

# GUFI User's Guide

GUFI Developers

September 11, 2023

## Contents

<b>1</b>	<b>License</b>	<b>2</b>
<b>2</b>	<b>Introduction</b>	<b>4</b>
<b>3</b>	<b>Environment</b>	<b>5</b>
<b>4</b>	<b>gufi_find</b>	<b>6</b>
4.1	Flags . . . . .	6
<b>5</b>	<b>gufi_ls</b>	<b>7</b>
5.1	Flags . . . . .	7
<b>6</b>	<b>gufi_stat</b>	<b>8</b>
6.1	Flags . . . . .	8
<b>7</b>	<b>gufi_stats</b>	<b>9</b>
7.1	Flags . . . . .	9
7.1.1	Recursive . . . . .	9
7.1.2	Cumulative . . . . .	10
7.1.3	Other . . . . .	10

# 1 License

This file is part of GUF1, which is part of MarFS, which is released under the BSD license.

Copyright (c) 2017, Los Alamos National Security (LANS), LLC  
All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

1. Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
2. Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
3. Neither the name of the copyright holder nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

From Los Alamos National Security, LLC:  
LA-CC-15-039

Copyright (c) 2017, Los Alamos National Security, LLC All rights reserved.  
Copyright 2017. Los Alamos National Security, LLC. This software was produced under U.S. Government contract DE-AC52-06NA25396 for Los Alamos National Laboratory (LANL), which is operated by Los Alamos National Security, LLC for the U.S. Department of Energy. The U.S. Government has rights to use, reproduce, and distribute this software. NEITHER THE GOVERNMENT NOR LOS ALAMOS NATIONAL SECURITY, LLC MAKES ANY WARRANTY, EXPRESS OR IMPLIED, OR ASSUMES ANY LIABILITY FOR THE USE OF THIS SOFTWARE. If software is

modified to produce derivative works, such modified software should be clearly marked, so as not to confuse it with the version available from LANL.

THIS SOFTWARE IS PROVIDED BY LOS ALAMOS NATIONAL SECURITY, LLC AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL LOS ALAMOS NATIONAL SECURITY, LLC OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

## 2 Introduction

Over the years, the amount of data we store and use has grown exponentially to the point that petabytes of storage is not uncommon. What used to be a simple task of accessing and sorting through information has been compounded into an arduous task with the size and scale of super-computing data centers. Being able to query data effectively, while also taking into account user permissions becomes paramount into accomplishing daily tasks. This is what the Grand Unified File Index (GUFU) tool aims to accomplish.

This process of efficiently accessing data is accomplished by recreating the tree structure via indexing. Each directory contains an SQL database file that stores the metadata of the files as well as summary information for that directory and optionally summary information for the entire tree below that directory.

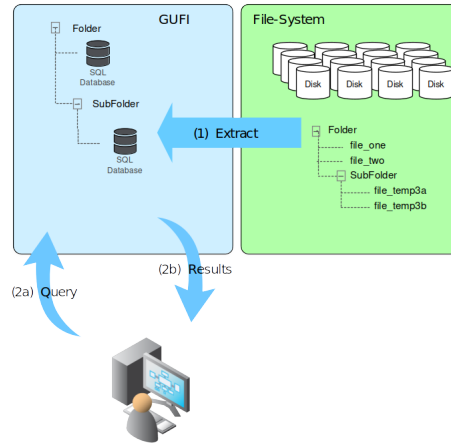


Figure 1: Layout and user interaction with GUFU

### 3 Environment

Users interact with GUFi using `gufi_find`, `gufi_ls`, `gufi_stat`, and `gufi_stats`. These are wrapper scripts that access the actual implementations of these files.

There should be a file readable (but not necessarily modifiable) by all GUFi users called `/etc/GUFi/config`. This file specifies where the GUFi server is, as the actual GUFi trees are not expected to be locally available.

## 4 gufi\_find

`gufi_find` is a wrapper script for `gufi_query` that attempts to recreate a large portion of GNU `find(1)`'s functionality.

One major difference between them is how arguments are parsed. In `find(1)`, expression order matters. In `gufi_find`, expression order does not matter.

### 4.1 Flags

The following `find(1)` test expression options are replicated. See `find(1)` for details.

`-amin`, `-atime`, `-cmin`, `-ctime`, `-empty`, `-executable`, `-false`, `-fprint`, `-gid`, `-group`, `-help`, `-iname`, `-inum`, `-links`, `-lname`, `-maxdepth`, `-mindepth`, `-mmin`, `-mtime`, `-name`, `-newer`, `-path`, `-printf`, `-readable`, `-samefile`, `-size`, `-true`, `-type`, `-uid`, `-user`, and `-writable`

The action expressions `printf` and `fprint` are also available. `printf` writes to `stdout`, and `fprint` writes to per-thread output text files. All format specifiers other than `ACFTYZ` have been implemented.

Additionally, a few GUFi extensions have been added:

Option	Description
<code>--size% num num</code>	Extends the <code>size</code> option to search for sizes in the range <code>[size + min%, size + max%]</code> . e.g. <code>-size 100c --size% 10 20</code> $\rightarrow$ <code>[110, 120]</code> <code>-size 100c --size% -10 20</code> $\rightarrow$ <code>[90, 120]</code>
<code>--num-results num</code>	Limit the number of results printed
<code>--smallest</code>	Output by size, ascending.
<code>--largest</code>	Output by size, descending.
<code>--in-memory-name name</code>	Change the name of the tables used to store intermediate results when aggregating. Generally not used.

