

Task Parallelization as a Service

A Runtime System for Automatic Shared Task Distribution

Luc Bläser

HSR University of Applied Sciences Rapperswil
IFS Institute for Software
lblaeser@hsr.ch

Abstract. The native parallelization support in predominant programming languages focuses on local cores, while mostly neglecting the huge parallelization potential of distributed execution. We have therefore developed a runtime system that automatically distributes ordinary shared memory parallel tasks into the cloud, to execute them on a large number of cores, e.g. on a cluster. The runtime mechanism takes care of transmitting the necessary task code and data to the service, as well as of propagating task results and side-effect changes back to the client memory. For programs with long-running tasks or a high amount of tasks, the system is able to soon compensate the network transmission overheads and thereafter scale with the number of tasks up to the available server-side cores.

1 Introduction

Today’s mainstream programming languages principally target only on local multi cores with their institutionalized parallelization features, as it fits well to their underlying shared memory model. However, moving from local to distributed parallelization, imposes the substantial extra burden for the developers to leave their conventional programming model and close the gap to distribution on their own, e.g. by redesigning the programs for specific distribution frameworks (e.g. service, remoting, or grid computing architectures) or by engaging a different more distribution-friendly programming paradigm (e.g. descriptive dataflow models).

It is therefore of no surprise that parallelization in daily programming practice largely concentrates on utilizing local processors, while mostly neglecting distributed parallelization - unless there is a strong performance urge justifying the efforts of a corresponding dedicated solution. Another strong distribution obstacle is the fact that developers and users normally do not have remote processor power at hands, although many clusters and distributed systems would have free capacities – but unfortunately without an easy-to-use “parallelization-as-a-service” interface.

There has been intensive research done in this area, principally going into three main directions: (1) new/different programming models inherently suited

for distribution, such as Actors/MPI [5, 1], or dataflow/query models [13, 3], (2) distributed task/thread frameworks, e.g. in the grid computing area [6, 10, 12], and (3), distributed shared memory systems [9, 2, 15, 4]. While the first direction takes the more radical approach of tackling the distribution impedance already at its roots, it usually compels programmers to apply the different paradigm on top of or aside their ordinary imperative shared memory programming language. The second direction, represented by the grid computing or distributed thread/task frameworks, typically leads to visible seams in the program design: it necessitates explicit task and data offloading, marking data serializable, wrapping data in specific sub-classes and so on. The third direction, distributed shared memory systems, transparently enables distribution of normal shared memory programs across machines. Although, it is usually heavy-weight for this purpose, operating on the whole program rather than selectively on the distributed parallel parts. A more detailed survey of related work is given in Section 5.

Our research goal is to significantly ease distributed parallelization by providing a new runtime system that allows scaling up parallel programs *seamlessly* on massive processor power in the cloud, i.e. with the same programming model as for local parallelization and without requiring any explicit code and data transmissions. For this purpose, we propose an enhanced thread pool mechanism that transparently integrates remote multi-processing: Although described like conventional local parallel tasks operating on shared memory, the runtime system automatically distributes the tasks to a web service into the cloud. Behind the service, the tasks are to be executed on a large number of cores before their results and their potential effected memory changes are finally sent back to the client runtime. The web service is of our design and can abstract an arbitrary parallel processing infrastructure behind the interface, e.g. a high-performance computer cluster. We have realized this system for the .NET framework, to demonstrate the concept by the example of a popular shared memory programming platform.

The remainder of this paper is structured as follows: Section 2 describes the programming model of distributed parallel tasks. Section 3 explains the design and implementation of the runtime system. Section 4 presents performance and scalability results. Section 5 compares distributed task parallelization to related work. Section 6 finally draws a conclusion of this work.

2 Programming Model

Task parallelization as a service encourages programmers to implement and start distributed parallel tasks that can be dispatched and executed on remote processors.

2.1 Distributed Tasks

In our system, which is based on .NET, distributed tasks can be programmed like conventional local thread pool tasks offered by the .NET task parallel library [8].

```

var distribution = new Distribution(ServiceURL, AccessCode);
var taskList = new List<DistributedTask<long>>();
foreach (var number in inputs) {
    var task = DistributedTask.New(
        () => Factorize(number)
    );
    taskList.Add(task);
}
distribution.Start(taskList);
foreach (var task in taskList) {
    Console.WriteLine(task.Result);
}

```

```

long Factorize(long number) {
    for (long k = 2; k * k <= number; k++) {
        if (number % k == 0) { return k; }
    }
    return number;
}

```

Fig. 1. Distributed parallel tasks factorizing numbers.

