

SW and HW Speculative Nelder-Mead Execution for High Performance Unconstrained Optimization

Artur Mariano
Institute for Scientific Computing
Technische Universität Darmstadt
Darmstadt, Germany
artur.mariano@sc.tu-darmstadt.de

Paulo Garcia, Tiago Gomes
Centro Algoritmi
Universidade do Minho
Guimarães, Portugal
paulo.garcia@algoritmi.uminho.pt,
tiago.a.gomes@algoritmi.uminho.pt

Abstract—This paper addresses the performance assessment of a new Nelder-Mead variant, that speculatively executes the *simplex* operations. This new variant was implemented as x86 parallel and sequential CPU versions as well as in handwritten and automatic C-to-RTL FPGA designs. As the execution flow is the same on every version, the efficiency of the synchronization by software and hardware is also accessed.

Performance trials of these versions were performed using (i) a last-generation FPGA and a last generation multi-core CPU-chip to run the software versions and (ii) relatively simple objective functions in \mathbb{R}^2 . Results show that performance of the handwritten hardware design is relatively equivalent to the sequential software version of the algorithm, even running at a much lower clock frequency (average of 1.9Mhz vs 3.4GHz). They also suggest that the synchronization methods employed to control the speculative execution are too expensive when managed by software, but efficient if managed by hardware.

I. INTRODUCTION

One of the most fundamental principles in our world is the search for an optimal state [1]. In applied mathematics and numerical analysis, this is often called optimization, i.e., the process of trying to find the best possible elements \mathbf{x}^* in \mathbb{X} in such a way that an objective function $F(\mathbf{x}^*)$ is either maximized or minimized, depending on the target goal. Factors such as discontinuity and multiplicity of both local maxima and minima increase the complexity of the problem.

One particular class of optimization is unconstrained optimization. The goal is to locate a minimizer \mathbf{x}^* of a given (nonlinear) function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. If f is nonsmooth or even discontinuous at some points in \mathbb{R}^n , the optimization method should only use the function values of f , since the derivatives of f may not exist for a particular point. Methods within this category are usually called Direct Search Methods (DSMs). Some of these derivative-free methods have been proposed in the past decades [2], including Spendley-Hext-Himsworth's method [3], Powell's method [4] and the Nelder-Mead algorithm [5], the focus of this paper.

The Nelder-Mead algorithm is one of the best known algorithms for multidimensional unconstrained optimization without derivatives [5]. Beyond solving the classical unconstrained optimization problem, it is also used to solve parameter estimation and similar statistical problems [6], [7].

As the Nelder-Mead algorithm becomes computationally expensive for high *simplex* dimensions and discontinuous

objective functions, high high performance versions are demanded by industry and science communities. This motivates the implementation and assessment of parallel versions of the Nelder-Mead algorithm for multi-core CPU-chips and FPGAs.

The goals of this paper are: (i) to implement and assess novel software parallel versions of the Nelder-Mead algorithm for shared-memory multi-core chips, (ii) to evaluate the suitability of re-configurable logic for this algorithm and (iii) to compare the two approaches, when considering relatively simple (although hard-to-optimize) functions in \mathbb{R}^2 . FPGAs are particularly suited for comparison since the execution flow of the software versions can be replicated in an FPGA design, therefore enabling the comparison of the efficiency of the fork-and-join synchronization by software and hardware. According to the best knowledge of the authors, there are neither disclosed parallel versions of the Nelder-Mead algorithm for shared-memory CPU-chips, nor implementations for FPGAs.

The main contributions of this paper are the following:

- The specification and assessment of a new parallel Nelder-Mead algorithm's version, which speculatively calculates the *simplex* basic operations in parallel;
- The implementation of the specified version, using two different methods for managing the fork-and-join mechanism that enables speculative execution;
- The assessment and comparison of the new parallel software Nelder-Mead's version with a sequential version, and two hardware designs, based on manual and automatic C-to-RTL synthesis.

The remainder of the paper is organized as follows. Section 2 presents the implementation details of the devised software versions whereas the hardware implementation is covered in Section 3. Section 4 presents the results of the performed trials. Section 5 describes some related work and Section 6 concludes the paper. Due to space issues, the algorithm is not described in the paper, but can be consulted in the original paper [5].