## 5 gufi\_ls

`gufi_ls` is the equivalent of the POSIX command `ls` applied to a GUFi tree. As with `gufi_find`, there are a multitude of options available listed below

### 5.1 Flags

Flags	Functionality
<code>--help</code>	displays help menu
<code>-v, --version</code>	show program's version number and exit
<code>-a, --all</code>	do not ignore entries ending with <code>.</code>
<code>-A, --almost-all</code>	do not list implied <code>.</code> and <code>..</code>
<code>--block-size &lt;block_size&gt;</code>	with <code>-l</code> , scale sizes by <code>block_size</code> when printing them
<code>-B, --ignore-backups</code>	do not list implied entries ending with <code>~</code>
<code>-G, --no-group</code>	in a long listing, don't print group names
<code>-i, --inode</code>	print the index number of each file
<code>-l</code>	used a long listing format
<code>-r, --reverse</code>	reverse order while sorting
<code>-R, --recursive</code>	list sub-directories recursively
<code>-s, --size</code>	print the allocated size of each file in blocks
<code>-S</code>	sort by file size, largest first
<code>--time-style &lt;TIME_STYLE&gt;</code>	time/date format with <code>-l</code>
<code>-t</code>	sort by modification time, newest first
<code>-U</code>	do not sort; list entries in directory order

Table 1: `gufi_ls` Flags and Functionality

Additionally, several GUFi extensions have been added:

Option	Description
<code>--delim &lt;c&gt;</code>	delimiter separating output columns
<code>--in-memory-name &lt;name&gt;</code>	Name of in-memory database when <code>-R</code> is used
<code>--aggregate-name &lt;name&gt;</code>	Name of final database when <code>-R</code> is used
<code>--nlink-width &lt;chars&gt;</code>	Width of <code>nlink</code> column
<code>--size-width &lt;chars&gt;</code>	Width of <code>size</code> column
<code>--user-width &lt;chars&gt;</code>	Width of <code>user</code> column
<code>--group-width &lt;chars&gt;</code>	Width of <code>group</code> column
<code>--skip-file &lt;filename&gt;</code>	Name of file containing directory basenames to skip

## 6 gufi\_stat

`gufi_stat` is analogous to `stat(1)` in file status mode.

### 6.1 Flags

Option	Description
<code>-c FORMAT, --format FORMAT</code>	use the specified FORMAT instead of the default; output a newline after each use of FORMAT
<code>-t, --terse</code>	print the information in terse form
<code>--help</code>	display this help and exit
<code>--version</code>	output version information and exit



## 7 gufi\_stats

`gufi_stats` is used to analyze a tree and retrieve statistics from an index. `gufi_stat` does not have an analogous standard utility.

### 7.1 Flags

`gufi_stats` is called with at minimum one positional argument specifying the statistic desired. Each statistic has a default version, Most also have a more in-depth version that is specified by their category: recursive, cumulative, or other.

Optional Flags	Functionality
<code>--help</code>	displays help menu
<code>--version, -v</code>	display program's version number and exits
<code>--recursive, -r</code>	run command recursively
<code>--cumulative, -c</code>	return cumulative values
<code>--order &lt;order&gt;</code>	sort output (if applicable)
<code>--delim &lt;c&gt;</code>	delimiter separating output columns
<code>--num-results &lt;n&gt;</code>	first n results
<code>--uid &lt;u&gt;, --user &lt;u&gt;</code>	restrict to user
<code>--in-memory-name &lt;name&gt;</code>	Name of intermediate database
<code>--aggregate-name &lt;name&gt;</code>	Name of final database
<code>--skip-file &lt;filename&gt;</code>	Name of file containing directory basenames to skip

Table 2: `gufi_stats` Flags and Functionality

#### 7.1.1 Recursive

These statistics will run a computation on a single directory. Specifying `-r` or `--recursive` will recursively descend the starting directory and return each subdirectory's statistic.

Statistic	Description
<code>depth</code>	Get the depths of the provided directory relative to the root directory
<code>filesize</code>	Get the size the files in the immediate directory
<code>filecount</code>	Get the number of files in the immediate directory
<code>linkcount</code>	Get the number of links in the immediate directory
<code>dircount</code>	Get the number of directories in the immediate directory
<code>leaf-dirs</code>	Get the leaf directories immediately under the current directory
<code>leaf-depth</code>	Get the depth of the leaf directories immediately under the current directory
<code>leaf-files</code>	Get number of files in the leaf directories immediately under the current directory
<code>leaf-links</code>	Get number of links in the leaf directories immediately under the current directory

### 7.1.2 Cumulative

These statistics will run a computation on an entire subtree. By default, the results will be grouped together by UID. Specifying `-c` or `--cumulative` will combine all of the results into a single sum.

Statistic	Description
total-filesize	Get the total size taken up by the entire directory
total-filecount	Get the total number of files under the directory
total-linkcount	Get the total number of links under the directory
total-dircount	Get the total number of directories under the directory
total-leaf-files	Get total number of files in the leaf directories under this directory
total-leaf-links	Get total number of links in the leaf directories under this directory
files-per-level	Get counts of how many files are in each level of the tree
links-per-level	Get counts of how many links are in each level of the tree
dirs-per-level	Get counts of how many dirs are in each level of the tree
average-leaf-files	Get average number of leaf files under the provided directory
average-leaf-links	Get average number of leaf links under the provided directory

### 7.1.3 Other

Statistics under the “other” category do not have extra flags associated with them.

Statistic	Description
median-leaf-files	Get median number of leaf files under the provided directory
duplicate-names	Find files with matching names and sizes