A distributed task is implemented as an ordinary .NET delegate¹ or lambda². In principle, working with distributed tasks remains analogous to using local parallel tasks, i.e. they can be instantiated, started, and joined. Certain restrictions apply for distributed tasks: Inner synchronization and calls to IO are for example forbidden. A detailed explanation of restrictions is provided in subsequent sections.

Figure 1 shows a code example for factorizing a set of numbers in parallel, each number being factorized as a separate distributed task. No extra compilation step is involved here; adding a reference to the library of our cloud task parallelization is sufficient for the runtime mechanism. The code sample looks very similar to the local task parallelization, as depicted in Figure 2. The URL and access authorization code need to be specified in advance for the remote task parallelization service, before a set of tasks can be started. Accessing task results blocks as long as the corresponding task is not terminated, where task faults are propagated as exceptions. We deliberately did not unify the local and distributed task class because we would like to encourage explicit combined starts of multiple distributed tasks for reducing network roundtrips, whereas the existing local task class promotes starting one-by-one.

For increased convenience, distributed tasks can also be applied in the form of data parallelism (Figure 3), by using parallel invocations or parallel loops. As for a parallel loop, each loop step starts a distributed task that executes the body and is joined again at the loop end.

Alike local tasks, distributed tasks are allowed to also perform side-effect changes on disjoint memory locations in shared memory, as also illustrated in Figure 3. The modifications in the array `outputs` of the example become

¹ A .NET delegate is a reference to a method and an associated object.

² A .NET lambda is an anonymous delegate in the form of an inline statement or expression with access (closure) to variables of its lexical scope.

```

var taskList = new List<Task<long>>();
foreach (var number in inputs) {
    var task = Task.Factory.StartNew(
        () => Factorize(number)
    );
    taskList.Add(task);
}
foreach (var task in taskList) {
    Console.WriteLine(task.Result);
}

```

Fig. 2. Analogous solution with local parallel tasks.

```

distribution.ParallelFor(0, inputs.Length, (i) => {
    outputs[i] = Factorize(inputs[i]);
});

```

Fig. 3. Distributed data parallelism.

automatically visible after task completion. For this purpose, the runtime system collects side-effect changes of tasks at the server side and propagates them back to the client-side memory. The system detects certain data races, as described in the next section.

2.2 Task Isolation

Distributed tasks are required to be independent of all other active tasks and threads, i.e. read-only accesses on shared variables and arbitrary accesses on non-shared variables are allowed. The granularity of variable accesses is per field or array element. Notably, this does not constitute a strong limitation because for local task code, synchronization in task execution is usually also avoided for highest possible performance. This applies for both synchronization primitives and memory model atomicity/visibility. The demanded task isolation eases the distribution significantly, since it excludes information flow between active distributed tasks, as well as, between active distributed tasks and the remaining program code. Our system detects certain violations of task isolation, namely when tasks employ synchronization, or when write-write conflicts happen due to data races. Read-write conflicts are not detected though: The reading task will not see the change of another concurrent task (snapshot isolation). In contrast, data races in local concurrency are not detected at all, while our system detects at least write-write conflicts.

2.3 Security Concerns

The runtime system prevents distributed tasks from directly or indirectly executing IO operations, system calls, reflection or unsafe/unmanaged .NET code. IO and system calls are not allowed because we do not delegate these calls back to the client, such that they would otherwise become effective on the remote machines. If reflection and unmanaged code would not be forbidden, programmers

could accidentally or intentionally inspect or modify arbitrary program state or corrupt memory safety at the remote side.

3 Runtime System

The system for distributed task parallelization consists of three components: a client runtime library, the cloud processing web service, and a server runtime library.

