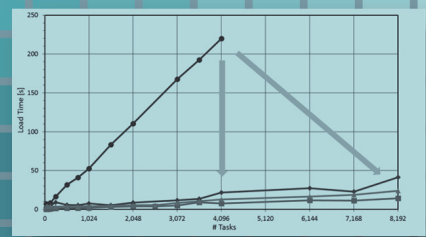
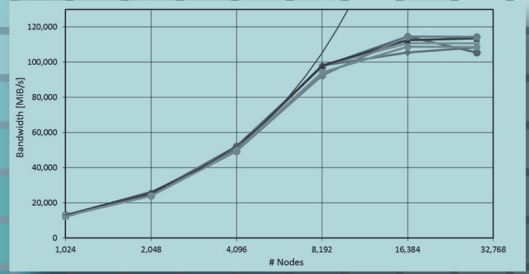


SIONlib

SPINDLE



## Efficient Task-Local I/O Operations of Massively Parallel Applications

Wolfgang Frings

Forschungszentrum Jülich GmbH  
Institute for Advanced Simulation (IAS)  
Jülich Supercomputing Centre (JSC)

# Efficient Task-Local I/O Operations of Massively Parallel Applications

Wolfgang Frings

Schriften des Forschungszentrums Jülich  
IAS Series

Band 30

---

ISSN 1868-8489

ISBN 978-3-95806-152-1

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Parallel Computing . . . . .	1
1.1.1	Parallel architectures . . . . .	2
1.1.2	Parallel software . . . . .	3
1.2	Parallel I/O . . . . .	6
1.2.1	Characteristics . . . . .	7
1.2.2	Access methods . . . . .	10
1.2.3	Parallel file systems . . . . .	13
1.3	Parallel Task-Local I/O . . . . .	17
1.3.1	Serial I/O based on the POSIX I/O interface . . . . .	17
1.3.2	Checkpointing with task-local I/O . . . . .	18
1.3.3	Performance-tools using task-local I/O . . . . .	20
1.4	Dynamic Linking and Loading . . . . .	21
1.4.1	Dynamic shared objects . . . . .	22
1.4.2	The dynamic loader . . . . .	23
1.4.3	Optimization of symbol binding . . . . .	24
1.4.4	Dynamic linking and loading in parallel applications . . . . .	25
1.4.5	Dynamic loading at runtime and within scripting languages . . . . .	26
1.5	Contribution of this Thesis . . . . .	27
<b>2</b>	<b>I/O Limitations at Large Scale</b>	<b>29</b>
2.1	Schematic View of the Parallel I/O Data Flow . . . . .	29
2.2	Scalability of Parallel Task-Local I/O . . . . .	31
2.2.1	Parallel file creation . . . . .	31
2.2.2	File management . . . . .	33
2.3	Scalability of Parallel Loading . . . . .	34
2.3.1	Scaling issues . . . . .	34
2.3.2	Dynamic loading on a Linux cluster . . . . .	36
2.3.3	Dynamic loading on JUQUEEN . . . . .	38
<b>3</b>	<b>SIONlib</b>	<b>43</b>
3.1	File Container for Task-Local Data . . . . .	43
3.2	Scalability of Shared-File I/O . . . . .	45
3.2.1	File locking (GPFS) . . . . .	45
3.2.2	Number of tasks per files . . . . .	47
3.2.3	Shared files and data size . . . . .	49
3.2.4	Shared files and threaded applications . . . . .	50

