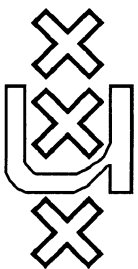


**Institute for Language, Logic and Information**

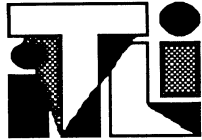
**MACHINE MODELS AND SIMULATIONS**

Peter van Emde Boas

ITLI Prepublication Series  
For Computation and Complexity Theory CT-89-02



University of Amsterdam



Institute for Language, Logic and Information  
Instituut voor Taal, Logica en Informatie

## MACHINE MODELS AND SIMULATIONS

Peter van Emde Boas  
Department of Mathematics and Computer Science  
University of Amsterdam

Received February 1989

To appear in  
J. van Leeuwen (ed.)  
*Handbook of Theoretical Computer Science*  
North Holland, Amsterdam, 1989

---

Correspondence to:

Faculteit der Wiskunde en Informatica  
(Department of Mathematics and Computer Science)  
Roetersstraat 15  
1018WB Amsterdam

or

Faculteit der Wijsbegeerte  
(Department of Philosophy)  
Grimburgwal 10  
1012GA Amsterdam

# Machine models and simulations\*

Peter van Emde Boas

*Departments of Mathematics and Computer Science, University of Amsterdam,  
Nieuwe Achtergracht 166, 1018 WV Amsterdam  
and CWI-AP6,  
Kruislaan 413, 1098 SJ Amsterdam*

## Abstract

In complexity theory we must cope with the discrepancy between a machine based foundation of complexity notions and a machine independent intuition on what computational complexity means. This discrepancy is resolved by investigating to which extent machine models simulate each other in an effective and efficient manner. Two theses are introduced which express the right kind of invariance needed. The invariance thesis states that reasonable sequential models simulate each other with polynomial overhead in time and constant factor overhead in space. The parallel computation thesis states that parallel time is the equivalent of sequential space.

Rather than using these theses as absolute truths we use them as heuristic tools by which two machine classes of well-behaved devices are separated from the more irregular devices. The first machine class consists of the sequential devices which obey the invariance thesis whereas the second machine class consists of the parallel devices which satisfy the parallel computation thesis.

Our survey of existing machine models takes us from the weakest universal models introduced by the mathematicians to the most powerful parallel models considered today where time no longer is a complexity measure.

---

\*to appear in: J. van Leeuwen (ed.), *Handbook of Theoretical Computer Science*, North Holland, Amsterdam, 1989

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	The invariance of complexity theory . . . . .	3
1.2	Formalization of machine models . . . . .	6
1.2.1	Machines and computations . . . . .	6
1.2.2	Simulations . . . . .	9
1.2.3	Complexity measures and complexity classes . . . . .	10
1.2.4	Simulation overheads and the machine classes . . . . .	13
1.2.5	Reductions and completeness . . . . .	15
<b>2</b>	<b>Sequential models</b>	<b>18</b>
2.1	Turing machines . . . . .	18
2.1.1	Description of the Turing machine models . . . . .	19
2.1.2	Simulation overheads . . . . .	20
2.1.3	Other Turing machine simulations . . . . .	22
2.2	Register machines . . . . .	24
2.2.1	The variety of models of register machines . . . . .	24
2.2.2	Time measures for <i>RAM</i> 's . . . . .	27
2.2.3	Space measures for the <i>RAM</i> . . . . .	29
2.2.4	The problematic simulation of a <i>RAM</i> on a Turing machine . . . . .	32
2.3	Storage modification machines . . . . .	34
2.3.1	Complexity measures for the <i>SMM</i> . . . . .	35
2.3.2	Simulations for the <i>SMM</i> model . . . . .	36
2.4	Networks and non-uniform models . . . . .	37
2.4.1	The network model . . . . .	39
2.4.2	Relation with Turing machine complexity . . . . .	40
<b>3</b>	<b>The second machine class</b>	<b>42</b>
3.1	<i>PSPACE</i> and transitive closure . . . . .	43
3.1.1	Transitive closure algorithms and <i>PSPACE</i> -complete problems . . . . .	43
3.1.2	Establishing membership in the second machine class . . . . .	46
3.2	The alternation model . . . . .	46
3.2.1	The concept of alternation . . . . .	46
3.2.2	Relation with sequential models . . . . .	48
3.3	Sequential machines operating on huge objects in unit time . . . . .	49
3.3.1	The <i>EDITRAM</i> model . . . . .	50
3.3.2	The <i>EDITRAM</i> is a second machine class device . . . . .	51

3.4	Machines with true parallelism . . . . .	53
3.4.1	The <i>SIMDAG</i> model . . . . .	54
3.4.2	The array processing machine . . . . .	55
3.4.3	Models with recursive parallelism . . . . .	56
<b>4</b>	<b>Parallel machine models outside the second machine class</b>	<b>58</b>
4.1	A weak parallel machine . . . . .	58
4.2	Beyond the second machine class . . . . .	59
4.2.1	The <i>MIMD-RAM</i> . . . . .	60
4.2.2	The <i>LPRAM</i> . . . . .	62
4.2.3	Extending the <i>SIMDAG</i> with powerful arithmetic . . . . .	62
4.2.4	Arbitrary computations in constant time . . . . .	63
	<b>Index</b>	<b>76</b>

# Chapter 1

## Introduction

### 1.1 The invariance of complexity theory

This chapter deals with the subject of machine models and simulations. Machine models are a necessary ingredient of any formal theory on computational complexity. If one wants to reason about such properties as the time and space consumed by algorithms solving problems, then one needs to stipulate what time and what space one is talking about. The conventional way of talking about time and space within theoretical computer science is based on the implementation of these algorithms on rather abstract machines, called *machine models*. The alternative would be to talk about the same measures but then applied to real life computers. This would make the results dependent on technology and real life peculiarities rather than insight and mathematical experience. Moreover, it can be argued that the models provide a reasonable approximation of what one might expect if real life computers were used indeed. Finally, one always can perform experiments on real computers if one wants to do so; by doing so one just stops being a theoretician and one becomes an empirical computer scientist instead.

Still, even if we base complexity theory on abstract instead of concrete machines the arbitrariness of the choice of the model used remains. It is at this point that simulations enter our theory. If we present mutual simulations of one model on another and give estimates for the time and space overheads suffered by performing these simulations we show in fact that complexity in these two models doesn't differ by more than this overhead in the relevant time and space measures. The sizes of these overheads provide information which indicates to which extent assertions about complexity are machine based and to which extent they become machine independent.

In order to appreciate the existence of many machine models we consider the state of affairs in general computation theory. In this theory we know a large variety of computing devices or formal calculi for effective computation. This divergence has not lead to a proliferation of computation theories due to the basic observation that the resulting formalisms are equivalent in the following sense: each computation in one formalism can be simulated in the other formalism as well. Computation theory originated from the need to show that specific problems were unsolvable by effective means. The above equivalence now suffices for providing the researchers freedom in choosing their model. For if one can prove that a problem is unsolvable in one particular model it is also unsolvable for all other formalized computing devices to which this particular model is related by a pair of

mutual simulations.

Another consequence of this state of affairs is that one can work in constructive mathematics without having to bother about the precise formal model of effective computability used. The existence of a single formal concept in many disguises has allowed us to return to our informal, intuitive style of working, relying on what has become known as the inessential use of CHURCH'S THESIS: *whatever is felt to be effective can be brought within the scope of our formal models*. The basic models themselves remain on the shelf to be used if when needed. They are also useful for educational purposes; their discovery has become history [23] [59].

Church's thesis is the kind of assertion of which the validity seems to be based on the lack of counterexamples. History has shown that all proposed models indeed are equivalent in the manner as explained above. Alternative models of computation are dreamt about in science fiction. I leave the question on whether the human brain can be modeled by one of our standard models or not to the philosophers and the neurophysiologists. In this chapter I restrict myself to models for which Church's thesis indeed is valid.

When in the early sixties the theory of computational complexity was founded [43] the subject of machine models was revived. Models which had been shown to be computationally equivalent became distinct again since they could be separated by their time and space behavior. Also a large collection of new models was introduced, based on the experience in the world of practical computing.

The classical models from recursion theory, like the unary Turing machine [22] [130] and Minsky's multi-counter machine [81] [119] turned out to be too unwieldy for modeling real-life computations. Depending on the nature of the objects one likes to manipulate in the computations (numbers (non-negative integers) or alphanumeric strings), two models have obtained a dominant position in machine-based complexity theory. The off-line multi-tape Turing machine [1] represents the standard model for string oriented computation and the random access machine (*RAM*) as introduced by Cook and Reckhow [21] has become the idealized von-Neumann number cruncher. But other models entered the theory as well and they were absorbed without much problems. At the same time during the early seventies theoretical computer science became aware of the fact that some complexity notions seemed to have a model independent meaning. The most important of these notions had to do with the characterization of *feasibility*.

Characterizing a problem to be feasible in case it can be solved in polynomial time (as was done for the first time by Edmonds [26]) requires that one identifies the machine on which this time is measured. On the other hand, if it can be shown that reasonable machines simulate each other with polynomial time overhead, it follows that the particular choice of model in such a definition of feasibility is irrelevant, as long as one remains within the realm of reasonable machine models. And indeed such polynomial time overhead simulations are abundant in the textbooks [1] [75] [138].

The fundamental problem which was discovered during this period is the  $P \stackrel{?}{=} NP$  problem, which asks for the power of nondeterminism. Are there problems which a non-deterministic machine can solve in polynomial time which can't be solved in polynomial time by a reasonable deterministic device? Stated otherwise: is it reasonable to assume that we have nondeterministic devices at our disposal?

The fundamental complexity classes like  $P$  and  $NP$  became part of a fundamental hierarchy, together with classes like  $LOGSPACE$ ,  $NLOGSPACE$ ,  $PSPACE$ ,  $EXPTIME$ , ...

And again theory faced the problem that each of these classes has a machine-dependent definition and that efficient simulations are needed in order to show that these classes are in fact machine-independent and represent fundamental concepts.

It seems therefore that complexity theory, as we know it today is based on the following assumption:

INVARIANCE THESIS: *Reasonable machines simulate each other with polynomially bounded overhead in time and constant factor overhead in space.*

A further development arose out of the study of parallelism. It turned out that there was a deep connection between time in the parallel world and space in the sequential world. This relation is known under the name PARALLEL COMPUTATION THESIS. The thesis states that whatever can be solved in polynomially bounded space on a sequential device can be solved in polynomially bounded time on a parallel device and vice versa. Using the notations which will be developed below this thesis is expressed by the equation  $PSPACE = || PTIME$ .

This chapter of the handbook of theoretical computer science deals with the evidence which has been gathered for the validity of the above two theses. As I mentioned before the study of simulations with polynomially bounded time overhead is a traditional subject in textbooks on introductory complexity theory. The additional clause on constant factor space overhead seems to have been taken for granted and/or self evident. However, as we will see in chapter 2 of this chapter this issue is not as evident as it seems to be. The situation with respect to the parallel computation thesis is again better; there are many results in the literature which show the thesis to be valid for some particular parallel model but on the other hand the thesis has been under attack as well.

For the interpretation of the invariance thesis it makes a difference whether one requires that a single simulation achieves both bounds on the overheads involved (orthodox interpretation) or whether one accepts a liberal interpretation which allows for a time-efficient simulation and an entirely different space-efficient simulation (which then may turn out to require an exponential overhead in time). As long as one investigates space or time bounded complexity classes independently the liberal interpretation suffices in order to make the classes in the fundamental hierarchy machine independent. Stockmeyer [126] has argued that one reason for using the Turing machine model in the foundation of complexity theory is precisely the invariance it provides. On the other hand nearly all efficient simulations known in the literature achieve both bounds at the same time. Since classes defined in terms of simultaneous time and space bounds are introduced and studied as well, I believe that the orthodox interpretation is the correct one. It is also evident that on the basis of the orthodox interpretation a subject like structural complexity theory [3] obtains a sound machine-independent basis.

The escape in defending these theses clearly is presented by the word *reasonable*. This can be easily inferred from the observed behavior of the theoretical computer science community when faced with evidence against the validity of the invariance thesis. For example, when in 1974 it was found that a *RAM* model with unit-time multiplication and addition (together with bitwise Boolean operations) became as powerful as a parallel device this model (the *MBRAM*) was thrown out of the realm of reasonable (sequential) machines, and it was considered to be a parallel device instead.

It seems therefore that the standard strategy is to adjust the definitions of reasonability when needed. The thesis becomes a guiding rule for specifying the right class of models



rather than an absolute truth. Therefore this thesis, once being accepted, will never be invalidated. This strategy is made visible if we exchange the word *reasonable* by some more neutral phrase. I have proposed for this purpose the use of the notions of *machine classes*. The **FIRST MACHINE CLASS** consists of those sequential devices which satisfy the invariance thesis with respect to the traditional standard device: the Turing machine. The **SECOND MACHINE CLASS** consists of those (parallel or sequential) devices which satisfy the parallel computation thesis [131] [132].

Not all devices we encounter in the literature are included in one of these two machine classes. The classical mathematical models like the unary Turing machine and the multi counter machine are too weak to keep in pace with the first machine class members. There are intermediate models which could be located in between the two classes, and there are parallel models which are even more powerful than the second machine class members, up to the models which accept everything in constant time. But so far I have not found reasons for introducing a third machine class at this level.

In the course of this chapter these devices will be introduced and their power will be explained and clarified. We will do so in a rather informal manner. Machine models are explained rather than defined formally, and their behavior is explained in words rather than in terms of computation sequences. An outline of the formal aspects of a theory of machine models will be given in the sequel of this chapter. Our treatment moreover necessarily will be incomplete. The reader who wants to obtain an impression on the extension of the field of machine models and related subjects is referred to the encyclopaedic volume written by Wagner and Wechsung [138], a book which had not yet appeared when the proposal for the present chapter was drafted.

## 1.2 Formalization of machine models

### 1.2.1 Machines and computations

A *machine model*  $M$  is a class of similarly structured devices  $M_i$ , called *machines*, which can be described as mathematical objects in the language of set theory. It is common to define these objects as a tuple of finite sets. For example, for the standard single tape Turing machine we will present a definition as a seven-tuple consisting of three finite sets, one finite relation and three elements in these sets.

The above set theoretical object provides only partial information on how the machine will behave and what its computations will look like. Common in these definitions is the presence in the tuple of a finite object, called *program* or *finite control*, which is supposed to operate on a potentially infinite structure called *memory*. One possible structure for the memory is a depository for finite chunks of information consisting of symbols from one or more finite sets called *alphabets*; these alphabets are specified by being members of the tuple which defines the machine. In an alternative structure the memory is filled with finite sets of numbers, which are in general taken to be elements from the set of nonnegative integers  $\omega$ . The memory is modeled as a regular structure of *cells* where this information is stored. This structure is mostly a linear and discrete order but alternative regular structures like higher dimensional grids or trees are also allowed.

Although the structure of cells modeling the memory is infinite in principle, it will always be the case that only a finite part of the memory is *used*. In the formalization this

is achieved by storing a special value (some selected *blank symbol* which is listed as such in the tuple, respectively the number 0) in those locations where nothing has happened so far.

The finite control is represented as a program in a suitable type of a rather primitive assembly code. Typically there will be instructions for fetching information from or storing information in memory, for modifying the visible cells in memory, for testing of conditions and performing conditional or unconditional jumps in the program, and for performing in- and output. In the example of the Turing machine model all these types of instructions have been fused into a single read-test-write-move-goto instruction; for models like the *RAM* the program does indeed resemble a primitive assembly code.

For the purpose of in- and output two special sections of memory are dedicated; one from which information is read one symbol at a time and one by which information is communicated to the outside world, again one symbol at a time. With regard to the input section of memory two interpretations offer themselves. One can regard the input characters as entities which offer themselves for being processed only once like signals received from another galaxy; if one wants to read the input symbol for a second time one must first store it in some other part of the memory. This leads to the mode of *on-line computations*. If on the other hand the input symbols are permanently available for inspection and can be read as often as one likes (presumably in some other order than the sequential order of the input string itself) one uses the mode of *off-line computations*. Some models, like for example the single tape Turing machine, have no separate in- or output devices at all; here the input is stored in memory at the start of the computation and the result can be retrieved from memory when the computation has come to an end.

In order to present a formalization of computations we first need the concept of a *configuration* of the machine. A configuration is a full description of the state of a machine and its memory. Typically a configuration will consist of a location in the program in the finite control, the structure and contents of the finite part of the memory which has participated in the computation so far, and descriptions of the states of the in- and output sections of memory. For all components of memory the configuration must moreover indicate on which cells the interaction between finite control and memory is currently taking place.

Next one introduces the so-called *transition relation* between two configurations. Configuration  $C_2$  is obtained by a transition from configuration  $C_1$ , notation  $C_1 \vdash C_2$ , if performing a single instruction of the program in the finite control of the machine takes us from  $C_1$  to  $C_2$ . Computations are described by the reflexive and transitive closure  $\vdash^*$  of this relation  $\vdash$ . The word *computation* refers primarily to the sequences of configurations, connected by  $\vdash$  which establish the presence of a pair consisting of two configurations in  $\vdash^*$ . The word may refer also to the pairs in  $\vdash^*$  themselves.