3.1 Processing Roundtrip

The processing of distributed tasks involves the following steps, as illustrated in Figure 4: (1) The potentially executed program code and accessible data of the invoked tasks are collected and serialized by the client side library at runtime. (2) The serialized code and data are shipped to our web service which represents the cloud processor resources. (3) The service distributes the tasks on server-side compute nodes, currently by launching a HPC cluster job consisting of a HPC task per input .NET task, i.e. by using the default task scheduling of the cluster. (4) The code and data are deserialized and instantiated by the server runtime library on each server compute node. (5) The remote tasks are executed on the compute nodes. (6) When terminated, the results and modified data of tasks are collected and serialized by the server runtime library. (7) The serialized data of task completion is sent back to the initiating client over the web service. (8) The updates are finally made effective in local memory of the client.

3.2 Task Serialization

When tasks are started for distribution, the client runtime component serializes the necessary code and data in two phases by way of reflection.

In a first phase, a conservative context-insensitive code analysis determines all reachable program code. Starting from the task delegate, the transitive closure of potentially directly or indirectly invoked methods is calculated. Additionally, it records all potentially used classes and accessed fields within the reachable methods. The code of each visited method is examined to only contain supported instructions and calls according to the security constraints (the system triggers a runtime exception if the code cannot be distributed). The set of reachable methods and their intermediate language code is eventually serialized.

In a second phase, all potentially accessed task data is collected. For this purpose, the system generates a partial heap snapshot, being the graph of objects that are reachable via references from the task delegate, by only considering the references occurring in potentially accessed fields according to the preceding code analysis. For the collected objects, only the state of accessible fields needs to be serialized. Besides the objects, the snapshot also includes static fields and constants that can be used by tasks. Because of the required task isolation, the runtime serialization delivers a consistent state without need of synchronization,

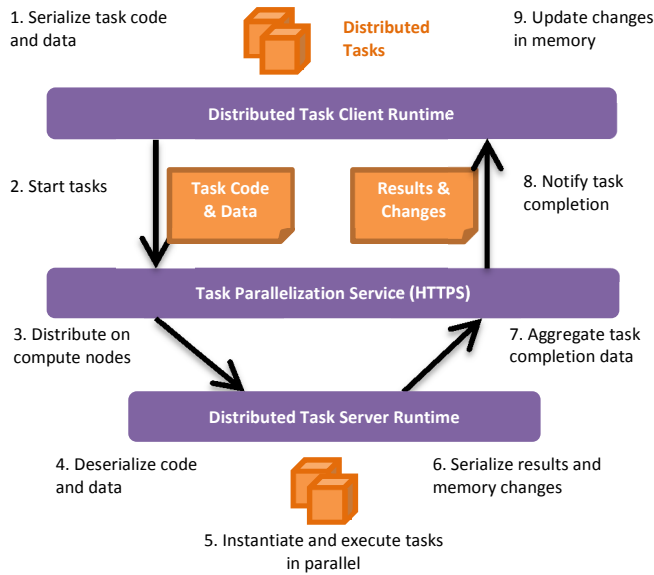


Fig. 4. Distributed task runtime system.

i.e. the system never blocks other running threads. Due to the conservative analysis, the snapshot may, however, include data that is not effectively accessed by the tasks and therefore also not required to be isolated: the state of this data may be inconsistent though but it is also never accessed by the distributed tasks.

3.3 Task Results

The server runtime component returns all necessary information of completed tasks, such as the task delegate result value, modifications on transmitted objects and static data (updates of fields and array elements), as well as, all reachable new objects that have been created by the remote task execution. The client in turn performs the in-place updates on arrival of the task completion information, i.e. modifications are applied to the corresponding objects and static fields of its input snapshot. We perform change detection by comparing the field and array element state before and after task execution. With this approach, the client runtime also detects certain data races, namely illegal write-write conflicts across distributed tasks. If remote task execution fails (e.g. due to a thrown exception), an aggregated exception is propagated to the client.

3.4 Service Design

Task code (program metadata and intermediate language code) and data (object graph and static fields) are encoded in an own binary format to reduce client-to-service traffic as much as possible. The service functionality basically comprises

two operations, one for starting a set of tasks and another for awaiting the termination of a set of tasks. To reduce network roundtrips, multiple tasks can be sent in one bunch, where the task instances can also share the same task code. To support secured network transmission, a HTTPS service binding can be used.

4 Experimental Results

Distributed task parallelization is intended for running computing-intensive tasks and/or a large amount of tasks, offering a high potential of parallelization.