II. SOFTWARE IMPLEMENTATIONS

A. Execution-flow analysis

The Nelder-Mead algorithm tries to improve the *simplex* (and consequently its solution) at each iteration, based on the calculation of the centroid. The computational requirements of

the algorithm become particularly high when a large number of iterations are required for convergence. No iterations can be processed in parallel, nor can an iteration be parallelized itself, unless it consists of a *simplex* shrinkage. However, this operation is too cheap to be worth computing it in parallel.

After initializing the process, the algorithm calculates the centroid, the reflection vertex and the reflection vertex test (RVT), which assesses the reflection vertex and determines which operations must be applied to the *simplex*. When the reflection vertex is considered *good*, only the *simplex* is updated, which involves low overhead. When the reflection vertex is considered *weak*, an outside contraction is performed. For the *very weak* case, the algorithm performs an inside contraction. When the reflection vertex is considered *very good*, an expansion is done instead. Except when the reflection vertex is considered *good*, which requires no further work, every case consists of a basic operation and an associated test.

As the complexity of the contraction (regardless its direction) and the expansion require the same number and type of operations, their associated tests determine how computationally expensive each case is. In that regard, classifications are ordered as $good < very\ good \leq weak \leq very\ weak$, from the least to the most expensive. The *weak* case becomes more expensive than the *very weak* case, if a shrinkage is required in the *weak* case but not required in the *very weak* one. The exact cost of each iteration depends primarily on the objective-function and, at a lesser extent, on the stage of the *simplex*.

B. Devised versions

Two parallel versions were developed to calculate the basic operations in a thread-level speculative execution fashion (see e.g., [8]), in addition to a sequential version that follows the original description. All code is written in C and parallel versions are implemented with pThreads. Although more convenient, OpenMP does not provide a way of killing tasks, a functionality one of the proposed versions depends on.

Similarly to branch prediction in computer architecture, but for a broader number of cases, all basic operations are computed before the result of the RVT is known. In each iteration, the RVT calculation is overlapped with the parallel speculative computation of the basic operations, referred to as “decision paths”, i.e., *good*, *very good*, *weak* and *very weak*. This is expected to boost performance when several iterations are required to convergence, as long as each computation is executed on a different core. A process p runs the algorithm and manages thread creation, destruction and synchronization. As soon as the result of the RVT is known, one of the four computations is committed (process p might need to wait for its completion), whereas the others are discarded.

The implemented parallel versions differ on how they discard unnecessary computation. In both versions, four threads are created by process p - one per decision path - each maintaining a local, auxiliary, *simplex*. When created, threads immediately block on a private condition. The *simplex* used by the algorithm, i.e., by process p , is referred to as global *simplex*, also because it is visible to all threads.

At the beginning of each iteration, conditions are signaled and threads released, so they can copy the global *simplex* to their local *simplices*, with which they perform their respective operations. When threads are released, local *simplices* are coherent copies of the global *simplex*. At each iteration and depending on the result of the RVT, one thread is marked as *valid* and the remaining ones as *invalid*. Also at each iteration, process p copies the local *simplex* of the *valid* thread to the global *simplex*, which might require waiting for its completion.

The implemented versions, that follow this workflow, differ on how thread management is handled. The first version, referred to as “thread killing”, kills and creates threads at each iteration, whereas the second version, referred to as “persistent threads”, uses heavy synchronization to control the execution of each thread. Figure 1 shows the second version’s workflow.

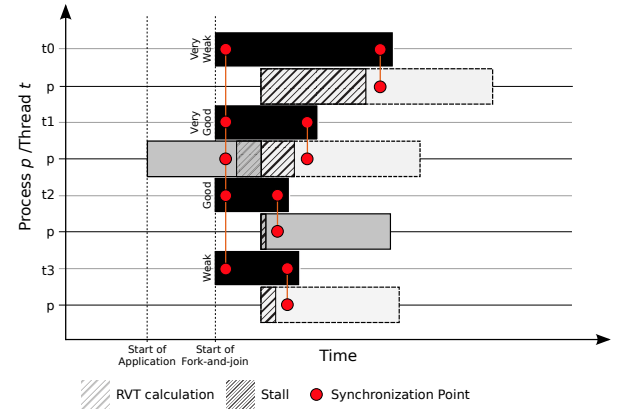


Fig. 1. Execution flow of the first iteration of a sample run in the “persistent threads” version, where the reflection vertex is considered *good*. Process p shown in light gray for remaining cases. Qualitative plot, not drawn to scale.