At this point the distinction between *deterministic* and *nondeterministic* machines can be made. For a deterministic machine there exists for every configuration  $C_1$  at most one configuration  $C_2$  such that  $C_1 \vdash C_2$ , whereas for nondeterministic devices there may exist several such configurations  $C_2$ . For virtually all devices in the literature the source of nondeterminism is the fact that the program itself indicates that more than one instruction can be performed in the given configuration, and this leads to *bounded nondeterminism*: the number of possible configurations  $C_2$  is finite and bounded by a number dependent on the machine's program only. Unbounded nondeterminism is obtained if also the interaction

with the memory can be nondeterministic. For example on a *RAM* one could consider an instruction which loads a random integer in a memory location. In this chapter we will restrict ourselves to bounded nondeterminism.

Given the transition relation one next introduces *initial*, *final*, *accepting* and *rejecting* configurations. Initial configurations are characterized by an initial state in the finite control (selected as such in the tuple describing the machine), a memory which is blank except for the input, an empty output memory, and all communication devices between program and memory located at some standard position. The input  $x$  completely determines an initial configuration for machine  $M_i$  which we denote by  $C_i(x)$ . A configuration  $C$  is called final if there exists no configuration  $C'$  such that  $C \vdash C'$ . If the state in the finite control in a final configuration equals a specific accepting state—which is designated as such in the tuple describing the machine—the final configuration is called accepting; otherwise the configuration is called rejecting. There exist alternative ways of defining configurations to be accepting or rejecting: for example, sometimes it is required that in an accepting configuration the contents in memory are reset to some "clean" situation.

A *full computation* is a computation which starts in an initial configuration and which does not terminate in a non-final one; either the computation is infinite in which case it is called a *divergent computation* or it *terminates* in a final configuration. The computation then is called accepting or rejecting depending on whether its final configuration has this property.

Machine computations can be used for several purposes. Machine  $M_i$  *recognizes* the language  $L(M_i)$  consisting of those inputs  $x$  for which an accepting computation starting with  $C_i(x)$  exists. Machine  $M_i$  *accepts* the language  $D(M_i)$  consisting of those inputs  $x$  for which some terminating computation starting with  $C_i(x)$  can be constructed. The machine  $M_i$  finally *computes* a relation  $F(M_i)$  consisting of those pairs  $\langle x, y \rangle$  which consist of an input  $x$  and an output  $y$  such that there exists an accepting computation which starts in  $C_i(x)$  and which terminates in a configuration where  $y$  denotes the contents of the output memory.

In the above formalization there exists only one finite control—or program—which is connected to the memory by possibly more than one device. As a consequence the machine will execute one instruction at a time (even in the case of a nondeterministic machine where several instructions can be performed only one of the possible instructions is chosen for execution). We emphasize this feature of the model by calling the above machine model a *sequential machine model*. A *Parallel machine model* is obtained if we omit in the above definition the condition that there exists just one finite control. In a parallel machine finite control is replaced by a set of processors the size of which may become infinite in the same way as memory is infinite: the number of processors has no fixed bound, but in every configuration only a finite number of them have been activated in the computation so far. Each processor has its private channels for interacting with memory. It can be the case that large amounts of memory are accessible to all processors (*shared memory*) or that processors have their own *local memory*.

The condition that only one symbol at a time gets communicated between the input or output section of memory and the finite control is relaxed for the case of parallel machine models. In parallel models processors also have the possibility to exchange information directly by communication channels without the use of intermediate passive memory. In an extreme case like the cellular automata model [128] there exist only processors and

their finite state controls have taken over the role of the memory.

Transitions for a parallel processor are the combined result of transitions within each of the active processors. In one formalization a global transition is obtained as the cumulative effect of all local transitions performed by the active processors which operate at the same time (*synchronous computation*). An alternative formalization is that processors proceed at their own speed in some rather obscure way (*asynchronous computation*). In both cases two or more processors can interfere with each other while interacting with shared memory; this holds in particular when one processor attempts to write at a location where another processor is reading or writing at the same time. It belongs to the precise definition of a parallel model to stipulate what will happen if those read/write conflicts arise, and whether these conflicts may arise at all during a legal computation.

### 1.2.2 Simulations

Given the above global description of machines and computations we next have to explain what we mean by a simulation of machine model  $M$  by a machine model  $M'$ . Intuitively a simulation of  $M$  by  $M'$  is some construction which shows that everything a machine  $M_i$  in  $M$  can do on some input  $x$  can be performed by some machine  $M'_j$  in  $M'$  on this input as well. The evidence that the behavior on input  $x$  is preserved should be retrievable from the computations themselves rather than just from the input/output relations established by the devices  $M_i$  and  $M'_j$ .

As it turns out it is rather difficult to provide for a more specific formal definition for this concept which at the same time is sufficiently general; as soon as one proposes a definition one finds examples of simulations which are not covered by this formal definition.

A first problem is, that it is quite conceivable that  $M'$  cannot process the input  $x$  at all, since the structure of the objects which can be stored in the input part of memories of the two machines may be quite different. If  $M_i$  is a Turing machine which operates on strings and  $M'_j$  is a *RAM* operating on numbers the sets of possible inputs for the two devices are disjoint. In order to overcome this hurdle one therefore allows that  $M'$  operates on some suitable encoding of the input  $x$  rather than on  $x$  itself. The encoding has to be of a rather simple nature, in order to prevent spurious interpretations which would allow for the recognition of non-recursive sets using simulations.

Next one could introduce a relation between configurations  $C_1$  of  $M_i$  and  $C_2$  of  $M'_j$  which expresses the fact that  $C_2$  represents the configuration  $C_1$ . Let us denote such a relation by  $C_1 \approx C_2$ . Ideally one would have for a simulation a condition like:

$$(C_1 \vdash C_3 \wedge C_1 \approx C_2) \Rightarrow \exists C_4[(C_2 \vdash C_4 \wedge C_3 \approx C_4)]$$

This kind of lock-step simulation turns out to be far too restrictive for incorporating existing simulations. For example, such simulations are doomed to preserve the number of steps in a computation, and it is therefore no wonder that many simulations presented in this chapter will violate this condition.

A next attempt is to require:

$$(C_1 \vdash C_3 \wedge C_1 \approx C_2) \Rightarrow \exists C_4[(C_2 \vdash^* C_4 \wedge C_3 \approx C_4)]$$

Still this is far too restrictive. For example it requires that each configuration of  $M_i$  has its analogue in the computation of  $M'_j$  and that the general order of the computation by

$M_i$  is preserved by the simulation. The first requirement will turn simulations into an asymmetrical concept which is not unreasonable. The second condition would however exclude such constructions as backward reconstruction, or recursive simulations where earlier configurations are re-computed many times. The famous simulation which proves Savitch's theorem [104] on the relation between deterministic and nondeterministic space bounded complexity classes represents a typical example of such a simulation and it would be a pity if this simulation would be excluded by a definition.

A next attempt will consider an entire computation of  $M_i$  and require that for each configuration  $C_1$  occurring in its computation there exists a corresponding configuration  $C_2$  in the simulating computation of  $M'_j$ . This extension would allow both backward and recursive simulations but still it would exclude simulations where the entire transition relation of  $M_i$  becomes an object on which the computation of  $M'_j$  operates and where the effect of  $M_i$  on some input  $x$  is simulated by letting  $M'_j$  evaluate the reflexive transitive closure of this transition relation. As we will see in chapter 3 and 4 this latter type of simulation will turn out to be quite common in the theory of parallel machine models.

The above attempts show how hard it is to define a simulation as a mathematical object and still remain sufficiently general. The initial definition seems to be the best one can provide, and I leave it to the reader to develop a feeling for what simulations can do by looking at the examples given in the sequel and elsewhere in the literature.

### 1.2.3 Complexity measures and complexity classes

Given the above intuitive description of a computation it is hardly a problem to define the time taken by a terminating computation: take the number of configurations in the sequence which describes this computation. What we have done here is to assign one time unit to every transition. This measure applies both to sequential and parallel models.

In models where a single transition may involve manipulation of objects from an infinite set (like in the case of a *RAM*) this approach (*uniform time measure*) may turn out to be unrealistic, and one uses a measure where each transition is weighted according to the amount of information manipulated in the transition; the *logarithmic time measure* for the *RAM* represents an example.

A crude measure for the space consumed by a computation is the number of cells in the memory which have been affected by the computation. A more refined measure takes also in account how much information is stored in such a cell. It is moreover tradition not to charge for the memory used by the input and the output, in case the corresponding sections of the memory are clearly separated from the rest of the memory. This convention is required for talking about functions which can be computed in logarithmic space, as is the case with the important notion of a logspace reduction. However, in structural complexity theory [3] the alternative, where output is charged is used as well; it leads to a theory which differs from the one obtained in the usual interpretation.

Having assigned a time and space measure to an individual computation, the next topic is to assign such measures to a machine. The way these measures are extended depends on what the machine is intended to do with the input: recognize, accept or compute. Moreover, one is interested not so much in the dependence of time and space on the precise input  $x$  but rather in a more global relation between time and space of computation and the length of the input  $x$  which is denoted by  $|x|$ .

The length  $|x|$  of the input depends on the encoding of  $x$  for the particular machine at hand. It is common that  $x$  represents some mathematical object like a number, a vector or an even more complex object like a finite graph. Traditionally one looks at codes where strings over a finite alphabet denote themselves, and where numbers are denoted using a binary or decimal notation (unless it is explicitly stated that a unary notation is used). Since we may assume in general that inputs are represented using an alphabet of at least two symbols, it is the case that the number of inputs  $x$  of length  $n$  increases as an exponential function in  $n$ .

**Definition 1** If  $f(n)$  is a function from the set of non-negative integers  $\omega$  to itself, machine  $M_i$  is said to terminate in time  $f(n)$  (space  $f(n)$ ) in case every computation on an input  $x$  of length  $n$  consumes time (space)  $\leq f(n)$ . The deterministic machine  $M_i$  accepts in time  $f(n)$  (space  $f(n)$ ) if every accepting computation on an input of length  $n$  consumes time (space)  $\leq f(n)$ . The nondeterministic machine  $M_i$  accepts in time  $f(n)$  (space  $f(n)$ ) if for every accepted input of length  $n$  an accepting computation can be found which consumes time (space)  $\leq f(n)$ .

**Definition 2** A set of inputs  $L$  is recognized in time  $f(n)$  (space  $f(n)$ ) by machine  $M_i$ , if  $L = L(M_i)$  and  $M_i$  terminates in time  $f(n)$  (space  $f(n)$ ). The set  $L$  is accepted in time  $f(n)$  (space  $f(n)$ ) by machine  $M_i$  if  $L = L(M_i) = D(M_i)$  and  $M_i$  accepts in time  $f(n)$  (space  $f(n)$ ).

An alternative definition is to let  $L$  be recognized by  $M_i$  in time  $f(n)$  (space  $f(n)$ ) in case  $L = L(M_i)$  for a machine which accepts in time  $f(n)$  (space  $f(n)$ ). This definition turns out to be equivalent to the one given for nice models and nice functions  $f(n)$  for which it is possible to shut-off those computations which exceed the amount of resources allowed. In the general case the definitions become different.

**Definition 3** The set of languages recognized by some machine model  $M$  in time  $f(n)$  (space  $f(n)$ ) consists of all languages  $L$  which are recognized in time  $f(n)$  (space  $f(n)$ ) by some member  $M_i$  of  $M$ .

We denote this class by  $M-TIME(f(n))$  ( $M-SPACE(f(n))$ ).

The class of languages recognized simultaneously in time  $f(n)$  and space  $g(n)$  will be denoted  $M-TIME\&SPACE(f(n), g(n))$ . The intersection of the classes  $M-TIME(f(n))$  and  $M-SPACE(g(n))$  consisting of the languages recognized in both time  $f(n)$  and space  $g(n)$  will be denoted  $M-TIME, SPACE(f(n), g(n))$ .

Standard functions for resource bounds in this theory are logarithms<sup>1</sup>, polynomials, exponential functions and their combinations. These functions involve machine or simulation dependent constants which are denoted by  $k$  in the sequel. In general we consider sequences of such functions, describing the generic behavior of still larger classes of orders of magnitude. We will look in particular to the following sequences  $\mathcal{F}$ :

$$\begin{aligned} Log &= \{k \cdot \log(n) | k \in \omega\} \\ Lin &= \{k \cdot n | k \in \omega\} \end{aligned}$$

---

<sup>1</sup>all logarithms in this paper are to the base 2

$$\begin{aligned}
Pol &= \{k \cdot n^k + k \mid k \in \omega\} \\
Expl &= \{k \cdot 2^{k \cdot n} \mid k \in \omega\} \\
Exp &= \{k \cdot n^k \mid k \in \omega\}
\end{aligned}$$

Next we define that a set  $L$  is recognized (accepted) in time (space)  $\mathcal{F}$  if it is recognized (accepted) in time (space)  $f$  for some  $f \in \mathcal{F}$ . For a machine model  $M$  this leads then to the following fundamental complexity classes:

**Definition 4**

$$\begin{aligned}
M\text{-LOGSPACE} &= \{L \mid L \text{ is recognized by } M_i \in M \text{ in space } Log\} \\
M\text{-PTIME} &= \{L \mid L \text{ is recognized by } M_i \in M \text{ in time } Pol\} \\
M\text{-PSPACE} &= \{L \mid L \text{ is recognized by } M_i \in M \text{ in space } Pol\} \\
M\text{-EXPTIME} &= \{L \mid L \text{ is recognized by } M_i \in M \text{ in time } Expl\} \\
M\text{-EXPTIME} &= \{L \mid L \text{ is recognized by } M_i \in M \text{ in time } Exp\} \\
M\text{-EXPLSPACE} &= \{L \mid L \text{ is recognized by } M_i \in M \text{ in space } Expl\} \\
M\text{-EXPSPACE} &= \{L \mid L \text{ is recognized by } M_i \in M \text{ in space } Exp\}
\end{aligned}$$

The above classes clearly are machine dependent. It is possible, however, to obtain the hierarchy of fundamental complexity classes as indicated in the introduction by selecting for  $M$  the model of standard off-line multi-tape Turing machines as described for example in [1]. Denoting the deterministic Turing machines by  $T$  and their non-deterministic counterpart by  $NT$  we obtain:

**Definition 5**

$$\begin{aligned}
LOGSPACE &= T\text{-LOGSPACE} & NLOGSPACE &= NT\text{-LOGSPACE} \\
P &= T\text{-PTIME} & NP &= NT\text{-PTIME} \\
PSPACE &= T\text{-PSPACE} & NPSPACE &= NT\text{-PSPACE} \\
EXPTIME &= T\text{-EXPTIME} & NEXPTIME &= NT\text{-EXPTIME} \\
EXPSPACE &= T\text{-EXPSPACE} & NEXPSPACE &= NT\text{-EXPSPACE}
\end{aligned}$$

From the elementary properties of the Turing machine model we now obtain the hierarchy

$$\begin{aligned}
LOGSPACE \subseteq NLOGSPACE \subseteq P \subseteq NP \subseteq PSPACE = NPSPACE \subseteq \\
\subseteq EXPTIME \subseteq NEXPTIME \subseteq EXPSPACE = NEXPSPACE \dots
\end{aligned}$$

Moreover, all the fundamental complexity classes mentioned above have obtained their standard meaning; note however the notational distinction between the classes described by exponential bounds with linear and polynomials in the exponent, and compare with the definitions used in the chapter by D.S. Johnson [56]. All inclusions represent notorious open problems in complexity theory, except for the equalities  $PSPACE = NPSPACE$  and  $EXPSPACE = NEXPSPACE$  which follow from Savitch's theorem [104]. The  $P \stackrel{?}{=} NP$  problem is included as the most notorious of these open problems.

### 1.2.4 Simulation overheads and the machine classes

Having defined the standard notions of complexity we return to the quantitative aspects of simulations. The unanswered question from chapter 1.2.2 *what does it formally mean to be a simulation* remains open but we now can at least quantify the simulation overheads by:

**Definition 6** We say that model  $M'$  simulates  $M$  with time (space) overhead  $f(n)$ , if the following holds:

for every machine  $M_i$  of  $M$  there exists a corresponding machine  $M'_{s(i)}$  of  $M'$  such that  $M'_{s(i)}$  simulates  $M_i$ , and such that moreover for every input  $x$  of  $M_i$ , if  $c(x)$  is the encoded input of  $M'_{s(i)}$  which represents  $x$ , and if  $t(x)$  is the time (space) needed by  $M_i$  for processing  $x$ , then the time (space) required by  $M'_{s(i)}$  for processing  $c(x)$  is bounded by  $f(t(x))$ .

Notation:  $M \leq M'(time\ f(n))$  or  $M \leq M'(space\ f(n))$

In this definition processing stands for either recognizing, accepting, rejecting the input or evaluating some (possibly partial and/or multi-valued) function on this input. As explained before the encoding  $c(x)$  must be recursive (and even of a reasonably low complexity). If  $s(i)$  is recursive as well the simulation is said to be effective. In abstract complexity theory [10], it is shown that simulations can be made effective, and that recursive overheads always exist.

If  $\mathcal{F}$  denotes a class of functions the above definition can be extended to stating that the model  $M'$  simulates  $M$  with time (space) overhead  $\mathcal{F}$  in case the simulation holds for some overhead  $f \in \mathcal{F}$ .