4.1 Measurement Setup

For an experimental evaluation of our current system version, a set of synthetic parallel problems have been implemented on the basis of distributed tasks and eventually run in an environment with a MS HPC computer cluster behind the cloud service. The cluster comprises 32 nodes with 12 Intel Xeon cores, 2.6 GHz each (of which we were allowed to use 100 cores, 8 nodes with 12 cores plus 1 node with 4 cores for our experimental study). The client and web service each run on an Intel processor, 2 cores, 2.9 GHz machine, with 100Mbit/sec bandwidth and 1ms network delay between client, service and the cluster. All measurements have been performed by using compiler-optimized 64-bit .NET program assemblies. For all runtime results, the minimum of three repeated runtimes is considered, to reduce negative influences of temporary network speed fluctuations.

4.2 Performance Scalability

To study the performance scalability, we measure the runtime for a set of independent computation tasks. To start with a first scenario, we compute the factorization of a set of sample 64-bit numbers, each composed of two random prime factors around 2^{32} . Naturally, the same random seeds and numbers are used to ensure reproducible measurements. Each number is factorized independently in a parallel task. The comparison involves three processing approaches: (1) with distributed tasks, (2) with local tasks, and (3) sequential execution. Figure 5 shows the runtime in seconds depending on the amount of input numbers, which is equal to the number of parallel tasks. As expected, the parallel speedup of distributed tasks scales linearly with the number of tasks – in this scenario, each task runs on a separate instance of the 100 available cores. Local parallelization only offers a speedup of 2 on the two core client machine. Of course, the speedup also depends on the number of free cores available in the cluster.

4.3 Cost Breakdown

The runtimes for distributed task parallelization involve different performance cost factors, which vary from problem to problem: (1) the effective task execution time on the server side, (2) the network transfer time from the client over

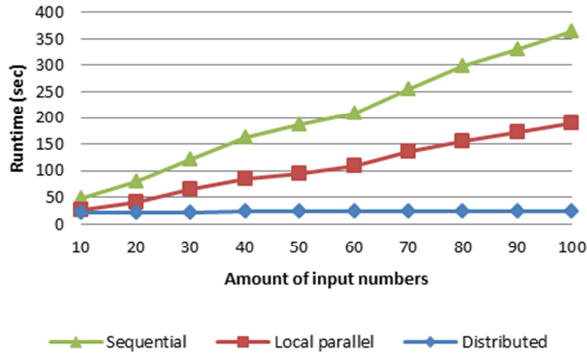


Fig. 5. Performance scaling by number of tasks.

Runtime costs (seconds)	Factorization (10 numbers)
Node execution	18.4
Network transfer	1.9
Cluster dispatching	0.3
Task serialization	0.3

Table 1. Runtime costs breakdown.

the service to the cluster, (3) the cluster dispatching costs, and (4) the accumulated effective overheads of our runtime mechanism, that is task serialization, deserialization, and change/result propagation. Table 1 depicts the breakdown of runtime costs in seconds for the factorization of 10 numbers, i.e. for 10 tasks. In this scenario, task execution represents the most significant part: this is the time where tasks are executed in parallel on the cluster. Network transfer constitutes the second most substantial portion. The remaining cost factors, including our own runtime mechanism, are relatively small.

4.4 Performance Comparisons

For a more general performance comparison, we evaluate the runtimes for different problem cases: (1) Mandelbrot fractal computation for a specified image size, as a representative of a parallel problem with a higher data amount compared to the task computation time. (2) Knight tours computation on a chess board of a specified size, as a representative for relatively long-running task computations. (3) Primes scanner counting primes in a specified number range, as a representative of relatively short-running tasks. Table 2 shows the runtimes in seconds, rounded to two significant figures, for the specified instances of these problems. We again compare distributed task parallelization (100 cores), local task parallelization (2 cores) and sequential execution. Once more, significant performance improvements can be achieved with the runtime support of dis-

Runtime (seconds)	Mandelbrot (10000 x 1000)	Primes Scanner (range 10 ⁷)	Knight Tours (6 x 6 board)
Distributed	8.0	4.7	120
Local parallel	20	9.2	1100
Sequential	37	19	2200

Table 2. Performance comparison of parallel problems.

tributed tasks. The involved data is around 1.25 MB for Mandelbrot (0.25 MB compressed input and 1 MB compressed output), while it is around 100 KB for knight tours and 80 KB for primes scanner.