The second version, process p has a set of four pThread conditions in which it blocks as a means of waiting for the completion of a thread. Process p waits until the valid thread awakes it by releasing the condition. When this happens, the valid thread blocks itself on its condition until the next iteration starts (similarly to the “thread killing” version). As the remaining threads might still be computing their associated basic operations, process p resets them, by changing a variable read by every thread. Changing this variable is a lock-free operation. As a result, invalid threads might execute (a few) additional operations, very likely overlapped with the end of the iteration, which includes the stopping criterion calculation. Threads execute their respective operations on two nested loops where the outermost loop is executed indefinitely. When the “stopping” variable changes, threads break the inner-loop, which is executed again, since the outermost loop is infinite.

III. HARDWARE IMPLEMENTATION

A. RTL Nelder-Mead Implementation

The Nelder-Mead algorithm was implemented on a Xilinx Virtex-7 FPGA. Every variable (e.g., vertices and intermediate

values used for calculations) is implemented in both the 32 and 64-bit floating point representation defined by the IEEE 754-2008 standard, and all operations are performed using the embedded FPGA DSP blocks. The system implements input ports to set initial *simplex* values and output ports to display the result, as well as input/output control/status signals, e.g., start operation, calculation done. In order to maximize processing speed, the implementation attempted to parallelize computations as much as possible. Hence, the system is composed of four execution paths computing concurrently, in a similar fashion to the speculative execution in software versions. Common operations to all execution paths at the beginning of each iteration, such as vertex sorting, are executed by a single module that propagates results to the execution paths. The evaluation function multiplexes the computation path's output that is written back in the *simplex* registers. A simplified block diagram of the system is shown in Figure 2.

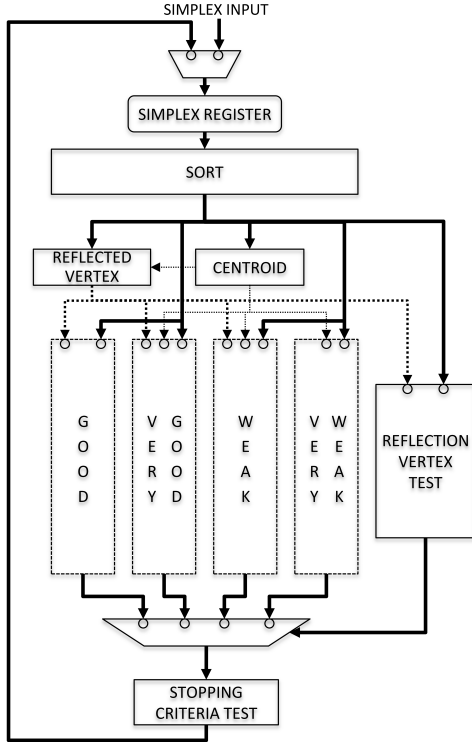


Fig. 2. Simplified block diagram of the FPGA Nelder-Mead design.

Table I presents the resource utilization rates for every implemented function. The design does not use BRAMs or DFFs. DSP blocks, on the other hand, are considerably used, especially in the 64-bit version, due to the wide number of 64-bit floating point arithmetic operators in the algorithm.

The design was implemented using Xilinx ISE 14.3 for design and implementation, including mapping, placing and routing, whereas Xilinx ISIM 14.3 was used for simulation.

Although the hardware designs run at a low frequency (< 3.6 MHz), the system performs one algorithm iteration per clock cycle, thus yielding very short execution times. The current operating frequency is due to the large number of

TABLE I
FPGA SYNTHESIS RESULTS. f REPRESENTS FREQUENCY.

Function	1	2	3	4	5	6
32-bit						
f (MHz)	3.591	3.582	3.730	2.821	2.380	2.701
Registers	<1%	<1%	<1%	<1%	<1%	<1%
LUTs	4.3%	4.5%	3.3%	9.0%	9.8%	9.6%
IOBs	40.1%					
DSPs	28.5%	24.8%	18.1%	27.0%	20.3%	30.7%
64-bit						
f (MHz)	2.147	2.148	2.123	1.880	1.398	1.713
Registers	<1%	<1%	<1%	<1%	<1%	<1%
LUTs	13.1%	15.0%	10.8%	16.9%	9.8%	18.4%
IOBs	80.2%					
DSPs	87.9%	59.0%	49.4%	87.9%	58.9%	97.5%

floating point arithmetic modules cascaded on the datapath. Using 64-bit precision also contributes to the low frequency; However, this design allows fair comparison between hardware and software implementations since they follow the same workflow. The 32-bit implementation results in a much more area-efficient implementation and provides moderate performance improvements. The smaller resolution showed equivalent results for five functions, failing only on one function, where for certain inputs, the difference in resolution yielded close, different results. Experiments up to date indicate that 32 bit resolution is sufficient for most function optimizations using Nelder-Mead, thus future work will encompass reducing resolution in order to further increase performance, including using alternative encodings for increased throughput [9]. Performance results were obtained through Xilinx's ISIM timing (post place and route) simulation.