3.3	Objectives and Strategy . . . . .	51
3.4	Separation of I/O Streams . . . . .	53
3.5	File Organization . . . . .	54
3.6	Application Requirements . . . . .	57
3.7	Software Layers and APIs . . . . .	58
3.7.1	Parallel write . . . . .	59
3.7.2	Parallel read . . . . .	60
3.7.3	Serial write . . . . .	60
3.7.4	Serial read . . . . .	61
3.7.5	Fortran interface . . . . .	61
3.7.6	Commandline Utilities . . . . .	62
3.8	Coalescing I/O . . . . .	62
3.9	Key-Value Containers . . . . .	63
3.10	Support for Tools . . . . .	65
3.10.1	Generic API . . . . .	65
3.10.2	Reinit . . . . .	66
3.10.3	Mapped open . . . . .	67
3.11	Related Work . . . . .	68
<b>4</b>	<b>Evaluation of Task-Local I/O with SIONlib</b>	<b>71</b>
4.1	Architecture and I/O-infrastructure of Blue Gene/Q . . . . .	71
4.1.1	I/O-Forwarding on Blue Gene/Q . . . . .	72
4.2	I/O Benchmarks . . . . .	74
4.2.1	Parallel file creation . . . . .	75
4.2.2	Baseline measurements . . . . .	76
4.2.3	Alignment and file locking . . . . .	78
4.2.4	Shared file I/O with SIONlib . . . . .	80
4.2.5	Scalability . . . . .	82
4.2.6	Small data I/O at large scale . . . . .	84
4.3	File System as a Shared Resource . . . . .	86
4.4	Applications . . . . .	87
4.4.1	Checkpointing in MP2C . . . . .	88
4.4.2	Trace-file generation with Scalasca . . . . .	94
<b>5</b>	<b>Spindle</b>	<b>97</b>
5.1	Objectives and Overall Architecture . . . . .	97
5.2	Client Adapter . . . . .	99
5.3	Load Server . . . . .	100
5.3.1	Operation . . . . .	101
5.3.2	Memory overhead . . . . .	103
5.4	Overlay Network . . . . .	104
5.5	Bootstrapping and Bulk Preloading . . . . .	106
5.6	Python Module Loading . . . . .	106
5.7	Flexible Caching Algorithms . . . . .	107
5.7.1	Variable network topologies . . . . .	109

---

5.7.2	Optimization for Blue Gene/Q architecture . . . . .	110
5.8	Related Work . . . . .	111
5.8.1	Parallel file system . . . . .	111
5.8.2	Caching and staging solutions . . . . .	112
5.8.3	Peer-to-Peer solutions . . . . .	112
5.8.4	Loading as a parallel service . . . . .	113
<b>6</b>	<b>Evaluation of Parallel Loading with Spindle</b>	<b>115</b>
6.1	Simple Loader Benchmark . . . . .	115
6.2	Pyname . . . . .	116
6.3	Scaling Spindle on JUROPA . . . . .	117
6.4	Scaling Spindle on Sierra . . . . .	120
6.5	Memory overhead of Spindle . . . . .	123
<b>7</b>	<b>Conclusion &amp; Outlook</b>	<b>125</b>
	<b>Contribution to Publications</b>	<b>129</b>
	<b>Glossary</b>	<b>131</b>
	<b>Bibliography</b>	<b>133</b>

On current large-scale HPC systems often occur I/O patterns that produce a high load on the file system during access to checkpoint and restart files. Applications running on systems with distributed memory will often perform such I/O individually by creating task-local file objects on the file system. At large scale, these task-local I/O patterns impose substantial stress on the metadata management components of the I/O subsystem. Such metadata contention occurs also at the startup of dynamically linked applications while searching for library files.

The reason for these limitations is that the serial I/O components of the operating system do not take advantage of application parallelism. To avoid the above bottlenecks, this work describes two novel approaches which exploit the knowledge of application parallelism, the underlying I/O subsystem structure, the parallel file system configuration, and the network between HPC-system and I/O system to coordinate and optimize access to file-system objects. The underlying methods are implemented in two tools, SIONlib and Spindle, which add layers between the parallel application and the corresponding POSIX-based standard interfaces of the operating system, eliminating the need for modifying the underlying system software.

SIONlib is already applied in applications to implement efficient checkpointing and is also integrated in the performance-analysis tools Scalasca and Score-P to efficiently store trace data. Latest benchmarks on the Blue Gene/Q in Jülich demonstrate that SIONlib solves the metadata problem at large scale by running efficiently up to 1.8 million tasks while maintaining high I/O bandwidths of 60-80% of file-system peak with a negligible file-creation time. The scalability of Spindle could be demonstrated by running a benchmark on a cluster of Lawrence Livermore National Laboratory at large scale. The results show that the startup of dynamically linked applications is now feasible on more than 15000 tasks, whereas the overhead of Spindle is nearly constantly low.

With SIONlib and Spindle, this work demonstrates how scalability of operating system components can be improved without modifying them and without changing the I/O patterns of applications. In this way, SIONlib and Spindle represent prototype implementations of functionality needed by next-generation runtime systems.

This publication was edited at the Jülich Supercomputing Centre (JSC) which is an integral part of the Institute for Advanced Simulation (IAS). The IAS combines the Jülich simulation sciences and the supercomputer facility in one organizational unit. It includes those parts of the scientific institutes at Forschungszentrum Jülich which use simulation on supercomputers as their main research methodology.