIGOPS Oper. Syst. Rev.* 36.
- COX, L. P. AND NOBLE, B. D. 2003. Samsara: Honor among thieves in peer-to-peer storage. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP'03)*. ACM, New York. 120–132.
- DAHLIN, M. D., WANG, R. Y., ANDERSON, T. E., AND PATTERSON, D. A. 1994. Cooperative caching: Using remote client memory to improve file system performance. In *Proceedings of the 1st Symposium on Operating Systems Design and Implementation (OSDI'94)*.
- DOUCEUR, J. R. AND BOLOSKY, W. J. 1999. A large-scale study of file-system contents. In *Proceedings of the Joint International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS)*. 59–70.
- DOWNNEY, A. B. 2001. The structural cause of file size distributions. In *Proceedings of the 9th International Symposium on Modeling Analysis, and Simulation of Computer-Telecommunications Systems (MASCOTS'01)*.
- EBLING, M. R. AND SATYANARAYANAN, M. 1994. Synrgen: An extensible file reference generator. In *Proceedings of the ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems (SIGMETRICS'94)*.
- FU, K., KAASHOEK, M. F., AND MAZIÈRES, D. 2002. Fast and secure distributed read-only file system. *ACM Trans. Comput. Syst.* 20, 1, 1–24.
- GOPAL, B. AND MANBER, U. 1999. Integrating content-based access mechanisms with hierarchical file systems. In *Proceedings of the 3rd Symposium on Operating Systems Design and Implementation (OSDI'99)*.
- GRIBBLE, S. D., MANKU, G. S., ROSELLI, D. S., BREWER, E. A., GIBSON, T. J., AND MILLER, E. L. 1998. Self-similarity in file systems. In *Proceedings of the Joint International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS)*. 141–150.
- HUTCHINSON, N. C., MANLEY, S., FEDERWISCH, M., HARRIS, G., HITZ, D., KLEIMAN, S., AND O'MALLEY, S. 1999. Logical vs. physical file system backup. In *Proceedings of the 3rd Symposium on Operating Systems Design and Implementation (OSDI'99)*.
- IRLAM, G. 1993. Unix file size survey—1993. <http://www.base.com/gordoni/ufs93.html>.
- KATCHER, J. 1997. PostMark: A new file system benchmark. Tech. rep. TR-3022, Network Appliance Inc.
- MESNIER, M. P., WACHS, M., SAMBASIVAN, R. R., LOPEZ, J., HENDRICKS, J., GANGER, G. R., AND O'HALLARON, D. 2007. Trace: Parallel trace replay with approximate causal events. In *Proceedings of the 5th USENIX Symposium on File and Storage Technologies (FAST'07)*.
- MCDUGALL R. Filebench: Application level file system benchmark. <http://www.solarisinternals.com/si/tools/filebench/index.php>.
- MITZENMACHER, M. 2002. Dynamic models for file sizes and double pareto distributions. In *Internet Mathematics*.
- MPLAYER. The MPlayer movie player. <http://www.mplayerhq.hu/>.

- MULLENDER, S. J. AND TANENBAUM, A. S. 1984. Immediate files. *Softw. Practice Exper.* 14, 4, 365–368.
- MUTHITACHAROEN, A., CHEN, B., AND MAZIÈRES, D. 2001. A low-bandwidth network file system. In *Proceedings of the 18th ACM Symposium on Operating Systems Principles (SOSP'01)*. 174–187.
- NIST. 2007. Text retrieval conference (trec) datasets. <http://trec.nist.gov/data>.
- OUSTERHOUT, J. K., COSTA, H. D., HARRISON, D., KUNZE, J. A., KUPFER, M., AND THOMPSON, J. G. 1985. A trace-driven analysis of the UNIX 4.2 BSD file system. In *Proceedings of the 10th ACM Symposium on Operating System Principles (SOSP'85)*. 15–24.
- PADIOLEAU, Y. AND RIDOUX, O. 2003. A logic file system. In *Proceedings of the USENIX Annual Technical Conference*.