B. Manual vs automatic hardware generation

The software version of the Nelder-Mead algorithm was translated to hardware through the Xilinx HLS C-To-RTL generation tool. This approach presented worse results than the *ad hoc* implementation, due to several reasons:

- The tool unrolls only a few loops. The majority of computations are implemented as multi-cycle hardware paths, probably due to the wide number of function calls on several loops' iterations.
- Albeit the C-To-RTL version yields much higher operating frequency than the manual implementation, at the cost of multi-cycle algorithm iterations, this results in no performance increase. This is due to the fact that subsequent Nelder-Mead iterations cannot be parallelized, therefore there is no gain from a pipelined implementation. The higher number of intermediate registers increases the setup-hold times in the critical path, resulting in an overall slower execution.

Software versions could be modified to achieve better results in automatic synthesis, using code re-factoring techniques encompassing hardware synthesis estimation [10], to be done as future work.

TABLE II
TESTED FUNCTIONS AND THEIR RESPECTIVE INPUTS.

Function	Input Simplex
$F_1(x, y) = -40000x - 60000y + 5x^2 + 10y^2 + 10xy$	$S = < (0.99, -0.34), (0.61, 1.39), (1.05, -1.895) >$
$F_2(x, y) = 20000 \times ((x + 70)^2 + (y + 275)^2) + y^2 + (y + 195)^2$	$S = < (234.55, 8.32), (23.343, 34.33), (0.992, 2.23) >$
$F_3(x, y) = 100(y - x^2)^2 + (1 - x)^2$	$S = < (0.081, 0.912), (92.2, 0.21), (18.11, 0.01) >$
$F_4(x, y) = x \times y \times 0.7623 + \frac{1}{4000} \times \prod_{i=1}^{i=2} x_i^2 - \sum_{i=1}^{i=2} \cos(x_i)$	$S = < (10.23, 0.16), (1.24, 0.7), (0.1, 0.1) >$
$F_5(x, y) = 418.9820 \times 2 - \sum_{i=1}^{i=2} x_i \times \sin(\sqrt{ x_i })$	$S = < (2.234, 0.832), (1.118, 304), (1.999, 0.354) >$
$F_6(x, y) = 20 + ((x^2 - 10\cos(2\pi x)) + y^2 - 10\cos(2\pi y))$	$S = < (20.667, 340.832), (1.2318, 200), (10.54, 0.7354) >$

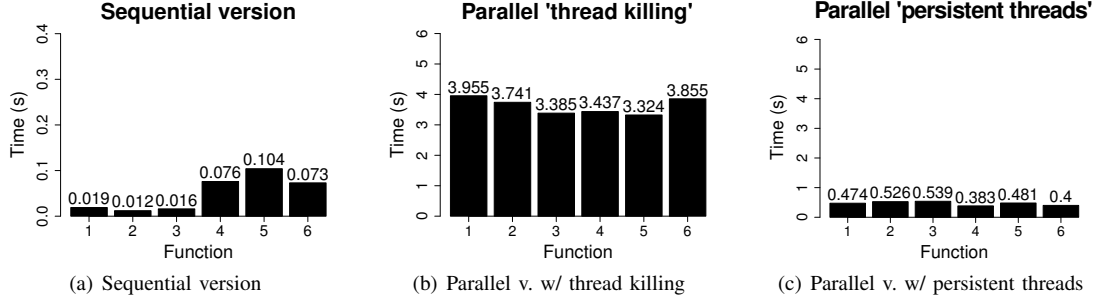


Fig. 3. Total runtime, in seconds, for the devised Nelder-Mead CPU versions, for the six presented objective functions.

IV. RESULTS

The developed software versions were tested on a last generation CPU whereas the hardware versions were simulated on a last generation FPGA. The characteristics of the tested platforms are summarized in Table III. Software versions were compiled with `gcc -O3`. The execution time was measured with the OpenMP `omp_get_wtime()` flag.