Notation:  $M \leq M'(time\ \mathcal{F})$  or  $M \leq M'(space\ \mathcal{F})$

This enables us to identify some important classes of simulations:

$M \leq M'(time\ Pol)$  *polynomial-time simulation*

$M \leq M'(time\ Lin)$  *linear-time simulation*

$M \leq M'(space\ Lin)$  *constant factor space overhead simulation*

Special cases of a linear time simulation are the so-called *real-time* and *constant-delay* simulation. These simulations behave according to one of the attempted definitions which we presented in chapter 1.2.2, where the configurations of  $M$  are represented by corresponding configurations of  $M'$  preserving the order; moreover the number of configurations of  $M'$  between the representatives of two successive representations of  $M$  configurations is bounded by a constant in the case of a constant-delay simulation; for the real-time case this constant equals 1. For machine models like the Turing machine which have a constant factor speed-up property it is possible to transform a constant-delay simulation into a real-time simulation. We denote the existence of such a real-time simulation by

$M \leq M'(real - time)$

If a single simulation achieves both a time overhead  $f(n)$  and a space overhead  $g(n)$  this can be expressed by the notation:

$M \leq M'(time\ f(n) \ \& \ space\ g(n))$

If both overheads can be achieved but not necessarily for the same simulation we express this fact by the notation:

$M \leq M'(time\ f(n),\ space\ g(n))$



Again these notations are extended to classes of functions where needed.

If simulations exist in both directions, achieving the same overheads we replace the symbol  $\leq$  by the equivalence symbol  $\approx$ . For example  $M \approx M'$  (*time Pol*), expresses the fact that the models  $M$  and  $M'$  simulate each other with polynomial time overhead.

The above machinery enables us to define what I mean by the *first machine class* and the *second machine class*.

**Definition 7** Let  $T$  be the model of Turing machines and let  $M$  be another machine model. Then  $M$  is a first class machine model if  $T \approx M$  (*time Pol & space Lin*). We say that  $M$  is a second class machine model if  $M-PTIME = PSPACE$ .

So first class machine models are exactly those models which satisfy the invariance thesis under the orthodox interpretation, whereas second class machines are those devices which make the parallel computation thesis true.

This chapter illustrates how we have turned these two theses from dogmatic specifications of the truth about machine models into a heuristic tool for classifying machines. It would be nice if all models which have been proposed in the literature would turn out to be either first or second class, but as it turns out that is not the case. On the other hand it would be disappointing if would be the case that, beside the Turing machines, there exist no other first class machine models. This would present strong evidence that our intuitive approach towards complexity theory where we are used to migrate freely from Turing machines to *RAM*'s and vice versa is unsound. Since we believe that what we are doing in our daily life when practising complexity theory is basically correct this is another unlikely result of our investigations. The truth turns out to be somewhere in the middle.

The first class machine models are precisely the machines for which the fundamental hierarchy:

$$\begin{aligned} LOGSPACE \subseteq NLOGSPACE \subseteq P \subseteq NP \subseteq PSPACE = NPSPACE \subseteq \\ \subseteq EXPTIME \subseteq NEXPTIME \subseteq EXPSPACE = NEXPSPACE \dots \end{aligned}$$

represents the complexity hierarchy for this model. Here the liberal interpretation of the invariance thesis would suffice also. However, for the models considered in this chapter it is the case that a model which satisfies the liberal interpretation of the invariance thesis for all modes of computation also satisfies the orthodox interpretation. Therefore we use the orthodox reading for the definition of the first machine class. For us the first machine class represents the class of reasonable sequential devices.

For the second machine class it is questionable whether these devices can be considered to be reasonable at all. It seems that the marvelous speed-ups by the parallel models of the second machine class require severe violations of basic laws of nature [136] [137]; stated otherwise, if physical constraints are taken into account, all gains of parallelism seem to be lost [17] [117]. So naturalness can never be a motivation for selecting precisely this particular class of parallel devices to become the representative one. Instead the motivation is drawn from the intrinsic stability of this class, its large number of members, and the variety of mechanisms by which its power can be obtained. And, as we shall see in the sequel, there exists at least one second machine class model which is a reasonable model of a device we use in our daily life [123].

### 1.2.5 Reductions and completeness

For lack of a better definition we have chosen a notion of simulation where one machine model simulates another model as soon as the second model can do everything the first model can do. We did stipulate when we presented this definition that this property should be inferred from the behavior of the machines and not just on the basis of the input-output relationship. If the latter condition is omitted we approach a notion which plays a crucial role in complexity theory: the notion of a *reduction*.

Even though the concept of a reduction as described for example in D.S. Johnson's contribution to this handbook [56] is defined in terms of the input-output behavior of machines only, it turns out that the study of reductions is a valuable tool for establishing the existence of simulations. In fact we will invoke in the study of parallel models frequently the argument that any device which solves some *PSPACE*-complete problem in polynomial time must embody at least the power of the second machine class models.

In this section I will present this technique and I will explain why it is not in violation of the constraints on simulations which I have stipulated in chapter 1.2.2.

First we need the notion of a reduction. This concept originates from the theory of recursive functions. There it deals with the relative computability of one set given another set for free. In complexity theory relative computability is exchanged against relative feasibility but the technical notions remain more or less the same.

**Definition 8** Let  $A$  and  $B$  be subsets of  $\Sigma^*$  and  $T^*$  respectively, and let  $f$  be some function from  $\Sigma^*$  to  $T^*$ . We say that  $f$  *reduces*  $A$  to  $B$ , notation  $A \leq_m B$  by  $f$ , in case  $f^{-1}(B) = A$ .

The above notion of a reduction is the well known *many-one reduction* from recursive function theory. There exist many more reducibilities but for our purposes this notion suffices. In the context of complexity theory the mapping  $f$  is subjected to the additional restriction that  $f$  can be evaluated in polynomial time (with respect to the length of the input). This additional condition is expressed by the notation  $A \leq_m^P B$ . If  $f$  satisfies the even stronger condition that  $f$  can be evaluated in logarithmic space one speaks about a *logspace reduction*.

In recursion theory the reduction  $A \leq_m B$  establishes the connection that  $A$  is recursive provided  $B$  is and consequently if  $A$  is undecidable then so is  $B$ . In complexity theory such a reduction  $A \leq_m^P B$  establishes the fact that  $A$  is easy in case  $B$  is, and consequently that  $B$  is hard provided  $A$  is hard.

It is not difficult to see that almost all fundamental complexity classes in the hierarchy

$$\begin{aligned} LOGSPACE \subseteq NLOGSPACE \subseteq P \subseteq NP \subseteq PSPACE = NPSpace \subseteq \\ \subseteq EXPTIME \subseteq NEXPTIME \subseteq EXPSPACE = NEXPSPACE \dots \end{aligned}$$

are closed under polynomial-time reductions. So we have for example the basic observation that  $(B \in P \wedge A \leq_m^P B) \Rightarrow A \in P$ . Exception with respect to this property are the classes *EXPLTIME* and *NEXPLTIME* which are not closed under polynomial time reductions; this is due to the fact that the composition of  $2^n$  with a polynomial  $p(n)$  yields a function  $2^{p(n)}$  which is a member of *Exp* rather than *Expl*. Only at the *LOGSPACE* and *NLOGSPACE* level one needs the stricter condition that the reductions are logspace reductions.

Given the notion of a reduction the concepts of hardness and completeness are defined by:

**Definition 9** Let  $B$  be some set in  $\Sigma^*$ , let  $\mathcal{X}$  be some class of sets and suppose that  $A \leq_m^P B$  for every  $A \in \mathcal{X}$ , then  $B$  is called  $\mathcal{X}$ -hard; if moreover  $B \in \mathcal{X}$  the  $B$  is called  $\mathcal{X}$ -complete.

The best-known instances of these notions are the concepts of  $NP$ -completeness and  $PSPACE$ -completeness. Cook's theorem [18] established in 1971 the  $NP$ -completeness of the set  $SATISFIABILITY$  whereas the existence of  $PSPACE$ -complete problems follows from [125].

It turns out that the original computational behavior of machines can be encoded in the specific decision problems which are shown to be hard or complete in the literature. This implies that some machine model  $M$  which can solve instances of such a complete problem in polynomial time can simulate all devices in the corresponding machine class.

In order to understand this observation we first observe that with hindsight the existence of complete problems for the fundamental classes in our hierarchy is easy to obtain. For example, a very simple  $NP$ -complete problem is given by the *universal* language:  $L_{NP} = \{\langle M_i, x, w \rangle \mid \text{nondeterministic Turing machine } M_i \text{ accepts } x \text{ in time } \leq |w|\}$ . If we replace time by space and consider deterministic instead of nondeterministic machines we obtain similarly a universal language  $L_{PSPACE}$  for  $PSPACE$ .

These universal languages are encoded by so-called *master reductions* into combinatorial problems. For example, in the proof of Cook's theorem [18] one constructs for every triple  $\langle M_i, x, w \rangle$  a propositional formula  $\Phi(M_i, x, w)$  such that a satisfying truth-value assignments to its propositional variables corresponds to an accepting computation of  $M_i(x)$  of length  $\leq |w|$ . From the solution of the combinatorial problem the original computation can be reconstructed. Other examples of master reductions are based on the use of tiling problems [65] [103] (a construction that works both at the level of  $NP$  and  $PSPACE$ ), but also Reif's ingenious construction of a skeleton which can only be removed from a closet in case some accepting  $PSPACE$ -bounded computation exists is an example of a master reduction [96].

Furthermore virtually all reductions which are used in proofs establishing hardness or completeness of specific problems have the property that they are *parsimoneous*: they transform not only solvable instances into solvable instances but also solutions into solutions. The property of master reductions that solutions encode accepting computations is therefore preserved under parsimoneous reductions. Finally the traditional complete problems have the property that one can compute efficiently a solution as soon as one is capable of establishing the existence of such solutions.

All these specific properties of the traditional complete problems together now yield the observation that solving some complete problem establishes the existence of simulations. From deciding the solvability of some combinatorial encoding of some machine computation one first computes the solution to the combinatorial problem and subsequently one obtains the original computation by decoding this solution. Therefore we obtain for example that, if  $M$  solves a  $PSPACE$ -complete problem in polynomial time then all polynomial space bounded machines can be simulated by  $M$  with polynomial time overhead. This observation will be used frequently in chapters 3 and 4 of this paper. It leads occasionally to much shorter proofs for the power of some parallel models than the proofs originally

given by the designers of these models.

## Chapter 2

# Sequential models

### 2.1 Turing machines

Turing machines have established themselves as the standard machine model in contemporary complexity theory. Only in the related area of analysis of algorithms this role is taken over by the *RAM* model. In the next section we present a formal definition of the basic version of the Turing machine model. I will abstain from giving formal mathematical definitions of the alternative variants of Turing machines. An informal explanation of the additional bells and whistles by which the basic model is extended will suffice in order to explain what these models are about. The additional features deal with the structure, dimension and number of tapes, the numbers of heads on a tape, restrictions on the use of tapes, or the availability of head-to-head jumps. Every particular selection of some particular set of available features leads to another version of the Turing machine model.

This basic machine model already shows a huge proliferation of subspecies. One might ask therefore why this model is so popular, and how it is possible that complexity theory as a whole seems not to suffer from the variability of this model.

It becomes clear at this point that I have still some unpaid debts. In chapter 1.2.4 I have defined the first machine class using Turing machines as a reference class, but, as indicated above, there exists no such standard model but an entire family of models. So I must first provide the reader with a specific model of the Turing machine which will serve as basis for the definition of the fundamental classes in chapter 1.2.4. For this version I follow the widespread convention of selecting the familiar version where all tapes are one dimensional, where one has one head on each tape, and where one has two special purpose tapes: one read-only input tape, and one write-only, one-way output tape; the remaining  $k \geq 2$  two-way infinite tapes are called work-tapes.

On the basis of this choice one must investigate whether the other Turing machine models satisfy the invariance thesis. These models must simulate each other with polynomial overhead in time and constant factor overhead in space; moreover the required simulations must achieve both overheads at the same time. Such simulations will make it clear that the concept of the first machine class indeed does not depend on the particular type of Turing model chosen.

Hence there remains just one troublesome question: Is the invariance thesis indeed true for all Turing models?

### 2.1.1 Description of the Turing machine models

In the Turing machine model the memory is structured in the form of a finite collection of tapes which consist of tape-cells each capable of storing a single symbol from the work-tape alphabet . The cells are arranged in the structure of a finite dimensional grid in Euclidean space. So cells have addresses which can be represented by integer vectors. For each tape there exist one or more communication devices connecting the finite control to the tape. These devices are called heads. Each head is positioned at some tape cell. Several heads may be located at the same cell.

In the program of the Turing machine the various types of instruction have been integrated into a single type: the machine reads the ordered set of symbols on the work-tapes which are currently scanned by the ordered sequence of heads. Depending on these symbols and the internal state of the finite control, the machine will overwrite the scanned symbols with new ones, move some (possibly all) of the heads to an adjacent cell on the corresponding tape and proceed to a new state in the finite control.

In the simplest version of the Turing machine the machine is represented by a tuple  $M = \langle K, \Sigma, P, q_0, q_f, b, \Delta \rangle$ . The machine has only a single one-dimensional tape, with tape alphabet  $\Sigma$ , and a set of internal states  $K$ . The program  $P$  of the finite control consists of a set of quintuples  $\langle q, s, q', s', m \rangle \in K \times \Sigma \times K \times \Sigma \times \Delta$ . Here the set  $\Delta = \{L, 0, R\}$  denotes the set of possible head moves : Left, stay put or Right . The meaning of this quintuple is: if in state  $q$  the head is scanning symbol  $s$  then print symbol  $s'$ , perform move  $m$  and proceed to state  $q'$ . The states  $q_0$  and  $q_f$  are two special elements in  $K$  denoting the initial and the final state respectively. The symbol  $b$  is a special tape symbol called blank which represents the contents of a tape-cell which never has been scanned by the head.

In this single tape model there is no special input or output tape. The input is written on the unique tape in the initial configuration with the unique head scanning the leftmost input symbol. If one wants this model to produce output one obtains such output by an ad-hoc convention from the final configuration (for example, the output consists of all non-blank tape symbols written to the left of the head in the final configuration).

If we denote configurations of the single tape machine in the format  $\$ \Sigma^* K \Sigma^* \$$ , with the state symbol written in front of the currently scanned tape symbol, the transitions between two successive configurations are described by a very simple context sensitive grammar. In this grammar one includes for example for the instruction  $\langle q, s, q', s', R \rangle$  the production rules  $(qst, s'q't)$  for every  $t \in \Sigma$ , together with the rule  $(qs$,  $s'q'b$)$  for the blank symbol  $b$ . Similar rules encode the behavior of left-moving instructions or instructions where the head doesn't move. This transformation of the program into a grammar has the properties that a one-one correspondence is established between computations of the machine and derivations in the grammar. This syntactic representation of the computations opens the way to encode Turing machine computations in a large variety of combinatorial structures, like tilings [42] [65] [103] or regular expressions [124], leading to a large collection of master reductions in computation theory. The popularity of the single tape model is, among other reasons based on this use in reductions and other simulations.$

In the more complicated models the set  $\Sigma$  is replaced by a power  $\Sigma^k$  where  $k$  denotes the number of tapes. The set of moves becomes also a  $k$ -fold Cartesian product, where moreover each set of moves is adjusted to the dimension of the tape the corresponding head is moving on. In case several heads are moving on a single tape a new coordinate

must be added in order to let the heads "feel" that they are looking at the same square: the coincidence pattern of the heads which can be modeled as a partition of the set of heads.

The coordinates of an instruction can always be split in an observation and an action part: in the above case the first state and symbol represent the observation and the rest represents the action. If a single observation can lead to at most one action the machine is deterministic and otherwise the machine is nondeterministic. Nondeterminism of Turing machines is bounded by definition since the set of possible actions is bounded by the size of the program which is finite.

Tapes being completely homogeneous, it is unimportant where in the initial configuration a head is positioned on a work-tape, provided the entire work-tape is blank, and that (in the case of many heads on a tape) all heads are located at the same cell.

Since it serves no purpose to read left of the first input symbol on a read-only input tape, and since the output tape (if present) is one-way it can be presumed that in- and output tapes are semi-infinite; their cells are indexed by the non-negative integers. The model can restrict work-tapes to be semi-infinite as well.

There exist several restricted types of work-tapes which have obtained special names: A *stack* is a semi-infinite work-tape with the special property that whenever the head makes a left move, the previous contents of the cell are erased (*popping the stack*). After a write the head can move right (*pushing a symbol on the stack*). The head can not move left of the origin but it can feel the bottom of the stack and in this way test whether the stack is empty. A *queue* is a semi-infinite tape with two right-only heads. The first head is a write-only head, whereas the second one is a read-only head. The second head therefore can read (only once) everything the first head has written before. Also a queue can be tested for emptiness. A *counter* is a stack with a single letter alphabet. Its purpose is to count a number  $n$  by storing a string of  $n$  copies of its tape symbol; the counter can test whether the number equals zero by checking whether the stack is empty. The push and the pop move of a counter now correspond to an increment and a decrement of the integer stored in the counter.

For higher dimensional tapes there are several possible choices for the set of possible head-moves. It makes a difference whether the heads can move only along the edges in the grid or whether they can proceed in diagonal steps as well. Other moves which have been considered are fast rewinds (a head proceeds in one step to its original position) or head-to-head jumps (one head moving in one step to the position of another).

### 2.1.2 Simulation overheads

Most results on simulations of enhanced Turing machines on more simple ones are well-known and documented in the literature. I will list a number of results of this type and make a few comments. I use the notation  $\leq$  introduced in definition 8 which must not be confused with the reducibility from chapter 1.2.5. I will omit in this list the trivial real-time constant factor space simulations which exist between a Turing machine model and a more enhanced version: adding new features will never increase time- or space consumption as long as the additional features are not used.