4.5 Result Discussion

As expected, the examples confirm that the runtime system is able to reach a high parallel speedup by the large amount of remote cores. However, the gain of parallelization needs to compensate the involved overheads, which are primarily the network transmission time, depending on the size of task serialization, the data bandwidth and network delay. Distributed task parallelization is therefore generally beneficial if a large amount of tasks is executed, tasks are running sufficiently long, or tasks entail relatively low data transfer.

5 Related Work

5.1 Distributed Data Parallelism

Microsoft DryadLINQ [13], built on the distributed runtime engine Dryad [7], has close relation to our system: It permits automatic distributed processing of .NET LINQ [11] queries on clusters. While this model is oriented on descriptive programming of distributed dataflows in terms of queries, our system promotes more imperative task or data parallelization. Therefore, our system also allows distributed tasks to perform side-effects or changes that are propagated to the client, while the only backflow in DryadLINQ evaluations are the query results. Side effects of delegates inside the queries of DryadLINQ are ignored.

MapReduce [3] and in particular, also the Hadoop MapReduce implementation³, are popular dataflow programming models for high-scale distributed parallelization. The integration in a client program is, however, less seamless than in our model: Data is to be explicitly passed to the map and reduce functions from files or serializeable key-value sets. This is different to the shared memory illusion of our model, where the data of distributed tasks is automatically transmitted. MapReduce does also not directly incorporate a cloud approach where clients can easily offload their dataflows to a service. Though, such architecture can be designed around.

³ <http://hadoop.apache.org>

Other grid computing systems such as Pegasus⁴ and Swift [14] also facilitate DAG-like task workflow distribution on cloud computing resources but again with explicitly programmed data and task transmission. CIEL [12] supports powerful parallel task workflows with dynamic task spawning implemented in a specific language (Skywriting). Tasks can trigger batch commands, or invoke Java/.NET code in a less transparent way than our system: by denoting the class name and passing arguments and results.

5.2 Distributed Task Parallelism

Existing distributed thread/task programming frameworks, such as JPPF⁵, Hadoop, ProActive Parallel Suite⁶ [6], the already mentioned CIEL Skywriting [12], Alchemi [10], Manjrasoft Aneka .NET Tasks⁷, and many others, make distribution significantly more visible than in our system: heap data from the client program is not automatically shared across distributed processes but must be passed as explicitly serializable objects within task parameters and results, or has to be managed in specific grid heaps or distributed data collections. However, the focus of our system is on enabling mostly seamless and convenient task parallelization on remote processor resources.

5.3 Message Passing Models

The Actor model [5, 1] facilitates inherent distribution of active instances (actors) across machines, because actors only interact via explicit message communication and do not share memory. If applied within a conventional shared memory language (e.g. with MPI), this indispensably provokes a semantic gap, since programmers need to think in a different paradigm than the native language and have to stick to particular conventions. For example, actor communication must not be bypassed by ordinary references.

5.4 Distributed Shared Memory

Various systems have realized virtual shared memory on distributed computers, be it at the operating system level [9, 4] or at the runtime system of a programming language [2, 15]. While this can establish automatic distribution of an entire program, our system employs distribution only selectively for task parallelization. Moreover, we provide the distribution as a service for use by a possibly open group of clients.

⁴ <http://pegasus.isi.edu/>

⁵ <http://jppf.org>

⁶ <http://proactive.activeeon.com>

⁷ <http://www.manjrasoft.com>

6 Discussions and Conclusions

The presented runtime system enables seamless distributed task parallelization with the illusion of shared memory. While the programming model remains principally identical to working with local parallel tasks, the runtime engine automatically dispatches tasks over a service onto remote processor resources in the cloud. In contrast to other less seamless systems, this liberates developers from any distribution-specific programming artefacts, such as developing explicit remote code, realizing explicit communication, implementing any serialization, or wrapping/marking/attributing code or data for distribution-awareness.

Of course, distributed task parallelization is not appropriate for all classes of parallel problems. It is rather designed for computing-intensive tasks or a large amount of tasks, where it can achieve very high speedups. Thereby, the total task execution time has to be significantly larger than the network-dependent transmission time of task data between the client and the service.

We see a high potential if programmers can use “parallelization-as-a-service” in a way that is as simple and convenient as our task parallelization in the cloud.