Six hard-to-optimize functions in \mathbb{R}^2 , presented in Table II and known for having multiple local optima, were chosen as case studies for benchmarking.

Table II also shows the initial *simplices* used for each function. In particular, F_3 is known as Rosenbrock's function, F_4 is a derivative of Griewank's function, F_5 and F_6 are respectively known as Schwefel and Rastrigin functions. The algorithm was limited to one hundred thousand iterations for every function and every version, both in the CPU and in the FPGA, the number of iterations taken by every function due to the strict stooeping criterion. The execution times presented for each trial (optimization of one function) on the CPU is the mean of five runs, whereas the hardware simulations are 100% accurate and have not been therefore statistically treated.

A. CPU

The run time of the implemented software versions is shown in Figures 3(a), 3(b) and 3(c). Both parallel versions are slower than the sequential one, due to the fork-and-join mechanism's overhead: either by creating and destroying threads in the first version or by thread synchronization in the second version. Not surprisingly, thread synchronization is more efficient than

TABLE III
TEST PLATFORM SPECIFICATIONS. IC AND DC STAND FOR INSTRUCTION AND DATA CACHE, RESPECTIVELY.

Device	CPU	FPGA
Manufacturer	Intel	Xilinx
Brand	Core 5 Ivy Bridge	Virtex 7
Model	i5-3570K	XC7VX485T
Max clock	3.4 GHz	300 MHz
Cores	4	-
System mem	16 Gbytes	68 Mbytes
L1 Cache	32kB iC+dC/core	-
L2 Cache	256kB/core unified	-
L3 Cache	6MB shared unified	-
Year	2012	2012

kill and create threads at each iteration, but it is still slower than the sequential version.

While trials with functions in \mathbb{R}^2 show that the overhead of the speculative execution mechanism kills any speedup, more complex functions and functions in higher dimensions, will mitigate the impact of the fork-and-join mechanism, highly noticeable since calculating the used functions is a very quick process. CPU versions are optimized, mostly due to `-O3 gcc` flag, and have less than 1% of L1 cache misses, as reported by *Cachegrind*, since the *simplex* data structures fit on cache.

B. FPGA

Figure 4 shows the results of the FPGA, both for the C-to-RTL version, in 4(a), and for the handwritten design in 4(b). The latter was implemented with Verilog, according to the description in Section III.

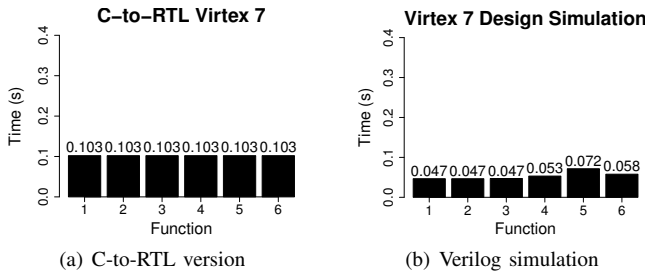


Fig. 4. Total runtime, in seconds, both for a C-to-RTL version and a 64-bit Verilog simulation of Nelder-Mead.

As shown in Figures 3 and 4, the FPGA implementation is more and less efficient than the equivalent software implementation, depending on the tested function. However, the FPGA runs at a considerably lower clock frequency (average of 1.9MHz of FPGA vs. 3.4GHz of the CPU). As the design replicates the parallel software versions, one can also conclude that the FPGA performance is not affected by synchronization issues, which are guaranteed by the used frequency.

V. RELATED WORK

A parallel version of the Nelder-Mead algorithm was proposed earlier, using parallelization at the parameter level [11]. However, the parallel version has a different search path though the parameter space than the non-parallel algorithm, in contrast to this paper. Moreover, the approach relies on even finer grained parallelism than the presented approach, thus likely unsuited for multi-core CPU-chips.

Both Nelder-Mead and Powell's methods were globalized and parallelized on a distributed memory environment with six single-core Pentium 4 machines running at 2.8 GHz, following a server-client fashion [12]. This paper, on the other hand, proposes parallel shared-memory and hardware implementations. A method for concurrent execution of the algorithm has also been proposed [13], with better results than the Nelder-Mead original algorithm on smooth, noisy, and functions with many local minima. This paper is rather focused on the original algorithm than on variants of it.

VI. CONCLUSIONS

This paper introduced a novel parallelization of the Nelder-Mead algorithm, to work on shared-memory multi-core CPU-chips. It is based on speculatively executing the operations to apply to the *simplex*, overlapping them with the RVT calculation to boost performance, by raising resource usage.