Rather than specifying the entire model, I will just mention the relevant features in the formulation of the overheads. Below I also provide a key reference.

**Theorem 1** *The extensions of the Turing machine model mentioned above can be simulated by more elementary features with the following overheads:*

1.  $1\text{-tape} \leq 2\text{-stacks}$  (real – time & space *Lin*)
2.  $1\text{-stack} \leq 2\text{-counters}$  (time *Expl* & space *Expl*)
3.  $1\text{-tape} \leq 2\text{-counters}$  (time  $2^{2^{k \cdot n}}$  & space  $2^{2^{k \cdot n}}$ )
4.  $m\text{-counters} \leq 1\text{-tape}$  (real – time & space *Log*)
5.  $m\text{-tapes} \leq 1\text{-tape}$  (time  $k \cdot n^2$  & space *Lin*)
6.  $m\text{-tapes} \leq 2\text{-tapes}$  (time  $k \cdot n \cdot \log(n)$  & space *Lin*)
7.  $\text{multi-headtapes} \leq \text{single-headtapes}$  (real – time & space *Lin*)
8.  $\text{multi-headtapes} + \text{jumps} \leq \text{single-headtapes}$  (real – time & space *Lin*)
9.  $2\text{-dimtapes} \leq 1\text{-tape}$  (time *Pol* & space  $k \cdot n^2$ )
10.  $2\text{-dimtapes} \leq 1\text{-tape}$  (time *Pol* & space  $k \cdot n \cdot \log(n)$ )
11.  $2\text{-dimtapes} \leq 2\text{-tapes}$  (time  $k \cdot n^{3/2}$  & space  $k \cdot \log(n)$ )
12.  $2\text{-dimtapes} \leq 1\text{-tape}$  (time *Pol* & space *Lin*)

**Proof:** The results 1) and 2) together imply 3) with 4 counters and a single exponential blow-up in time and space. But in fact two counters suffice (at the price of one more exponential blow-up). These results are due to Minsky and they imply that two-counter machines are universal [81].

Result 4), which is less well-known is due to Vitányi [135]; see also [118] for a tutorial exposé of this result. It is based on an oblivious simulation of a single counter on a single tape, where oblivious means that the position of the heads of the simulator depends on the number of steps in the computation so far only and not on the contents of the tape. The simulation uses a redundant number representation and an extremely ingenious method of simulating a recursive procedure without a stack. It is easy to see that—given an oblivious simulation of one counter—one can as well simulate  $k$  counters on the same tape, without loss of space or time.

Result 5) can be found in most textbooks [51], whereas 6) which is due to Hennie and Stearns [48] again is based on an oblivious simulation of a single tape on two tapes. This simulation uses a technique of tape segmentation where a tape is decomposed in regions of exponentially increasing sizes, which are located farther and farther away from the current head position and which are serviced with an exponentially decreasing frequency.

Result 7) originally was given by Fisher, Meyer and Rosenberg [30]. Later improvements require less tapes for the simulation (Leong and Seiferas [64]. The extension 8) with jumps has a long history but this problem was wrapped up finally by Kosaraju [61].

An overlooked case, as far as the invariance thesis is concerned is given by the simulations of multi-dimensional tapes on single-dimensional ones. The standard textbook simulations of two dimensional tapes on single dimensional tapes ( 9) and 10) ) require more than constant factor space overheads, in combination with time overheads  $k \cdot n^3$  and  $k \cdot n^2 \cdot \log(n)$ . The same is true for the time-efficient simulation 11) by Stoss [127] (also presented in [13]) which again is based on tape segmentation. The only reference known to me where the simulation 12) with a constant factor space overhead is given is



Hemmerling's report [47]; the only textbook I know which contains this result and this reference is the text by Wagner and Wechsung [138].

Simulation 12) is based on the so-called *history method*. The idea is to store address-content records but not using absolute but relative addresses related to the adjacent records on the simulated tape. This requires that the entire visited region of the higher dimensional tape is traced by a path which visits every cell, from every cell proceeds to one of its neighbors on the grid, and finally is not longer than a constant multiple of the number of cells visited. It should therefore not visit a cell more than a constant number of times. It turns out that a full traversal of a depth-first search tree in the subgraph of the grid which represents the visited area of the tape has these properties.

The simulations 9) -12) all extend to higher dimensions. The time overheads in 9),10) and 11) become  $k \cdot n^{d+1}$ ,  $k \cdot n^2 \cdot \log(n)$  and  $k \cdot n^{2-1/d} \cdot \log(n)$  respectively, whereas the space overhead in 9) becomes  $n^d$ . The space bounds in 10) and 11) and both bounds in Hemmerling's simulation are independent of the dimension.  $\square$

### 2.1.3 Other Turing machine simulations

An important property of the Turing machine model is the constant factor speed-up both in space and time. By replacing the tape alphabet  $\Sigma$  by its  $k$ -th power  $\Sigma^k$  one can compress  $k$  symbols into a single square. This makes it possible to design a machine which works  $k$  times as fast and uses  $k$  times less space. This transformation is a reverse of the simulations which prove that a tape alphabet with two symbols is universal, since in this simulation arbitrary symbols are encoded by binary strings.

Clearly the anomaly of the constant factor speed-up disappears if one assigns a weight to every step or tape-square of a Turing machine proportional to the amount of information processed. After compressing  $k$  symbols into one the manipulation of the symbol would become  $k$  times as expensive in a weighted time measure. Also the amount of information stored in a single symbol would be multiplied by the same factor  $k$ . There has been proposed a space measure called *capacity* which accounts for this anomaly [15]. On the other hand the constant factor speed-up is very convenient - one can get rid of all constant factors implied by  $O(f)$  terms.

The above constant factor speed-up for time and space is a well known property of Turing machines. What is less well known is the original proof given by Hartmanis and Stearns [46]. Most proofs you can find in textbooks start with a pre-processing stage where the input is condensed ( $k$  symbols into one); this pre-processing requires linear time provided the machine has two tapes; otherwise pre-processing requires quadratic time which leads to a constant factor speed-up in time which holds only for time bounds larger than  $n^2$ . Next the original computation is simulated, but in order to be certain that in a single move of the simulation one can simulate  $k$  steps of the original computation one must scan the neighborhood first and subsequently perform some updates to this neighborhood as well. This then leads to a speed-up where  $k$  moves are simulated by 6 moves. Hence, in order to achieve a speed-up by a factor  $k$  one must condense the tape by a factor  $6 \cdot k$  rather than  $k$ . In the original proof one keeps one block of the tapes in the finite control and this allows a  $k$  moves for 1 move simulation.

Another important question is whether overheads as indicated in the above theorem are optimal or not. Traditional folklore has the result that simulation of two tapes on a

single one requires quadratic time overhead, since palindromes can be recognized in linear time on two tapes, but for a single tape model time  $\Omega(n^2)$  is needed, as can be shown by a crossing-sequence argument. Note that this result does not imply anything on the optimality of a square overhead for simulating two work tapes by a single work tape or related problems. Such problems have been solved only recently using the technique of Kolmogorov complexity; see for example [66] [67] [73] [74] [87].

The following result expresses some well known relations between different resources:

**Theorem 2** *The Turing machine model satisfies the following inclusions*

1.  $T-TIME(f(n)) \subseteq NT-TIME(f(n)) \subseteq T-TIME(2^{k \cdot f(n)})$
2.  $T-SPACE(f(n)) \subseteq NT-SPACE(f(n)) \subseteq T-SPACE(f(n)^2)$
3.  $T-TIME(f(n)) \subseteq T-SPACE(f(n)/\log(f(n)))$

**Proof:** Relations 1) and 2) provide the known bounds on the relation between determinism and nondeterminism. For the time bounded classes the bounds are essentially trivial. The second inclusion in 2) represents Savitch's theorem [104]. For the case of linear time it is known that the first inclusion in 1) is a proper one [88]. If beside a time bound  $f(n)$  also a space bound  $s(n)$  is known Wiedermann [141] has recently improved the upper bound in 1) to:

$$NT-TIME \& SPACE(f(n), s(n)) \subseteq T-TIME(f(n) \cdot s(n)^2 \cdot 2^{k \cdot s(n)})$$

This also improves upon the time bound which results from the proof of Savitch's theorem:

$$NT-TIME \& SPACE(f(n), s(n)) \subseteq T-TIME(2^{k \cdot s(n) \cdot \log(f(n))})$$

Result 3) represents the unique result which indicates that space may be a more powerful resource than time, but beware: the proof given by Hopcroft, Paul and Valiant [50] is valid for one-dimensional tapes only. The result extends to higher dimensional Turing machines but then it requires a new proof [90]. There exist also versions for the *RAM* model [89] and for the Storage modification machine [41].

The proof of 3) involves several techniques: first the Turing machine computation is made block-respecting with a block-size  $k(n) \approx f(n)^{2/3}$ ; this means that both time and tapes are divided into segments of size  $k(n)$  such that segment boundaries are only crossed at the end of a time slot. The overhead for this part of the simulation is linear in time and space. Computing onwards for one block of time now only requires knowledge of the tape segments scanned at the beginning of a block period. The initial configurations of these blocks are collected into a computation graph which turns out to be a directed graph with indegree  $m + 1$  where  $m$  equals the number of tapes. A correspondence is made between solutions of a pebble game on this directed graph and space efficient simulations of the original computation. The fact that the graph of  $m := f(n)^{1/3}$  nodes can be pebbled with  $m/\log(m)$  pebbles now yields the logarithmic gain in space efficiency in 3) .  $\square$

For the case of single tape Turing machines , denoted *T1*, a better bound is known:

$$T1-TIME(f(n)) \subseteq T1-SPACE(\sqrt{f(n)})$$

provided  $f(n) \geq n^2$ ; see [86].

The result requires a constructibility condition on the timebound  $f(n)$ . Contrary to the Hopcroft, Paul and Valiant theorem the space-efficient simulation can be made polynomial time as well: Ibarra and Moran [52] have shown that

$$T1-TIME(f(n)) \subseteq T1-TIME \& SPACE(f(n)^2, \sqrt{f(n)})$$

The above result extends also to the case of a Turing machine with a two-way read-only input tape and one one-dimensional work tape. The simulation overheads become slightly different, due to the fact that the position of the input head must be stored. The space overhead there becomes  $\sqrt{f(n) \cdot \log(n)}$ . The time overhead remains polynomial  $f(n)^2$ , but the constructibility condition is still required.

Another difference with the Hopcroft, Paul and Valiant result is that the simulation extends to nondeterministic models as well. Ibarra and Moran [52] have shown that:

$$T1-NTIME(f(n)) \subseteq T1-NTIME \& SPACE(f(n)^{3/2}, \sqrt{f(n)})$$

Recently this result was further improved [70] [129]

$$T1-NTIME(f(n)) \subseteq T1-NTIME \& SPACE(f(n), \sqrt{f(n)})$$

The modified proof uses nondeterminism for a guess-and-certify strategy comparable to the strategies used in several efficient simulations of parallel models. It is interesting to combine the above result with the breadth-first simulation of Wiedermann. One obtains a result with a sublinear function in the exponent [141]

$$T1-NTIME(f(n)) \subseteq T1-NTIME(f(n)^2 \cdot 2^{k \cdot \sqrt{f(n)}})$$

## 2.2 Register machines

### 2.2.1 The variety of models of register machines

Register machines [21] have become the standard model in computation theory for the analysis of concrete algorithms. On the one hand these register machines can be recognized as the model of a standard computer which has been reduced to its minimal essential instruction repertoire. On the other hand these devices don't suffer from a doom hanging over all real computers: they have neither finite word length nor finite address space.

The popularity of the model seems to be based on the fact that, even though no real-life *RAM*'s exist, all models seem to be equipped with a highly efficient compiler for some ALGOL like language, which enables the authors of the various textbooks to use whatever high level language features they feel to be required for the description of their algorithms. It is clear that, in order to lead to realistic estimates on the complexity of the algorithms dealt with, some knowledge on the efficiency of the compiled code is presupposed but I won't go into the details of this issue. The model is sufficiently alike to real computers that we can rely on the collective experience of the designers of the compilers for real life programming languages.

In the basic *RAM* model the machine consists of a finite control where a program is stored, one (or more) accumulator register, denoted *Acc*, an instruction counter, and an infinite collection of memory registers  $R[0], R[1], \dots$ . The accumulator and the memory registers have an unbounded wordlength, but in any configuration a finite number will be stored in these registers.

The instruction repertoire of the *RAM* (and that of other machines we will meet in the sequel) can be divided in four categories:

1. instructions which influence the flow of control:
2. instructions for input and output: *read* and *print*.
3. instructions for transport of data between accumulator and memory.
4. instructions performing arithmetic.

Instructions of the first type are the unconditional instructions *goto*, *accept*, *reject*, *halt*, and various types of conditional jumps (*if condition then goto*, *if condition then skip*); here the condition is a simple test ( $Acc = 0?$  ,  $Acc \geq 0?$  ,  $Acc = R[i]?$  , ...). The meaning of these instructions is self-explanatory and they indicate themselves where the next instruction which will be executed is located in the program. For all other types of instruction the next instruction in the program is also the instruction which will be executed after completion of the previous one.

Depending on the precise model the transport instructions either can load entire integers to and from arbitrary registers, or these transports can require the accumulator as an intermediate storage location (*read Acc*, *print Acc*). An even more restricted form of transport is the model where a read instruction will input a single bit from the input channel and performs a conditional jump to a subsequent instruction label, depending on whether the bit equals 0 or 1.

The name *RAM* (*random access machine*) indicates the main feature in the instructions of the third type: the use of *indirect addressing*. There exist three types of load instructions and two types of store instructions:

$$\begin{array}{llll} \text{LOADD } i \text{ Acc} & := & i & \\ \text{LOAD } i \text{ Acc} & := & R[i] & \text{STORE } i \text{ R}[i] \quad := \text{Acc} \\ \text{LOADI } i \text{ Acc} & := & R[R[i]] & \text{STORE } i \text{ R}[R[i]] \quad := \text{Acc} \end{array}$$

In the weakest model the only arithmetical instructions is the Increment instruction:  $Acc := Acc + 1$ . The standard model has both addition and subtraction, where the second argument is fetched from memory:  $ADD \ j \quad Acc := Acc + R[j]$ . More powerful models extend this arithmetic with multiplication and division. For these powerful models it also becomes possible to treat the register contents as bit-strings rather than numbers; in this interpretation we have instructions for concatenation of bit-strings and bit-wise Boolean operations.

In the sequel we will denote by *SRAM* the model which only has the successor instruction; *RAM* denotes the standard model with addition and subtraction, whereas *MRAM* also has the multiplication and division. If the bitwise Boolean instructions are also available we introduce another B in the name (for example *MBRAM*). Prefixing with N denotes nondeterminism.

As with the Turing machines we have a huge collection of different models, but, as it turns out, also for every model at least two different methods of measuring time and space consumption of a *RAM* computation. In the *uniform measure* every instruction is counted as one step, regardless the size of the values operated on. In the *logarithmic measure* every instruction is given a weight equal to the sum of the logarithms of all the quantities involved in the instruction; this includes the lengths of addresses involved in direct or indirect addressing.

Nondeterminism can be incorporated in the model by legalizing a hideous error made by the beginning programmer: the multiple use of a label. If the machine jumps to a label which occurs twice in the program one of the occurrences of the label is chosen nondeterministically and the execution of the program continues there.

In real life computers a major breakthrough was von Neumann's idea to replace the specially wired program in the finite control by a program in memory which might even be modified by the computer during the course of its computation. A similar idea brings us to the *RASP* model [21] (*random access stored program*). Here the memory has been divided into the registers with even and with odd addresses. Two adjacent registers together store an instruction, with the operation code in the even addressed register and the operand address in its odd neighbor. By writing in such an odd location the machine can obtain the effect of indirect addressing without having it in its instruction repertoire. It has been shown that for the purpose of complexity theory the *RAM* and the *RASP* are fully equivalent (they simulate each other in real time with constant factor space overhead) and therefore I will not consider this *RASP* model in the sequel. The proof can be found in textbooks like [1].

It should be observed that even the weakest model (the successor *RAM*) remains universal as a computing device and that it even remains to be so when the instructions for indirect addressing are removed. In this case all addresses accessed during a computation are represented explicitly in the program. This model is equal to the *register machine* used as a model for basic recursion theory by Shepherdson and Sturgis [119]. The machine will use only a fixed finite set of its registers during the computation, and can be described as a finite control equipped with a fixed number of counters. Minsky has shown that two registers suffice for universal computing power [81]. On the other hand that very same result shows that the uniform space measure is not a proper complexity measure; it fails to obey the Blum axioms [10] which prohibit that with a fixed bounded set of resources you can compute everything.

The use of indirect addressing leads to two less desirable effects. It becomes possible to use registers in a very sparse way, leaving big gaps in the memory where nothing has happened. It also becomes less reasonable to presume that at the start of the computation all registers are properly initialized at zero; what if someone else has used the machine before you? If you have to initialize your memory who is going to pay for this? Both these problems have been considered and solved [1], but as we will see below it is not clear whether the solutions proposed which show that you can get around these problems with constant factor overheads in time and space are measure-independent.