- PATTERSON, D., GIBSON, G., AND KATZ, R. 1988. A Case for Redundant Arrays of Inexpensive Disks (RAID). In *Proceedings of the ACM SIGMOD Conference on the Management of Data (SIGMOD'88)*. 109–116.
- PRABHAKARAN, V., BAIRAVASUNDARAM, L. N., AGRAWAL, N., GUNAWI, H. S., ARPACI-DUSSEAU, A. C., AND ARPACI-DUSSEAU, R. H. 2005. IRON file systems. In *Proceedings of the 20th ACM Symposium on Operating Systems Principles (SOSP'05)*. 206–220.
- PRZYDATEK, B. 2002. A Fast Approximation Algorithm for the subset-sum problem. *Inter. Trans. Oper. Res.* 9, 4, 437–459.
- RIEDEL, E., KALLAHALLA, M., AND SWAMINATHAN, R. 2002. A framework for evaluating storage system security. In *Proceedings of the 1st USENIX Symposium on File and Storage Technologies (FAST'02)*. 14–29.
- ROWSTRON, A. AND DRUSCHEL, P. 2001. Storage management and caching in PAST, A large-scale, persistent peer-to-peer storage utility. In *Proceedings of the 18th ACM Symposium on Operating Systems Principles (SOSP'01)*.
- SATYANARAYANAN, M. 1981. A study of file sizes and functional lifetimes. In *Proceedings of the 8th ACM Symposium on Operating Systems Principles (SOSP)*. 96–108.
- SIENKNECHT, T. F., FRIEDRICH, R. J., MARTINKA, J. J., AND FRIEDENBACH, P. M. 1994. The implications of distributed data in a commercial environment on the design of hierarchical storage management. *Perform. Eval.* 20, 1–3, 3–25.
- SIGURD, B., EEG-OLOFSSON, M., AND VAN DE WEIJER, J. 2004. Word length, sentence length and frequency—Zipf revisited. *Studia Linguist.* 58, 1, 37–52.
- SMITH, K. AND SELTZER, M. I. 1997. File system aging. In *Proceedings of the Sigmetrics Conference*.
- SNIA. 2007. Storage network industry association: Lotta repository. <http://iota.snia.org>.
- SOBTI, S., GARG, N., ZHENG, F., LAI, J., SHAO, Y., ZHANG, C., ZISKIND, W., AND KRISHNAMURTHY, A. 2004. Segank: A distributed mobile storage system. In *Proceedings of the 3rd USENIX Symposium on File and Storage Technologies (FAST'04)*. 239–252.
- STORER, M. W., GREENAN, K. M., MILLER, E. L., AND VORUGANTI, K. 2008. Pergamum: Replacing tape with energy efficient, reliable, disk-based archival storage. In *Proceedings of the 6th USENIX Conference on File and Storage Technologies (FAST'08)*. USENIX Association, Berkeley, CA, 1–16.
- WIRZENIUS, L. 2009. Genbackupdata: Tool to generate backup test data. <http://braawi.org/genbackupdata.html>.
- WRIGHT, C. P., JOUKOV, N., KULKARNI, D., MIRETSKIY, Y., AND ZADOK, E. 2005. Auto-pilot: A platform for system software benchmarking. In *Proceedings of the Annual USENIX Technical Conference, FREENIX Track*.
- ZHANG, Z. AND GHOSE, K. 2003. yfs: A journaling file system design for handling large data sets with reduced seeking. In *Proceedings of the 2nd USENIX Conference on File and Storage Technologies (FAST'03)*. USENIX Association, Berkeley, CA, 59–72.
- ZHU, N., CHEN, J., AND CHIUEH, T.-C. 2005. Tbbt: Scalable and accurate trace replay for file server evaluation. In *Proceedings of the 4th USENIX Conference on File and Storage Technologies*. USENIX Association, Berkeley, CA, 24–24.

Received August 2009; accepted August 2009