Trials with a 2D implementation of these versions, running on a quad-core CPU-chip, showed that even though speculative execution is applicable, performance is degraded due to (i) thread creation, destruction and synchronization costs to manage the fork-and-join mechanism that maintains the speculative execution and (ii) small computation overlap between the RVT and decision paths. As a result, this approach is not profitable, unless the objective functions take more time to compute, thus reducing the relative communication overhead.

Re-configurable logic has shown to deliver similar performance to the CPU, but at a much lower clock frequency. The FPGA design also calculates all the four decision paths in parallel, but the execution time is as long as the longest path. This suggests that synchronization for fork-and-join mechanisms would be considerably more efficient if implemented by hardware. The results of a C-to-RTL version showed that automatic conversion is less efficient for this particular algorithm, especially due to the low resource utilization, a common handicap of automatic synthesis.

The performance of both devices is strictly related with the characteristics of the performed trials. FPGA designs benefit from the use of both simple objective functions and small dimensions. While complex objective functions would favor the CPU's parallel proposed variant, due to bigger computation overlaps between the RVT and basic operations, higher *simplex* dimensions would require more area in the FPGA, already at $\approx 90\%$ utilization rates for DSPs, in some cases. At some point, area would be completely used and performance would be decreased. Benchmarks of both more complex objective functions and higher *simplex* dimensions are scheduled for future work, as well as the assessment of other methods to parallelize the algorithm by software.

REFERENCES

- [1] T. Weise, *Global Optimization Algorithms - Theory and Application*, 2nd ed. Self published, May, 2009.
- [2] R. M. Lewis, V. Torczon, and M. W. Trosset, "Direct Search Methods: Then And Now," *Journal of Computational and Applied Mathematics*, vol. 124, pp. 191–207, 2000.
- [3] W. Spendley, G. R. Hext, and F. R. Himsworth, "Sequential Application of Simplex Designs in Optimization and Evolutionary Operation," *Technometrics*, vol. 4, pp. 441–461, 1962.
- [4] M. J. D. Powell, "An efficient method for finding the minimum of a function of several variables without calculating derivatives," *The Computer Journal*, vol. 7, no. 2, pp. 155–162, January 1964.
- [5] J. A. Nelder and R. Mead, "A Simplex Method for Function Minimization," *The Computer Journal*, vol. 7, no. 4, pp. 308–313, January 1965.
- [6] T. G. Kolda, R. M. Lewis, and V. Torczon, "Optimization by direct search: New perspectives on some classical and modern methods," *SIAM Review*, vol. 45, pp. 385–482, 2003.
- [7] M. J. D. Powell, "Direct Search Algorithms for Optimization Calculations," *Acta Numerica*, vol. 7, pp. 287–336, 1998.
- [8] H. K. Pyla, C. Ribbens, and S. Varadarajan, "Exploiting coarse-grain speculative parallelism," in *Proceedings of the 2011 ACM international conference on Object oriented programming systems languages and applications*, ser. OOPSLA '11. New York, NY, USA: ACM, 2011, pp. 555–574.
- [9] F. de Dinechin, M. Joldes, B. Pasca, and G. Revy, "Multiplicative square root algorithms for fpgas," in *Field Programmable Logic and Applications (FPL) 2010*, 2010, pp. 574–577.
- [10] A. Cilaro, P. Durante, C. Lofiego, and A. Mazzeo, "Early prediction of hardware complexity in hll-to-hdl translation," in *Field Programmable Logic and Applications (FPL), 2010 International Conference on*, 2010, pp. 483–488.
- [11] D. Lee and M. Wiswall, "A Parallel Implementation of the Simplex Function Minimization Routine," *Comput. Econ.*, vol. 30, no. 2, pp. 171–187, September 2007.
- [12] A. Kosciński and M. A. Luersen, "Globalization and Parallelization of Nelder-Mead and Powell Optimization Methods," *Innovations and Advanced Techniques in Systems, Computing Sciences and Software Engineering*, pp. 93–98, 2008.
- [13] A. Lewis, D. Abramson, and T. Peachey, "RSCS: a parallel simplex algorithm for the Nimrod/O optimization toolset," in *Third International Symposium on Algorithms, Models and Tools for Parallel Computing on Heterogeneous Networks*, 2004, July 2004, pp. 71–78.