### 2.2.2 Time measures for RAM's

As indicated above there exist two time measures for the *RAM* model. In the uniform measure every instruction is counted for one unit of time. In the logarithmic measure every instruction is charged for the sum of the lengths of all data manipulated implicitly or explicitly by the instruction. This length is based on a size function on numbers based on the binary logarithm with some ad-hoc clause for the small values 0 and 1:

$$size(n) = \text{if } n \leq 1 \text{ then } 1 \text{ else } \lceil \log(n) \rceil + 1 \text{ fi}$$

For example, the cost of an *ADDI j* instruction, with the effect  $Acc := Acc + R[R[j]]$  will be something like  $size(x) + size(j) + size(R[j]) + size(R[R[j]])$ , where  $x$  denotes the current value of the accumulator *Acc*.

In the sequel the use of uniform or logarithmic measure will be indicated by prefixing in complexity class denotations the words *TIME* or *SPACE* with U and L respectively.

It is clear that a computation in the uniform time measure costs less than in the logarithmic measure. The gap between the two measures may become as large as a multiplicative factor equal to the size of the largest value involved in the computation. How large this value can grow depends on the available arithmetic instructions. This leads to the following simulations:

**Theorem 3** *The uniform and the logarithmic time measure are related by*

1.  $SRAM-Utime \leq SRAM-Ltime(\text{time } n \cdot \log(n))$
2.  $RAM-Utime \leq RAM-Ltime(\text{time } n^2)$
3.  $MRAM-Utime \leq MRAM-Ltime(\text{time } Exp)$

The same inclusions hold for the corresponding nondeterministic classes. As can be seen the gap between uniform and logarithmic measure for the *MRAM* is so large that no simulation with polynomial time overhead can be guaranteed. Indeed, the *MRAM* in the uniform time measure has to be discarded from the realm of reasonable models: it is a member of the second machine class [9].

It is an elementary programming task to see that the bitwise Boolean instructions can be simulated in time polynomial in the length of the operands by the models which don't have them. This leads to the inclusions:

**Theorem 4** *The simulation overheads for the bitwise Boolean instructions are given by*

1.  $SBRAM-Utime \leq SRAM-Utime(\text{time } n \cdot \log(n))$
2.  $SBRAM-Ltime \leq SRAM-Ltime(\text{time } n \cdot \log(n))$
3.  $BRAM-Utime \leq RAM-Utime(\text{time } Pol)$
4.  $BRAM-Ltime \leq RAM-Ltime(\text{time } Pol)$
5.  $MBRAM-Utime \leq MRAM-Utime(\text{time } Pol)$
6.  $MBRAM-Ltime \leq MRAM-Ltime(\text{time } Pol)$

**Proof:** For 1) and 2) the standard tricks to translate between numbers and their bit-patterns don't work; one must explicitly store and manipulate the binary strings in arrays. The simulation 5) is not an invocation of the above observation but a rather deep result about the second machine class [9]: bitwise Boolean instructions are not needed for the power of parallelism in the presence of both addition and multiplication. The remaining simulations are straightforward and left to the reader. Again the results hold also for nondeterministic models.  $\square$

The absence or presence of multiplicative and parallel bit-manipulation operations is of relevance for the correct understanding of some results in the area of analysis of algorithms. Frequently for the analysis of concrete problems with low complexities it is assumed that the machine used can perform innocent looking multiplicative instructions on small values, since these instructions occur as well on real world computers. It is also common to use the uniform time measure. At the same time one is interested to obtain complexity bounds which are precise up to logarithmic factors. It is therefore relevant to know to which extent these complexity bounds depend on the precise instruction repertoire provided.

Assume by way of example that we are given the instructions which perform left- and right shifts of bit-patterns (multiplications and divisions by powers of 2). Assume moreover that we are allowed to perform these operations on arguments  $\leq n^k$ , where  $n$  denotes the largest value in the input and  $k$  denotes a fixed constant. It turns out to be possible to obtain a large collection of other multiplicative or bit-manipulation instructions at a time overhead  $O(1)$  in the uniform time measure by use of the basic technique of table look-up. Any operation can in principle be stored in a table of size  $n^{2k}$ ; the restricted multiplication allows us to simulate indexing in an  $n^k \times n^k$  two-dimensional array, and therefore the operator can be performed in time  $O(1)$ , once the table has been precomputed. But for many specific operations of a sufficiently "local" nature the same result can be obtained using a table of size  $O(n)$  (or even size  $O(\sqrt{n})$ ) by splitting the operands in pieces of length  $\log(n)/2$  or even shorter and reconstructing in constant time the entire result from the piecewise results by some suitable combining function which depends on the instruction. Since the table size becomes negligible compared to  $n$  so does the time required for precomputing its values. The parallel bit-wise logical operators and ordinary multiplication have this locality property, and so does the instruction which counts the number of 1-bits in a binary number.

This observation is a generalization of the strategy used by Schmidt and Siegel in [110]. It shows that there hardly exists such a thing as an "innocent" extension of the standard *RAM* model in the uniform time measure; either one only has additive arithmetic, or one might as well include all reasonable multiplicative and/or bitwise Boolean instructions on small operands at once.

The relation between the various *RAM* models and Turing machines can be obtained from a number of well known results found in the textbooks. There exist a number of improvements which are less well known. First the overhead for simulating a Turing machine on a *RAM*:

**Theorem 5** *The various RAM models simulate the Turing machine model with polynomial time overhead*

1.  $T \leq \text{SRAM-}U\text{time}(\text{time real} - \text{time})$

2.  $T \leq \text{SRAM-Ltime}(\text{time } k \cdot n \cdot \log(n))$
3.  $T \leq \text{RAM-Utime}(\text{time } k \cdot n / \log(n))$
4.  $T \leq \text{RAM-Ltime}(\text{time } n \cdot \log \log(n))$

The first of these simulations is due to Schönhage and results as a corollary of his results on Storage modification machines [114]. The second simulation is standard and can be found for example in the textbook by Aho, Hopcroft and Ullman [1]. Result 3) is an observation by Hopcroft, Paul and Valiant [49] which was not presented in the journal version [50] of that conference paper. The fourth simulation was given by Katajainen, Penttonen and van Leeuwen [58].

Next the overheads for the reverse simulations:

**Theorem 6** *Overheads for simulation of various RAM models on a Turing machine are given by*

1.  $\text{SRAM-Utime} \leq T(\text{time } n^2 \cdot \log(n))$
2.  $\text{RAM-Utime} \leq T(\text{time } n^3)$
3.  $\text{MRAM-Utime} \leq T(\text{time Exp})$
4.  $\text{SRAM-Ltime} \leq T(\text{time } n^2)$
5.  $\text{RAM-Ltime} \leq T(\text{time } n^2)$
6.  $\text{RAM-Ltime} \leq T(\text{time } n^2 / \log \log(n))$
7.  $\text{MRAM-Ltime} \leq T(\text{time Pol})$

The above results are obtained by allocating address/value pairs of the *RAM* on Turing machine work-tapes and estimating the time needed for processing these structures. Clearly most time in the simulation is spent by searching for a record on the linear tape. The first two simulations are due to Cook and Reckhow [21]. For simulation 3) there seems to be no alternative for an exponential time overhead simulation due to the possible growth of the length of the values involved. Results 4) and 5) again are due to Cook and Reckhow [21]. Simulation 6) is obtained by slightly rearranging the information contained in the *RAM* registers on the tape such that the shorter data are closer to the left end of the tape and are therefore easier to access. This result is due to Wiedermann [139]. Result 7) is easy.

If higher dimensional tapes are used the overheads are reduced to the effect that a factor  $n^2$  in the simulation overhead for a linear tape can be replaced by  $n^{1+1/d}$  where  $d$  denotes the dimension of the work-tapes used [138].

### 2.2.3 Space measures for the *RAM*

It is well known that *RAM* space should not be measured by counting the number of registers used. Minsky's result—a two counter machine has universal computing power [81]—shows that counting the number of *RAM* registers used during a computation does not yield a complexity measure satisfying the Blum axioms [10]. Still it is the case that this uniform space measure is commonly used in the theory of analysis of algorithms. This use can be justified on the basis that in concrete algorithms the intermediate results stored in



registers always are bounded in terms of the input values: either by a polynomial, or by a simple exponential function. Estimating the sizes of intermediate results therefore represents an essential part of the analysis of algorithms in the area of algebraic and number theoretical computations, since the above assumptions need justification.

In the theory of machine models every *RAM* register is charged for the size of its contents. Such a size function is invoked in the definition of the logarithmic time measure anyway, so why not use the same function for the space measure as well. But even then, having available such a size function, there are several ways to proceed. A rather crude way is to charge every register used for the size of the largest value produced during the entire computation. A more refined method is to charge every register for the largest value stored there during the entire computation. The latter heuristic leads to an expression:

$$space = \sum_{i=0}^{maxaddr} size_2(i, \max(i)) \quad (i)$$

where *maxaddr* is the index of the highest address accessed during the computation, and  $\max(i)$  is the largest integer ever stored in  $R[i]$  during the computation. In the above expression for the space measure I have added the additional parameter  $i$  to the size function  $size_2$  in order to make it possible to consider size functions which depend on the address of the register.

An even more refined method would be to look at individual *RAM* configurations. Charge every register used in some configuration for the size of the value currently stored there and compute the sum of these register costs for the registers which are currently used. The maximum of this sum, taken over all the configurations then would become the space measure. In this measure a copy of a large value which is repeatedly transported from one register to another one, after which the first register gets cleared is counted only twice, whereas in the measure expressed by the above formula  $i$  it gets charged for every register it visits. This difference, however, turns out to be non-problematic, and therefore I will restrict myself to space measures described by the above expression  $i$ . See Wagner and Wechsung [138] for several more alternative ways of defining space measures for the *RAM*.

The traditional *RAM* model supports the use of uninitialized storage: registers which have never been accessed before store the value 0. It is quite possible that there are registers with index in the range  $0, \dots, maxaddr$  which are never accessed during the computation, and therefore in the above formula something must be said about the space consumed by those unused registers. In a survey paper [105] W. Savitch considers the size function:

$$size_w(i, x) = \text{if } x \leq 1 \text{ then } 1 \text{ else } \lceil \log(x) + 1 \rceil \quad \text{fi} = size(x) \quad (ii)$$

The resulting measure charges unused registers for an amount of one bit for the value 0 stored there. As a consequence it becomes possible to consume exponential space during polynomial time by performing a program like:

```
addr := 1; for i from 1 to n
  addr := addr + addr;
  R[addr] := 1
```

od;

A more generally accepted measure (see for example [1]) uses the same size function for used registers but gives the unused registers for free:

$$size_s(i, x) = \text{if } R[i] \text{ is unused then } 0 \text{ else } size(x) \text{ fi} \quad (iii)$$

This measure solves the anomaly stated above but introduces a new problem which seems to have been overlooked in the literature: How to simulate a *RAM* on a Turing machine with constant factor overhead in space?

The standard trick of storing address-value records on a work tape requires additional space for the addresses, whereas a sequential allocation of all registers in the range  $0, \dots, maxaddr$  requires space proportional to Savitch's measure *ii*. So none of the two proposed measures leads to invariance of space in an evident way. Hence, taking the standard simulation of a *RAM* on a Turing machine, as a point of departure we arrive at the following size measure:

$$size_b(i, x) = \text{if } REG[i] \text{ is unused then } 0 \text{ else } size(i) + size(x) \text{ fi} \quad (iv)$$

We claim that  $size_b$  represents the intuitively correct way of measuring space on a *RAM*. It is however not the definition of the logarithmic space measure for *RAM*'s one encounters in the literature. Most textbooks provide either no formal definition at all, or they provide a definition based on  $size_s$ . This is most likely explained by the emphasis in these textbooks on time complexity. In the interesting case of a register access within an instruction using indirect addressing, the size of the address is charged for in the time measure. This is achieved by charging for the contents of the register used for the indirect addressing in the course of assigning a time measure to this indirection, but by doing so this charge has not become an intrinsic part of the cost the register which is reached by the indirection itself.

Basing the *RAM* space measure on  $size_b$  *iv* has several advantages. In the first place we observe that the constant factor space overhead for simulation of a Turing machine on a *RAM* remains intact, although the standard simulation which stores tape cells in consecutive registers has to be abandoned. This simulation would introduce an  $\Omega(S \cdot \log(S))$  space overhead due to the lengths of the addresses of  $S$  registers. But by using the standard trick of "one tape = two stacks", and by storing a single stack in a single register, a simulation with constant factor space overhead is obtained (at the price of increasing the time overhead by a factor  $S$  or  $S^2$  depending on the time measure used). So our proposal validates the invariance thesis.

Another advantage of the use of  $size_b$  is connected with the simulation of uninitialized storage as suggested in [1], exercise 2.12. When using  $size_s$  the space overhead becomes  $\Omega(S \cdot \log(S))$ , whereas it is a constant factor space overhead when using  $size_b$ . A similar observation can be made about the standard method of compacting sparsely used registers into a dense set by the creation of address-value pairs on the *RAM* itself.

Having chosen the right measure the space complexity of the *RAM* becomes fully equivalent to that of the Turing machine, and therefore I abstain from mentioning any specific further simulation results.

## 2.2.4 The problematic simulation of a *RAM* on a Turing machine

The reader might ask at this point whether the difference between  $size_s$  and  $size_b$  is really important. Surprisingly it is. As is shown in [121] [122] the invariance thesis, which is evidently valid if we use the measure based on  $size_b$ , becomes problematic if the definition based on  $size_s$  is used. Using the extremely complicated simulation which I will sketch below, it is shown that for deterministic off-line computations a Turing machine can simulate a *RAM* based on  $size_s$  with constant factor overhead in space. However, since this simulation requires exponential overhead in time, this is insufficient to show that the *RAM* model (again with the standard  $size_s$  space measure) is a first-class machine model! It can be shown that for on-line computations a constant factor overhead simulation does not exist, and the case of nondeterministic computations presents us with a wide open problem.

In this section I present a sketch of this complicated simulation. The straightforward simulation does not work here since no space is available for the address parts of the address-value records which must be stored on the work-tape of the Turing machine. On the other hand, the problem evaporates as soon as the size of the data stored in *RAM* memory exceeds or is just proportional to the space needed by these addresses. Therefore the problem will only arise in the unlikely situation where small chunks of data (say characters) are stored in registers with large addresses, encoding in this manner information about the address rather than about the data itself. This situation is illustrated by the following recognition problem:

$$L_0 = \{w_1 \# w_2 \# \dots \# w_k \# w_0 \mid w_i \in \{0,1\}^* \wedge w_0 \in \{w_i \mid 1 \leq i \leq k\}\} \quad (v)$$

It is not hard to see that on a *RAM* the above language can be recognized on-line by reading the words  $w_i$  as addresses, and storing a 1 in the corresponding registers. If we consider the typical case where the  $k$  words  $w_i$  in the input have length  $|w_i| = m$ , the above algorithm will consume space  $O(m + k)$  on a *RAM* with space measure based on  $size_s$ , whereas it requires space  $O(m \cdot k)$  if the measure based on  $size_b$  is used. The latter amount of space is a provable lower bound for Turing machine on-line recognition, since the Turing machine on-line acceptor must write a full description of its input on its work-tapes before it can process the last word  $w_0$ . Clearly this lower bound argument collapses if off-line processing is allowed, and it is not difficult to construct an off-line Turing machine acceptor which runs in space  $O(m + \log(k))$ .

The space measure based on  $size_s$  seems to underestimate the true space consumption when one uses a sparse table. If such a situation arises in actual computing, the problem may be solved by hashing techniques. Hashing will map a sparse set of logical addresses on a much denser set of physical addresses. Hence, rather than storing the address-value pairs directly, we represent the addresses by their hash codes for a suitable hash function, and reconsider the resulting space requirements. Two observations now can be made:

- In order for the computations not to be disrupted by confusion of registers the hash function used should be perfect with respect to the addresses used during the computation; it should be a 1-1 function if restricted to those arguments on which it will be evaluated.

- If the actual number of addresses used during the computation equals  $k$ , and if the hash function has a range  $c \cdot k$  for some constant  $c$ , the physical address-value records can be allocated sequentially on a work-tape, so the hash codes no longer need to be stored. Using a variable size format for the remaining value records, it can be achieved that inside this hash table an amount of space is consumed proportional to the space as measured by the *size*, space measure.

The above idea leads to the following problem on perfect hashing which must be solved: Given a set  $A$  of  $k$  elements of a Universe  $U = \{0, \dots, u - 1\}$  of size  $u$ . Find a perfect hash function  $f$  which scatters this set  $A$  completely into a hash table of size  $c \cdot k$  for some fixed constant  $c$ .

Functions satisfying the above requirements evidently exist as set theoretical objects. But for the problem at hand, these functions themselves should satisfy some space requirements as well. In the above situation, the *RAM* to be simulated may quite well consume no more space than  $O(k + \log(u))$ , and consequently, the manipulation of the hash function should not require more space than  $O(k + \log(u))$  as well. Otherwise the advantage gained by its use will be lost during its manipulation. So the final constraint on the above function  $f$  reads:

- the hash function  $f$  should be of program size  $O(k + \log(u))$  and require space  $O(k + \log(u))$  for its evaluation.

Hash functions satisfying the first two requirements were constructed in a very explicit manner by Fredman, Komlós and Szemerédi [34], but their construction did not satisfy the last requirement. Mehlhorn [78] showed that for the program size the upper bound could be met and even be improved to a tight  $\Theta(k + \log \log(u))$  bound, but his function required far more space for evaluation. In [121] [122] C. Slot and the present author succeeded in matching the Mehlhorn bound in combination with an  $O(k + \log(u))$  evaluation space. The required space-efficient logical-to-physical address translation therefore exists. See also [35] [53] [79] and the more recent paper by Schmidt and Siegel [110] where perfect hash functions are constructed which achieve the Mehlhorn bound and can be evaluated in unit time on a *RAM*.