Certainly, there is room for various improvements that we would like to address in future: (1) The runtime system could be enhanced to support more features, especially nested task starts, task chaining, task canceling, as well as, remote monitoring and debugging. (2) It could be investigated on alleviating task isolation by permitting well-defined synchronizations across tasks. (3) It would be interesting to offer a public parallelization service where users can directly consume and perhaps also offer multi-processor power on demand.

Acknowledgment

I gratefully appreciate the feedback from Peter Sommerlad, Svend Knudsen and Thomas Corbat, helping to improve this paper. I also express my thanks to Henrik Nordborg from the Microsoft Innovation Center for Technical Computing for offering access on the Microsoft HPC Cluster at the HSR during the experimental evaluations.

References

1. Agha, G.: *Actors: a model of concurrent computation in distributed systems*. MIT Press, Cambridge, MA, USA (1986)
2. Aridor, Y., Factor, M., Teperman, A.: *cJVM: a single system image of a JVM on a cluster*. In: *Parallel Processing, 1999. Proceedings. 1999 International Conference on*. pp. 4–11. IEEE (1999)
3. Dean, J., Ghemawat, S.: *MapReduce: simplified data processing on large clusters*. *Communications of the ACM* 51(1), 107–113 (2008)
4. Gupta, A., Ababneh, E., Han, R., Keller, E.: *Towards elastic operating systems*. In: *Proceedings of the 14th USENIX conference on Hot Topics in Operating Systems*. pp. 16–16. USENIX Association (2013)

5. Hewitt, C., Bishop, P., Steiger, R.: A universal modular ACTOR formalism for artificial intelligence. In: Proceedings of the 3rd international joint conference on Artificial intelligence. pp. 235–245. IJCAI'73, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (1973)
6. Huet, F., Caromel, D., Bal, H.E.: A High Performance Java Middleware with a Real Application. In: Proceedings of the Supercomputing conference. Pittsburgh, Pennsylvania, USA (Nov 2004)
7. Isard, M., Budiu, M., Yu, Y., Birrell, A., Fetterly, D.: Dryad: distributed data-parallel programs from sequential building blocks. *ACM SIGOPS Operating Systems Review* 41(3), 59–72 (2007)
8. Leijen, D., Schulte, W., Burckhardt, S.: The design of a task parallel library. In: *Acm Sigplan Notices*. vol. 44, pp. 227–242. ACM (2009)
9. Li, K., Hudak, P.: Memory coherence in shared virtual memory systems. *ACM Transactions on Computer Systems (TOCS)* 7(4), 321–359 (1989)
10. Luther, A., Buyya, R., Ranjan, R., Venugopal, S.: Alchemi: A. net-based enterprise grid computing system. In: *International Conference on Internet Computing*. pp. 269–278 (2005)
11. Meijer, E., Beckman, B., Bierman, G.: Linq: reconciling object, relations and xml in the. net framework. In: *Proceedings of the 2006 ACM SIGMOD international conference on Management of data*. pp. 706–706. ACM (2006)
12. Murray, D.G., Schwarzkopf, M., Smowton, C., Smith, S., Madhavapeddy, A., Hand, S.: CIEL: A universal execution engine for distributed data-flow computing. In: *Proceedings of the 8th USENIX Conference on Networked Systems Design and Implementation*. pp. 9–9. NSDI'11, USENIX Association, Berkeley, CA, USA (2011)
13. Yu, Y., Isard, M., Fetterly, D., Budiu, M., Erlingsson, Ú., Gunda, P.K., Currey, J.: DryadLINQ: A system for general-purpose distributed data-parallel computing using a high-level language. In: *Proceedings of the 8th USENIX conference on Operating systems design and implementation*. pp. 1–14 (2008)
14. Zhao, Y., Hategan, M., Clifford, B., Foster, I., Von Laszewski, G., Nefedova, V., Raicu, I., Stef-Praun, T., Wilde, M.: Swift: Fast, reliable, loosely coupled parallel computation. In: *Services, 2007 IEEE Congress on*. pp. 199–206. IEEE (2007)
15. Zhu, W., Wang, C.L., Lau, F.C.: Jessica2: A distributed java virtual machine with transparent thread migration support. In: *Cluster Computing, 2002. Proceedings. 2002 IEEE International Conference on*. pp. 381–388. IEEE (2002)