The next question is how the simulator can ever hope to find a suitable perfect hash function. Since we are looking for a deterministic simulation the method of guessing a hash function and using it is not allowed. One has to try out a large collection of hash functions and select a good one. Being good here means being perfect with respect to the set of addresses in the simulated computation. Whether a hash function is good must be certified. In the process of certification it is not possible to make a list of all addresses encountered and keeping track of possible collisions under the hash function, since this will re-introduce the space consumption we are trying to eliminate. However, it is possible within the available space bounds to check for every slot in the range of the hash function that there exists at most one address which is hashed onto this slot. Since the computation is deterministic this process can be repeated for every slot individually. This explains why the simulation breaks down for both probabilistic and nondeterministic modes of computation. Under these modes of computation one never knows whether a second run of the computation will access the same registers as the first run.

As long as the hash function has not obtained its certificate of being perfect, the simulation has to cope with the possibility that the hash function is in fact not perfect

and that the simulator therefore may confuse registers. It may provide wrong answers, it may consume far more space than it should, and it may even diverge on bounded memory. The first problem is not a serious problem since the answer of the simulation will be not be trusted before the hash function has been certified, and the second problem can be solved using Savitch's trick of incremental space [104]. It is the third problem which causes the greatest trouble. Since the space consumption  $O(k + \log(u))$  may turn out to be  $o(\log(n))$  where  $n$  denotes the input length, detection of loops on bounded storage by counting is not allowed. Luckily this problem has been solved already by Sipser [120], who proposed a simulation by backward search from the accepting configuration, and by this simulation loops are prevented and no extra space is consumed.

Details of the above simulation and the theory on perfect hashing on which it is based can be found in [122]. I hope that the sketch above suffices to convince the reader that the simulation is quite complicated, requires a lot of re-computations, causes an exponential time overhead, and fails to solve the problem for alternative modes of computation. As such it provides no evidence that the invariance thesis is satisfied if the *RAM* space measure is based on *size<sub>s</sub>*, *iii*. This measure should therefore be replaced in the literature by the measure based on *size<sub>b</sub>*, *iv*.

## 2.3 Storage modification machines

The Storage modification machine (*SMM*) is a model introduced by Schönhage in 1970. The model is similar to the model proposed by the Soviet mathematicians Kolmogorov and Uspenskii (*KUM*) in 1958 [60], but for a precise comparison I refer to the discussion in Schönhage's paper. For the purpose of this chapter I will base myself on the 1980 version published in *SIAM. J. Comput.* [114]. The author advocates his model as a model of extreme flexibility and therefore it should serve as a basis for an adequate notion of time complexity.

The model resembles the *RAM* model as far as it has a stored program and a similar flow of control. Instead of operating on registers in memory it has a single storage structure, called a  $\Delta$ -structure. Here  $\Delta$  denotes a finite alphabet consisting of at least two symbols. A  $\Delta$ -structure  $S$  is a finite directed graph each node of which has  $k = \#\Delta$  outgoing edges which are labeled by the  $k$  elements of  $\Delta$ . The main difference between the *SMM* and the *KUM* can now be explained: the *KUM* operates on undirected instead of directed graphs.

There exists a designated node  $a$  in  $S$ , called the *center* of  $S$ . There exists a map  $p^*$  from  $\Delta^*$  to  $S$  defined as follows: For the empty string  $\epsilon$  one has  $p^*(\epsilon) = a$ , and otherwise  $p^*(wa) =$  the end-point of the edge labeled  $a$  starting in  $p^*(w)$ . The map  $p^*$  does not have to be surjective; however, nodes which can not be reached by tracing a word  $w$  in  $\Delta^*$  starting from the center  $a$  will turn out to play no subsequent role during the computations of the *SMM*, and therefore nodes may as well be assumed to have disappeared when they become unreachable.

The program of the storage modification machine consists, similar to the *RAM* program, on the one hand of flow of control instructions (*goto*, *accept*, *halt*, ...), and transput instructions (*read* and *print* - in this case a *read* will input a single bit and act like a conditional jump, depending on the value of the bit read), and on the other hand instructions which operate on memory - in this case a  $\Delta$ -structure  $S$ . There exist three types of

instructions of the latter type:

1. *new w*: creates a new node which will be located at the end of the path traced by  $w$ ; if  $w = \epsilon$  the new node will become the center; otherwise the last edge on the path labeled  $w$  will be directed towards the new node. All outgoing edges of the new node will be directed to the former node  $p^*(w)$
2. *set w to v*: redirects the last pointer on the path labeled by  $w$  to the former node  $p^*(v)$ ; if  $w = \epsilon$  this simply means that  $p^*(v)$  becomes the new center; otherwise the structure of the graph is modified.
3. *if v = w (if v ≠ w) then ...*: the conditional instruction (conditional jump suffices); here it is tested whether the nodes  $p^*(v)$  and  $p^*(w)$  coincide or not.

Simple it may look like the effect of a storage modification machine in action is extremely hard to trace. An alternative name for the machine, coined by Knuth, is *pointer machine*, and indeed: its semantics is as complicated as any pointer based algorithm. The model suffers from the *paradoxes of assignment*:

- After *new w* it is not necessarily true that  $p^*(w)$  denotes the new node
- the Hoare formula  $p^*(v) = x \{ \text{set } w \text{ to } v \} p^*(w) = x$  is invalid

These problems—related to the issue of proving correctness of *SMM* programs—are however not related to the complexity properties of the model which form the subject of this chapter.

### 2.3.1 Complexity measures for the *SMM*

In order to discuss complexity we need a time and a space measure. For the time measure there exists only one candidate: uniform time measure. One might consider to charge an instruction according to the length of the paths traced during the instruction, but since the arguments of the instructions are denoted explicitly in the program this length is a constant, independent of the  $\Delta$ -structure currently in memory. Therefore such a weighted measure will differ from the uniform measure by no more than a constant factor which is fully determined by the program.

Schönhage has not introduced a space measure. The reasonable candidate seems to be the number of nodes in the current  $\Delta$ -structure, and, in fact, this corresponds to the definition given in [41]. This definition has, however, some problematic aspects. There exist simply too many  $\Delta$ -structures of  $n$  nodes. Schönhage presents for the number  $\mathcal{X}(k, n)$  of  $\Delta$ -structures of  $n$  nodes over an alphabet  $\Delta$  of size  $k$  the following estimates:

$$n^{n \cdot k - n + 1} \leq \mathcal{X}(k, n) \leq \binom{k \cdot n}{n} \cdot \frac{n^{n \cdot k - n + 1}}{n \cdot k + 1}$$

From the above estimate it follows that in space  $n$  one can encode  $O(n \cdot k \cdot \log(n))$  bits of information rather than  $O(n \cdot \log(k))$  bits as on  $n$  squares of a Turing machine tape with alphabet  $\Delta$ . This leads to a similar situation as we have seen in chapter 2.2.4: consider the following on-line recognition problem:

$L_1 = \{I_1\#\dots\#I_k \ \& \ v = w \mid I_1\#\dots\#I_k \text{ encodes a straight-line } SMM \text{ program such that after performing it } p^*(v) = p^*(w)\}$

It is evident that an *SMM* can recognize this language in space  $O(n)$  where  $n$  denotes the number of *new* instructions in the straight-line program  $I_1\#\dots\#I_k$ . On the other hand a Turing machine on-line recognizer must have written a full description of the  $\Delta$ -structure generated by program  $I_1\#\dots\#I_k$  when it reads the  $\&$ -symbol, and therefore by an information theoretical argument it must have consumed by that time space  $\Omega(n \cdot \log(n))$ .

This example leaves open the situation for the case of off-line computations but it is not difficult to see that also there  $O(n)$  nodes on an *SMM* are capable of storing as much information as  $O(n \cdot \log(n))$  tape squares on a Turing machine. The basic technique consists of a  $\Delta$ -structure which represents a cycle of  $O(n)$  nodes with in each node a pointer to some other node in the cycle. The latter pointer stores in this way  $O(\log(n))$  bits of information, and it is a matter of simple programming to show that this representation can be used and updated, using no more than  $O(\log(n))$  additional nodes. A similar technique works also for the *KUM*. For more details see [133].

More recently it has been shown by Luginbuhl and Loui [71] that the above gain in space can actually be obtained with a constant factor overhead in time by using a complete binary tree in stead of a cycle. The paths toward the leaves encode the  $O(\log(n))$  bits stored in this node and it is easy to see that they can be read-out with constant factor time overhead, provided that there after the crossing of the border between two tape blocks of size  $O(\log(n))$  the simulated Turing machine can compute onwards for at least  $O(\log(n))$  steps. This additional property can be obtained with constant factor overhead in time as well; compare with the idea of making a Turing machine computation block respecting in the proof of the time vs. space theorem of Hopcroft, Paul and Valiant [50].

Based on the above observation one obtains the space measure  $n \cdot \log(n)$  for a  $\Delta$ -structure of  $n$  nodes. Under this measure the *SMM* can be simulated with constant factor space overhead on a Turing machine and vice-versa. An alternative is to make the dependence on  $k$  explicit by assigning an  $n$ -node  $\Delta$ -structure over an alphabet  $\Delta$  of size  $k$  the measure  $n \cdot k \cdot \log(n)$ . This measure has been proposed by Borodin e.a. [15], where they call this measure the *capacity*. Introducing the factor  $k$  in this space measure is similar to charging each tape cell on a Turing machine for  $\log(k)$ , where  $k$  is the size of the tape alphabet. Note, however, that this choice for a space measure would destroy the constant factor speed-up property for Turing machine space.

### 2.3.2 Simulations for the *SMM* model

Schönhage [114] has compared his model both with (multi-dimensional) Turing machines and with some *RAM* models in the uniform time measure. The results are:

**Theorem 7** *The SMM is related to other sequential models by*

1.  $SMM \approx SRAM - Utime(real - time)$
2.  $T \leq SMM(real - time)$

As a consequence, the *SMM* can be simulated on a standard *RAM* in the logarithmic time measure with  $n \cdot \log(n)$  time-overhead. The simulations for 1) are straightforward. In the proof Schönhage introduces another simplified *RAM* model, the main advantage of

which is that the instruction code is address free (instead the model uses a special purpose address register). This reduced *RAM* model again is real-time equivalent to the *SRAM* in the uniform time measure.

Recently Schönhage [116] has published a corresponding lower bound which shows that there exists a gap between the *SMM* model and the *RAM* in the logarithmic time measure. He proves that storing a bitstring in a *RAM* requires nonlinear time, even if reading is for free and if arbitrary arithmetic instructions are available.

Simulation 2) is easy for one-dimensional Turing machines. In the higher dimensional case a tape-segmentation trick is used.

A final advantage claimed for the *SMM* model is that it supports an implementation of integer multiplication which runs in linear time. So if everybody in the area of analysis of algebraic and numerical algorithms would embrace this model all occurrences of the multiplication complexity function  $M(n)$  could be replaced by simple  $O(n)$  expressions. Attractive it might look like, it seems that this proposal so far has few followers in the literature. For a discussion on the connection between theoretical fast multiplication algorithms and practical machines I refer to the recent note by Schönhage in the EATCS bulletin [115].

There are a few theoretical results known about the *SMM* model. In [41] Halpern e.a. prove a version of the Hopcroft, Paul, Valiant result that time  $T(n)$  can be simulated in space  $O(T(n)/\log(T(n)))$ . For this purpose they use the space measure determined by the number of nodes, which, as we indicated above, differs from Turing machine space by a logarithmic factor. Since the time measure seems to differ from Turing machine time by the same logarithm one still can maintain that it is a similar result.

More recent is the result presented by Schnitger [111] which seems to indicate that there exists a non-linear gap in time between the *SMM* and the *KUM*. Since the main distinction between these two models is that Schönhage uses a directed graph as storage structure, where Kolmogorov and Uspenskii use an undirected graph, the implicit bound of  $\#\Delta$  on the out-degree for the *SMM* becomes a bound on the indegree for the *KUM* as well. Therefore there can't exist too many short paths leading to some particular node. Schnitger introduces a well designed on-line real-time recognition problem which exploits this feature, and establishes a non-linear lower bound for this problem on the *KUM*, depending however on some unproven conjecture about communication complexity.

The *SMM* represents an interesting theoretical model, but, in the context of the above observations, its attractiveness as a fundamental model for complexity theory is questionable. Its time measure is based on uniform time in a context where this measure is known to underestimate the true time complexity. The same observation holds for the space measure for the machine; under the plausible definition of the space measure the *SMM* does not belong to the first machine class as defined in chapter 1.2.4. This problem can be solved by the introduction of a logarithmic space measure where an  $n$ -node  $\Delta$ -structure is charged for space  $n \cdot \log(n)$ , but such an unnatural definition for the space measure would make the model even less attractive.

## 2.4 Networks and non-uniform models

All machine models discussed so far have the property that a single machine, or a single program operates on inputs of arbitrary length. Occasionally the machine, in performing



some simulation, will first have to decide the length of the input, and then, after having set up the proper simulation parameters, the “real” simulation will start.

Networks are a typical representative of an alternative approach to the complexity of decision problems and/or function evaluation. In this approach the complexity investigated is of a more structural kind. Given some problem, say deciding membership in the language  $L$ , one first reduces this problem to a family of finite problems  $L_n = L \cap \Sigma^n$ . These finite problems are easy with respect to machine computations since their solution can be programmed using table-look-up. Instead one investigates the size  $s_n$  of the devices which solve these finite problems. The quantity  $s_n$  as a function of the input size  $n$  now becomes a measure of the complexity of the problem  $L$ .

From the perspective of computation on arbitrary length input  $x$  the process of computation now clearly is decomposed in two steps: from the input length  $n = |x|$  one determines the  $n$ -th device  $M_n$ ; next  $x$  is processed by  $M_n$  in order to produce the required answer.

Intuitively the time complexity of  $L$  is bounded by the complexity of these two stages together. But note that neither of these two stages has a time complexity which is measured by  $s_n$  as a function of  $n$ ; for the first stage one obtains at best a lower bound: retrieving an object of size  $s_n$  requires time  $\Omega(s_n)$ . For the second stage  $s_n$  seems not to be related at all to the time required to operate  $M_n$  on  $x$ .

At this point we are saved by some restrictions which are enforced in order to make the theory meaningful. The first restriction is that the mapping  $n \Rightarrow M_n$  is *uniform* (it should be computable in polynomial time and/or logspace). The second restriction is that the devices  $M_n$  are extremely easy to evaluate, as is for example the case with both formulas and networks.

Without the first condition it is possible to have undecidable problems with trivial network complexity: take any undecidable problem  $U \subseteq \omega$ , and consider the language  $L_U$  defined by  $x \in L_U \Leftrightarrow |x| \in U$ . For every length  $n$  there exist very simple devices which accept or reject all inputs of length  $n$ ; it is the problem of deciding which of the two to use which is undecidable.

There exists a meaningful alternative for making the mapping  $n \Rightarrow M_n$  uniform; one can also consider Turing machine computability relative to some suitable oracle. See Schnorr [112] or Karp and Lipton [57] for more details.

Without the second restriction that the devices are extremely simple to operate the theory cannot provide any complexity bounds above  $\log(n)$ : take a universal Turing machine with a program for  $L$ ; this is an object of constant size. If one adds a description of the length  $n$  in  $\log(n)$  bits one can obtain a finite device prepared for dealing with inputs of length  $n$ . Asymptotically, the  $\log(n)$  bits of  $n$  dominate in the description of the size of this object.

There exists an extensive literature on the subject of network-like models and non-uniform complexity. A detailed treatment falls outside the scope of this chapter. I will restrict myself therefore to a few results which connect this theory to the machine models introduced elsewhere in this chapter. For more details I refer to the source references. An example of a textbook on complexity theory which is based on network-like models is the 1976 text by Savage [102]. Information on the role of nonuniform complexity in structural complexity theory can be found in chapter 5 of [3].

### 2.4.1 The network model

A *logical network* or *Boolean circuit* is a finite labeled directed acyclic graph. Input nodes (output nodes) are nodes without ancestors (successors) in the graph. Input nodes are labeled with the names of input variables  $x_1, \dots, x_n$ . The internal nodes in the graph are labeled with functions from a finite collection, called the basis of the network. Here the condition is enforced that the number of ancestors of an internal node is equal to the number of arguments of its label; moreover a suitable edge labeling establishes a 1-1 correspondence between these ancestors and the arguments of this label.

By induction one defines for every node in the network a function represented by this node: for an input variable labeled by  $x_j$  this is the projection function  $(x_1, \dots, x_n) \rightarrow x_j$ ; an internal node with label  $g(y_1, \dots, y_n)$  for which the  $k$  ancestors represent the functions  $h_1(x_1, \dots, x_n), \dots, h_k(x_1, \dots, x_n)$  represents the function  $g(h_1(x_1, \dots, x_n), \dots, h_k(x_1, \dots, x_n))$ . A function  $f(x_1, \dots, x_n)$  is computed by the network iff it is represented by an output node.

The most important size functions for networks are:

- *depth*: the length of a longest path in the graph between an input node and an output node
- *size* : the total number of nodes in the graph

It used to be common to restrict oneself to networks over a basis of logical functions consisting of monadic and binary functions only: There exists four monadic functions, among which only the negation function is important. There exist 16 binary logical functions, among which the *and*, the *or*, and the *xor* are the best-known ones in the network theory. Implication is another function but its role in network theory is incomparably small compared to its role in logic and mathematics. It is a known result that there exist binary functions like the *nand* function, which form a complete basis by themselves. Still the most common network basis consists of three functions: *and*, *or*, and *not*. If the negation is omitted the network only represents monotone functions and will be called a monotone network. Networks with unbounded fan-in *and* and *or* nowadays are also being considered.

Every node in the network can be ancestor of an arbitrary number of successors. If the number of successors is restricted to 1 the network becomes a tree, and one talks about a formula rather than a circuit.

The complexity measures introduced below depend on the logical basis, and on the fact whether one deals with networks or formulas. In the latter situation one speaks about formula complexity rather than network complexity. The logical basis is implied by the context; default is the standard three element basis, and it is moreover not hard to see that for network complexity the effect of the choice of a basis is at most a constant factor in depth and size; for formulas the story is different, see [63] [93]. If unbounded fan-in is disallowed in the basis the constant factor in size remains but a logarithmic penalty in depth may occur.

For a given finite collection of logical functions  $f_1(x_1, \dots, x_n), \dots, f_k(x_1, \dots, x_n)$  the network complexities  $size(f_1, \dots, f_k)$  and  $depth(f_1, \dots, f_k)$  are defined to be the minimal size or depth of a network such that the entire set is computed by this network. If  $\mathcal{F}_n$  for  $n \in \omega$  is a family of  $n$ -argument functions it is possible to investigate  $size(\mathcal{F}_n)$  and

$depth(\mathcal{F}_n)$  as a function of  $n$ . If  $L$  is a language over  $\{0, 1\}^*$  then the network complexities of  $L$  are defined as the network complexities of the family  $\mathcal{F}_n$  where  $\mathcal{F}_n$  for every  $n$  consists of the characteristic function  $L_n$  of the set  $L \cap \{0, 1\}^n$ .

The following elementary results are known:

**Theorem 8** *For arbitrary languages  $L$  the following estimates hold*

1.  $size(\mathcal{F}_n) \leq \#\mathcal{F}_n \cdot 2^n$
2.  $size(\mathcal{F}_n) \leq \frac{2^n}{n \cdot (1+\epsilon)}$  for  $\#\mathcal{F}_n = 1, n \geq N(\epsilon)$
3.  $depth(\mathcal{F}_n) \leq n + 1$
4.  $depth(\mathcal{F}_n) \leq size(\mathcal{F}_n)$
5.  $size(\mathcal{F}_n) \leq \#\mathcal{F}_n \cdot 2^{depth(\mathcal{F}_n)}$

Result 1) is easy and results 4) and 5) are trivial. Result 2) is the upper bound established by Lupanov [72], and result 3) has been obtained by McColl and Paterson [76]. It is known that the bound in 2) is asymptotically tight; by enumeration of all possible networks one shows that there for every  $\delta > 0$  only a fraction of all functions can be computed by a circuit of size  $\leq 2^n(1 - \delta)/n$  since there are not enough of those small circuits. It is a much harder problem to prove lower bounds in network complexity for explicit functions. There the largest lower bound in the literature for network complexity shown for an individual explicit function equals  $3 \cdot n + o(n)$  [11]. For formula complexity one has an  $\Omega(n^2/\log(n))$  lower bound [82]. For monotone networks exponential lower bounds for specific functions have been established in 1985 [2] [95].

It is generally assumed that establishing lower bounds for network complexities is about the strongest statement one can make about the intractability of some problem. It is also the most difficult task faced by the complexity theorist, and so far progress has been much smaller than one might hope for.

### 2.4.2 Relation with Turing machine complexity

The finite control of a Turing machine can be considered to be some finite automaton which can be implemented using a Boolean network. The size of this network tells something of the complexity of the program of the machine, whereas the depth of this network indicates how fast the machine can perform a single step of its computation, assuming that the machine has to be implemented using logical gates. Still these connections between networks and Turing machines are not the ones investigated in the literature; there one considers the relation between time and space measures for languages recognized by Turing machines and the network complexity for these languages as defined in the previous section.

Consider the time-space diagram which encodes the complete computation of some Turing machine  $M_i$  on some input of length  $n$ . Without loss of generality we can assume that the  $M_i$  consumes both maximal time and space on every input of length  $n$ , so the size of this diagram does not depend on the precise input. The structure of the diagram is such that it is filled with characters, and that each character is determined completely by the characters located at a small local set of neighbors. These dependencies are directed forward in time. It is not hard to design a Boolean network of size proportional to the size of the time-space diagram which evaluates (an encoding of) the character at position

$(s, t)$  in the space-time diagram by an amount of circuitry for this location which depends of the program only. In this way one obtains the following result from [102]

**Theorem 9**  $size(L_n) \leq P \cdot T(n) \cdot S(n)$

where  $T(n)$  and  $S(n)$  denote time- and space complexity of the Turing machine and  $P$  denotes its size.

In the above simulation a large amount of circuitry is wasted for calculations the purpose of which is to copy the contents of an unvisited tape cell from one configuration to the next. These local circuits could be eliminated, provided one would know in advance where the heads are residing on the tapes. But that is the problem which has been solved by the Hennie-Stearns simulation where a  $k$ -tape Turing machine was obviously simulated on a two-tape machine. Exploiting that idea one obtains the following improvement from Pippenger and Fischer [92], or Schnorr [112]

**Theorem 10**  $size(L_n) \leq P \cdot T(n) \cdot \log(T(n))$

It follows from the construction that both results are proved using a uniform family of circuits. Conversely, it is not hard to see that for logspace-uniform networks, a Turing machine can construct the network and subsequently evaluate in time polynomial in the size. A tighter bound has been established by Pippenger [91]. There exists also a tight relation between the depth measure and the space complexity of Turing machines. For spacebounds  $s(n) \geq \log(n)$  which are constructible in space  $O(\max(\log(n), \log(s(n))))$  Borodin [14] has shown

**Theorem 11**  $depth(L_n) \leq O(s(n)^2)$ , when  $L \in NSPACE(s(n))$

Moreover for these functions  $s(n)$  circuits of depth  $s(n)$  can be evaluated in space  $O(s(n))$ .

More recently networks have become used as a tool for characterizing complexity of parallel computation as well. This connection has lead to the introduction of three hierarchies of classes

- $NC_k$  = the class of languages  $L$  recognized by logspace uniform networks of polynomial size and depth  $O(\log(n)^k)$
- $AC_k$  = the class of languages  $L$  recognized by logspace uniform unbounded fan-in networks of polynomial size and depth  $O(\log(n)^k)$
- $SC_k = T-TIME \ \& \ SPACE(Pol, O(\log(n)^k))$

The families  $NC_k$ ,  $AC_k$  and  $SC_k$  form a hierarchy within  $P$ . The first two hierarchies, moreover can be combined to a single hierarchy:

$$NC_0 \subseteq AC_0 \subseteq \dots \subseteq NC_k \subseteq AC_k \subseteq NC_{k+1} \subseteq \dots$$

The  $NC$  and  $SC$  hierarchies have been named after Nick Pippenger and Steve Cook respectively. See Cook [19] [20] for further details.

## Chapter 3

# The second machine class

This chapter is dedicated to a survey of the second machine class which consists of those models which satisfy the parallel computation thesis. This thesis states that parallel time and sequential space are polynomially related. As will be seen in the sequel models don't have to use parallelism in order to become a second machine class member. The alternating models in chapter 3.2 are sequential devices with a modified mode of acceptance, whereas the models in chapter 3.3 are sequential devices which operate on huge objects in unit time. The really parallel models first appear in chapter 3.4 . As will be seen in chapter 4 the boundaries of the second machine class in the realm of parallel machines become rather vague. Slight modification of the rules of the game leads to even more powerful models. There are also weaker models of parallel devices, one of which is introduced in chapter 4.1 .

The current literature on parallel models is growing explosively and consequently in a chapter like the present one only a fraction can be covered. I am quite aware of the fact that even some main stream developments in the theory of parallel computation are not covered at all. I have not dealt with parallel models consisting of *RAM*'s or finite state devices in an arbitrary interconnection pattern [39], models for systolic computation [40] [77], cellular automata [128], or the subject of relative quality of one interconnection pattern with respect to another [36] [62] [80]. I have not looked into the relation between parallel models and network resource bounds [19] [20]. The subject of write-conflict resolution methods is only mentioned [25] [28] [29] [68]. I have also paid no attention to what happens if additional resource bounds are enforced on parallel devices, like polynomial bounds on the number of processors used or other hardware bounds [24]. The reader interested in such subjects is referred to the indicated literature.

For the class of models covered in the survey I establish a tight connection between simulations of sequential devices and transitive closure algorithms which perform in essence a brute force try-out everything simulation. This connection will be explained in chapter 3.1. Since all efficient simulations presented in chapters 3 and 4 of this chapter are based on this connection the "unreasonableness" of the parallel models presented is clearly exposed. If parallel models are efficient by simulating sequential devices in such a crude manner something must be wrong. Still the very same connection enables me to place these models in a common perspective.

### 3.1 *PSPACE* and transitive closure

*PSPACE* is the class of those languages which can be recognized in polynomial space on some sequential device, in particular on a standard single-tape Turing machine. There exists a generic translation between space bounded computations and paths in a suitably constructed configuration graph. For reasons of simplification we will assume from now onwards that all space bounds are  $\Omega(\log(n))$ . For a given input  $x$ , a given Turing machine  $M$  and a given space bound  $S$  we can form the graph  $G(x, M, S)$  of all configurations  $c$  of machine  $M$  which use space  $\leq S$ ; an edge connects two configurations  $c_1$  and  $c_2$  if there is a one step transition from  $c_1$  to  $c_2$ . This graph has the following properties:

1. For every input  $x$  and space bound  $S$  there exist a unique node corresponding to the initial configuration on input  $x$
2. Assuming that a suitable notion of acceptance has been chosen the accepting configuration is unique
3. The number of nodes in the graph  $G(x, M, S)$  is bounded by some exponential function  $2^{c \cdot S}$ , where the constant  $c$  depends on  $M$  but not on  $x$
4. The graph  $G(x, M, S)$  can be encoded in such a way that an  $S$ -space bounded Turing machine on input  $x$  and a description of  $M$  can write the encoding of  $G(x, M, S)$  on some write-only output tape
5. If  $M$  is deterministic, then every node in  $G(x, M, S)$  has outdegree  $\leq 1$ ; if  $M$  is nondeterministic then the outdegree of some nodes can be  $\geq 2$ , but for a suitable restriction of the Turing machine model the outdegree can be assumed to be  $\leq 2$  as well
6. The input  $x$  is accepted in space  $S$  by  $M$  iff there exists a path from the unique initial configuration on input  $x$  in  $G(x, M, S)$  to an (or with a suitable restriction on the model, the unique) accepting configuration. This path can be assumed to be loop-free, and therefore its length can be assumed to be  $2^{c \cdot S}$

The above properties suggest the following universal algorithm for testing membership in languages in *PSPACE*: assume that  $L$  is recognized by  $M$  in space  $n^k$ , then in order to test whether  $x \in L$  we first construct  $G(x, M, |x|^k)$ , next we compute the reflexive transitive closure of this graph, and finally investigate whether in this transitive closure there exists an edge between the unique initial configuration on  $x$  and the unique accepting configuration. Although it seems that this is a rather expensive method for simulating a single computation, it is exactly this transitive closure algorithm which is hidden in almost all proofs that some particular parallel machine obeys the parallel computation thesis.

#### 3.1.1 Transitive closure algorithms and *PSPACE*-complete problems

Various designs of a transitive closure algorithm lead to various insights. In the first place we can represent the graph  $G(x, M, S)$  by a Boolean matrix  $A(x, M, S)$  where 1 denotes the presence of an edge and 0 denotes the absence of an edge. The row and column indices denote configurations. Clearly these configurations can be written down in space  $S$ ; the

total number (and therefore also the size  $N$  of the matrix  $A(x, M, S)$ ) is exponential in  $S$ ;  $N = 2^{c \cdot S}$ . If we moreover let  $A(x, M, S)[i, i] = 1$  for all  $i \leq N$ , then the transitive closure of  $M$  can be computed by  $c \cdot S$  squarings of  $A(x, M, S)$  using the Boolean matrix multiplication. If  $N^3$  processors are available each squaring can be performed in time  $O(S)$  (needed for adding  $N$  Boolean values), so the entire transitive closure algorithm takes time  $O(S^2)$ . If moreover a time bound  $T$  on the computation is given this time is reduced to  $O(S \cdot \log(T))$ , due to the fact that no paths longer than  $T$  edges have to be investigated. Finally, if the parallel model has a concurrent-write feature the time needed for adding the  $N$  Boolean values can be reduced to  $O(1)$ , and in this case the time for the transitive closure algorithm becomes  $O(\log(T))$ . Note however that, in order to perform this algorithm, the matrix  $A(x, M, S)$  must be constructed first.

Next we investigate the following recursive function  $path(w, i, j)$  which evaluates to *true* in case there exists a path from node  $i$  to node  $j$  of length  $2^w$  :

```

proc path = (int w, node i, j) bool:
  if w = 0 then i = j or edge(i, j)
  else
    bool found := false ;
    forall node n while not found
      found := found or ( path(order-1, i, n) and path(order-1, n, j) )
    od;
  found
fi;

```

Existence of an accepting computation can be evaluated by the call  $path(c \cdot S, \text{init}, \text{final})$ , where *init* and *final* are the unique initial and final configurations in  $G(x, M, S)$ ; given a time bound  $T$  the initial parameter  $w$  can also be chosen to be  $\log(T)$ . With recursion depth  $c \cdot S$  (respectively  $\log(T)$ ) and parameter size  $O(S)$  this recursive procedure can be evaluated in space  $O(S^2)$  (respectively  $O(S \cdot \log(T))$ ); this is the main ingredient of the proof of Savitch's theorem [104] which implies that  $PSPACE = NPSPACE$ :

**Theorem 12** *If the language  $L$  is recognized by some nondeterministic Turing machine in space  $S(n) \geq \log(n)$ , then it can be recognized by some deterministic Turing machine in space  $S(n)^2$*

A third formulation of the transitive closure algorithm yields the *PSPACE-completeness* of the problem *QUANTIFIED BOOLEAN FORMULAS (QBF)* [125] :

**QUANTIFIED BOOLEAN FORMULAS:**

**INSTANCE:** A formula of the form  $Q_1 x_1 \dots Q_n x_n [P(x_1, \dots, x_n)]$ , where each  $Q_i$  equals  $\forall$  or  $\exists$ , and where  $P(x_1, \dots, x_n)$  is a propositional formula in the boolean variables  $x_1, \dots, x_n$ .

**QUESTION:** does this formula evaluate to *true*?

**Theorem 13** *The problem QUANTIFIED BOOLEAN FORMULAS is PSPACE-complete*

**Proof:** The idea is to encode nodes in  $G(x, M, S)$  by a boolean valuation to a sequence of  $k := c' \cdot S$  boolean variables. From the proof of Cook's theorem which establishes the *NP*-completeness of *SATISFIABILITY* ( see [18] or [37]) one obtains the existence of a propositional formula  $P_0$  in  $2 \cdot k$  variables  $P_0(x_1, \dots, x_k, y_1, \dots, y_k)$ , where the variables  $x_1, \dots, x_k$  encode node  $i$ , the variables  $y_1, \dots, y_k$  encode node  $j$  and  $P_0$  expresses that  $i = j$  or there exists an edge from  $i$  to  $j$ .

One now can define by induction a sequence of quantified propositional formulas  $P_d$  such that  $P_d(x_1, \dots, x_k, y_1, \dots, y_k)$  expresses the presence of a path of length  $\leq 2^d$  between node  $i$  and node  $j$ . In a naive approach the formula  $P_d$  would include only existential quantifiers, and  $P_d$  would include two copies of  $P_{d-1}$ ; by a standard trick from complexity theory we reduce the number of occurrences of  $P_{d-1}$  in  $P_d$  to one ; this trick, however, introduces universal quantifiers:

$$\begin{aligned} P_d(x_1, \dots, x_k, y_1, \dots, y_k) &= \exists z_1, \dots, z_k [\forall u_1, \dots, u_k [\forall v_1, \dots, v_k [ \\ &((u_1, \dots, u_k = x_1, \dots, x_k \wedge v_1, \dots, v_k = z_1, \dots, z_k) \vee \\ &(u_1, \dots, u_k = z_1, \dots, z_k \wedge v_1, \dots, v_k = y_1, \dots, y_k)) \\ &\Rightarrow P_{d-1}(u_1, \dots, u_k, v_1, \dots, v_k)]]] \end{aligned}$$

Substituting for  $x_1, \dots, x_k$  and  $y_1, \dots, y_k$  the codes of the initial and final node in  $G(x, M, S)$  in  $P_K(x_1, \dots, x_k, y_1, \dots, y_k)$  where  $K = c \cdot S$  denotes the logarithm of the size of  $G(x, M, S)$  we obtain a closed quantified Boolean formula, the truth of which expresses the existence of an accepting computation. It is not difficult to see that  $P_K$  is a formula of length  $O(k^2)$  in  $O(k^2)$  variables, since  $K = \Theta(k)$ . Here we have counted each variable as a single symbol ; clearly a representation in a finite alphabet will introduce another factor  $\log(k)$  in the length of the formula. From this one concludes that *QBF* is *PSPACE*-complete.  $\square$

The *PSPACE*-completeness of *QBF* explains the alternating nature of a number of known *PSPACE*-complete problems. *SATISFIABILITY* is the prototype of a solitaire game, where the player has to look for some configuration with a particular property or for a sequence of simple moves leading to some particular goal state, but the alternating quantifiers turn *QBF* into a two person game. Two players Elias and Alice in turn choose the truth values to be assigned to existentially or universally quantified variables in the order of their nesting inside the formula. Elias tries to establish the truth of the formula whereas Alice tries to show that the given formula is false. The truth of the entire formula is equivalent with the existence of a winning strategy for Elias in this game.

Starting with this game-theoretical interpretation of *QBF* several authors have investigated the end-game analysis of games inspired by real life games. A useful intermediate game is *GENERALIZED GEOGRAPHY* [109]; from there one can reach *HEX* on arbitrary graphs and even *HEX* on the traditional hexagonal board ; see [27] [97]. For people interested in the position of *CHECKERS*, *CHESS* and *GO*: these games are not on the list because, after earlier *PSPACE*-hardness results [32] [69], they turn out to be even more difficult than *PSPACE* [33] [98] [99].

From the above one should not conclude that in general solitaire games are at most *NP*-hard; beside the earlier *PSPACE*-completeness of the *BLACK PEBBLE GAME* [38], one has nowadays examples of group theoretical problems which have been shown to be



*PSPACE*-hard [54] [55]. Together with problems which encode *PSPACE*-bounded computations in a more direct way (like Reif's *GENERALIZED MOVER'S PROBLEM* [96]) this has led to an interesting zoo of *PSPACE*-complete problems. It should, however, be no longer a surprise that alternation forms a fundamental concept in one of the machine models in the second machine class.

### 3.1.2 Establishing membership in the second machine class

Suppose that we are faced with some parallel or otherwise powerful machine model  $\mathcal{X}$  and that we want to prove that  $\mathcal{X}$  is a member of the second machine class:  $\mathcal{X} - PTIME = \mathcal{X} - NPTIME = PSPACE$ . The above observations provide us with some tools for proving this result.

The inclusion  $PSPACE \subseteq \mathcal{X} - PTIME$  can be shown by inventing a polynomial time algorithm on  $\mathcal{X}$  which solves a *PSPACE*-complete problem like *QBF*. An alternative is to show that on  $\mathcal{X}$  one can implement one of the transitive closure algorithms from the previous chapter in polynomial time.

For the inclusion  $\mathcal{X} - NPTIME \subseteq PSPACE$  one uses in most cases a guess and verify method. A nondeterministic machine is used to guess a trace of an accepting  $\mathcal{X}$ -computation on the given input. Such a trace consists of instructions executed and memory values stored during this computation but in general not all information stored in memory during the  $\mathcal{X}$ -computation can be written down in polynomial space. Therefore one writes down enough information from which the remaining memory contents can be reconstructed. The reconstruction is done by a recursive procedure which reflects the machine architecture of model  $\mathcal{X}$ ; the recursive procedure tells how the memory contents at time  $t$  depend on those at time  $t - 1$ . Care has to be taken that all arguments for this procedure and also its values can be written down in polynomial space.

Using the trace and the recursive procedure it becomes possible to certify the trace as being consistent: for every conditional instruction the correct branch has been selected.

If one has established for two models  $\mathcal{X}$  and  $\mathcal{X}'$  that they are both members of the second machine class, one has indirectly shown that  $\mathcal{X}$  and  $\mathcal{X}'$  simulate each other with polynomial overhead in time. In the literature only a few instances of direct simulations establishing such polynomial time overheads are given. One example is the refined analysis of the power of various models of vector machines by Ruzzo in the unpublished report [100], where explicit overheads for simulation of alternating Turing machines and vector machines are given. Another example can be found in the paper by van Leeuwen and Wiedermann on array processing machines [134] where time overheads for simulation on the *SIMDAG* and a reverse simulation are determined.

## 3.2 The alternation model

### 3.2.1 The concept of alternation

The concept of alternation [16] leads to machine models which obey the parallel computation thesis without providing any intrinsic parallelism at all. As a computational device an alternating Turing machine is very similar to a standard nondeterministic sequential Turing machine; only the definition of accepting the input has been modified.

Since the machine is nondeterministic the computation can be represented as a computation tree the branches of which represent all possible computations. The leaves—the terminal configurations where the machine halts—are designated to accept or to reject as usual on basis of the designation of the included state as being accepting and rejecting. For a standard nondeterministic machine such a computation tree is considered to represent an accepting computation as soon as a single accepting leaf can be found. But for the alternating machine the notion of accepting is slightly more complicated .

The main idea is to equip states in the Turing machine program with labels *existential* and *universal*. Configurations inherit the label of the state included in this configuration. Next one assigns a quality *accept*, *reject* or *undef* to every node in the computation tree according to the following rules:

1. The quality of an accepting (rejecting) leaf equals *accept* (*reject*)
2. The quality of an internal node representing an *existential* configuration is *accept* if one of its successor configurations has quality *accept*
3. The quality of an internal node representing an *existential* configuration is *reject* if all of its successor configurations have quality *reject*
4. The quality of an internal node representing a *universal* configuration is *reject* if one of its successor configurations has quality *reject*.
5. The quality of an internal node representing a *universal* configuration is *accept* if all of its successor configurations have quality *accept*
6. The quality of a node with one successor equals the quality of its successor
7. The quality of any node the quality of which is not determined by application of the above rules is *undef*

Clearly the quality *undef* arises only if the computation tree contains infinite branches, but even nodes which have infinite offspring can obtain a definite quality since, for example, an accepting son of an *existential* node overrides the *undef* label of another son.

By definition an alternating device accepts its input in case the root node of the computation tree, representing the initial configuration on that input, obtains the quality *accept*.

The above treatment is a minor simplification of the presentation in [16] in so far as that the feature of negating states is not included. It should be clear that the notion of an alternating mode of computation makes sense for virtually every machine model and is not restricted to Turing machines. For example, it makes sense to consider alternating finite automata, alternating PDA's etc. For every device some power may be gained by proceeding to the level of alternation but whether the gain is substantial depends on the specific model considered. For example, in the case of finite automata, the languages recognized by alternating finite automata are still regular languages; the gain is a potentially doubly exponential reduction of the size of the automaton.

Time ( space ) consumed by an alternating computation is measured to be the maximal time (space) consumption along any branch in the computation tree. Another complexity measure which was considered by Ruzzo [101] for the alternating Turing machine is the

minimal size of an accepting computation tree. An accepting tree is a subtree of the full computation tree which yields a certificate that the input is accepted; it includes all successors of a *universal* node but for an *existential* node only one successor with quality *accept* must be included in the accepting computation tree. This measure is related both to nondeterministic time in the sequential model and parallel time and leads to various interesting classes with simultaneous resource bounds. I refer to the paper for more details.

### 3.2.2 Relation with sequential models

The alternating device, being an incarnation of a standard device in disguise, clearly inherits the simulation results for the first machine class devices. As a consequence there exists a device independent hierarchy for alternating classes:

$$ALOGSPACE \subseteq APTIME \subseteq APSPACE \subseteq AEXPTIME$$

Its behavior like a parallel machine model now is expressed by the following results which provide us with far more than the parallel computation thesis requires for becoming a second machine class member:

**Theorem 14** *The alternating Turing machine model is related to the sequential hierarchy by the equality*

$$APTIME = PSPACE$$

*but the other classes are shifted versions in the sequential hierarchy as well:*

$$ALOGSPACE = P, APSPACE = EXPTIME, AEXPTIME = EXPSPACE.$$

Note that the alternating devices have no nondeterministic mode of computation.

**Proof:** I will give a short indication why the above equalities are true. First consider the inclusion  $APTIME \subseteq PSPACE$ . It suffices to show that the quality of the initial configuration in a polynomial time bounded computation tree can be evaluated in polynomial space. Clearly this quality can be evaluated by a recursive procedure which traverses the nodes of the computation tree. This procedure has a recursion depth proportional to the running time of the alternating device, whereas each recursive call requires an amount of space proportional to the space consumed by the alternating device. Hence the space needed by the deterministic simulator is proportional to the space-time product of the alternating device, which in turn is bounded by the square of the running time.

The reverse inclusion  $PSPACE \subseteq APTIME$  follows as soon as we show how to *QBF* in polynomial time on an alternating machine. This is almost trivial: let a machine guess the valuation for the quantified variables where the universally (existentially) quantified variables are guessed in a universal (existential) state; these values are guessed in the order of the nesting in the formula. Next the formula is evaluated in a deterministic mode and the machine accepts (rejects) if the result becomes true (false) .

The equality  $AEXPTIME = EXPSPACE$  is obtained by a standard padding argument from the equality shown above.

The inclusion  $ALOGSPACE \subseteq P$  is shown as follows: note that a logspace-bounded alternating device for a given input has only a polynomial number of configurations. These configurations can be written on a work-tape. The terminal configurations obtain a quality

based on the included state. Next by repeatedly scanning the list of configurations the quality of intermediate configurations can be determined by application of the rules 2)–6). This scanning process terminates if during a scan no new quality can be determined. Since during each sweep either at least one quality is determined or the process terminates, and since the time needed for a single sweep is bounded by the square of the size of the list of configurations the running time for this procedure is polynomial.

The reverse inclusion  $P \subseteq ALOGSPACE$  is shown as follows. Assume that the language  $L$  is recognized in time  $T(n)$  by a standard single tape Turing machine  $M$ . It suffices to show how an alternating device can recognize  $L$  in space  $\log(T(n))$ . Consider therefore the standard computation diagram of the computation of  $M$  on input  $x$ . This diagram can be represented in the form of a  $K \times K$  table of symbols, where  $K = T(|x|)$ . The top row of this table describes the initial configuration on input  $x$ , and the bottom row describes the final configuration which should be an accepting one. Each intermediate symbol is completely determined by the three symbols in the row directly above it since the machine  $M$  is deterministic.

The alternating device now guesses the position in the bottom row of the occurrence of an accepting state, and certifies this symbol by generating in a universal state three offspring machines which guess in an existential state the symbols in the three squares above it. These guesses are certified in the same way, all the way up to the top row, where guesses are certified by comparison with the input  $x$ . The amount of space required by this procedure is proportional to the space required for writing down the position of the square considered in the diagram, which is  $O(\log(T(|x|)))$ . Since the machine  $M$  is deterministic, only those guesses which are correct can be certified (to be shown by induction on the row number) and therefore the guesses are globally consistent. This last observation breaks down for nondeterministic devices  $M$  and therefore we cannot obtain the inclusion  $NP \subseteq ALOGSPACE$  in this way.

Again by a simple padding argument one obtains the equality  $APSPACE = EXPTIME$ , as an easy consequence of the equality  $ALOGSPACE = P$ .  $\square$

### 3.3 Sequential machines operating on huge objects in unit time

The first machine model for which the validity of what later was to become known as the parallel computation thesis has been established is the vector machine model of Pratt and Stockmeyer [94], shortly later succeeded by the *MRAM* of Hartmanis and Simon [44] [45]. These models have in common that their power originates from the possibility to operate on objects of exponential size in unit time.

All these models are derived from the *RAM* model with uniform time measure by extending the arithmetic with new powerful instructions. In the vector machine this extension consists of the introduction of a new type of registers, called vectors, which can be shifted by amounts stored in the arithmetical registers of the *RAM*. The contents of the vector registers can also be subjected to parallel bitwise Boolean operations like *and*, *or*, or *xor*. This makes it possible to program the concatenation of the contents of two vectors and to perform various masking operations. The *MRAM* model was obtained by realizing that shifting a vector amounts to multiplication or division by a suitable power

of two. Hence the separation between vectors and arithmetic registers is an inessential feature in the model; the same power can be achieved by introducing multiplication and division in unit time, preserving the bitwise Boolean operations.

Restrictions of the model have been investigated. For example one can forsake one of the two shift directions (right shift of one register is simulated by left shifting all the others); as a consequence one can drop the division instruction, which yields the *MRAM* model as proposed in [44]. More recently it has been established that the combination of multiplication and division, in absence of the bitwise Boolean instructions suffices as well [9] [113]. This result shows the power of a purely arithmetical model.

In this section I will illustrate the power of these models by a model which has been obtained by moving in the other direction: stressing the pure symbol manipulation instructions and dropping the powerful arithmetic. This is the *EDITRAM* model proposed in [123] as a model of the text editor you may have in mind while editing texts behind your terminal. I will present an outline of the proof of the validity of the parallel computation thesis for the *EDITRAM* and indicate the connection with the proofs for the earlier models.

### 3.3.1 The *EDITRAM* model

In the *EDITRAM* we extend the standard *RAM* with a fixed finite set of text-files. Standard arithmetic registers can be used as cursors in a text-file. Beside the standard instructions on the arithmetic registers the *EDITRAM* has instructions for:

1. reading a symbol from a file via a cursor
2. writing a symbol into a file via a cursor
3. positioning a cursor at the end of a file (thus computing its length)
4. positioning a cursor into a file by loading an arithmetic value into the cursor
5. systematic replacement of *string1* by *string2* in a text-file
6. concatenation of text-files
7. copying of segments of text-files as indicated by cursor positions
8. deletion of segments of text-files as indicated by cursor positions.

In the systematic string replacement instruction 5) the arguments *string1* and *string2* are to be presented by literals in the program; substitution of the contents of an entire text-file for a single character would allow a doubly exponential growth of the size of text-files, which is more than we are aiming for.

The time complexity of the model is defined by using the uniform time measure for the arithmetic registers. So an edit instruction is charged one unit of time. A reasonable alternative would be to use the logarithmic time measure with respect to the arithmetic registers. Then an edit instruction would be charged according to the logarithm of the values of the involved cursors, and its cost therefore is proportional to the logarithm of the length of the (affected portion of) the text-file. Since in our model the growth

of a text-file is at most exponential the two measures will be polynomially related. This observation no longer would be valid if we would allow that entire text-files are substituted for strings by the replace operation; the doubly exponential growth of the text-files will require exponential time in the logarithmic measure so the uniform and the logarithmic measure no longer are polynomially related.

### 3.3.2 The *EDITRAM* is a second machine class device

**Theorem 15** *The EDITRAM is a second machine class device*

**Proof:** In order to verify that the *EDITRAM* obeys the parallel computation thesis we must prove the two inclusions  $EDITRAM-NPTIME \subseteq PSPACE$  and  $PSPACE \subseteq EDITRAM-PTIME$ .

The proof of the first inclusion is characteristic for the proof of this inclusion for similar models. Given an input we must test in polynomial space whether the given *EDITRAM* machine will accept this input or not. But by Savitch's theorem 12 our simulation may be nondeterministic. Therefore we first guess the trace of some accepting computation and write it down on some work-tape. The accepting computation being polynomially time bounded we can write down the sequence of instructions in the program of the *EDITRAM* which are executed. Moreover, since the length of the values of the arithmetic registers is linearly bounded by the time, we can also maintain a log on the register values in polynomial space.

We cannot maintain a log on the contents of the textfiles, since their length may grow exponentially. Instead we introduce a recursive procedure  $char(time, position, textfile)$  which evaluates to the character located at the given position in the given textfile after performing the instruction at the given time. The arguments of this procedure can be written down in polynomial space, due to the fact that the growth of the length of a textfile is bounded by a simple exponential function in the time (both systematic string replacement and concatenation will at most multiply the length of a textfile by a constant). Given this procedure it is possible to certify that the trace written on the work-tape indeed represents an accepting computation.

From the meaning of the individual instructions one can obtain a recursive procedure which expresses the value of  $char(time, position, textfile)$  in terms of similar values after the previous instruction at time  $time-1$ . In the case where the present instruction is a systematic string replacement we face the problem to figure out where the character at the given position was located before the replacement. Since this requires information on the number of occurrences of the replaced pattern preceding this position in the given textfile, the entire textfile, up to the given position must be re-computed by recursive calls. This is the most complicated case in the description of this recursive procedure. For details I refer to [123].

The total space required by the evaluation of this procedure is bounded by the product of the size of an individual call (which we indicated to be polynomial) and the recursion depth (which is bounded by the running time of the *EDITRAM* computation being simulated, which was also assumed to be polynomial). This completes the proof of the first inclusion.

Next we consider the inclusion  $PSPACE \subseteq EDITRAM-PTIME$ . It suffices to show how to solve the *PSPACE*-complete problem *QBF* in polynomial time on a determin-

istic *EDITRAM*. Consider a given instance  $Q_1x_1 \dots Q_nx_n[P(x_1, \dots, x_n)]$  of *QBF*. Our algorithm is performed in three stages:

Step 1: Remove the quantifiers in the order of their nesting from inside to outside by programming the transformations:

$$\begin{aligned}\forall x_i[P(\dots, x_i, \dots)] &\Rightarrow (P(\dots, 0, \dots) \wedge P(\dots, 1, \dots)) \\ \exists x_i[P(\dots, x_i, \dots)] &\Rightarrow (P(\dots, 0, \dots) \vee P(\dots, 1, \dots))\end{aligned}$$

Clearly each transformation preserves the truth of the involved formula; the involved formula  $P$  is a quantifier free formula, due to the order of the quantifier eliminations. After elimination of all quantifiers a formula of exponential size is obtained which still is equivalent to the given instance of *QBF*.

Step 2: Evaluate the resulting formula by systematic string replacements of the type:

$$\begin{array}{llll} (0 \vee 0) \Rightarrow 0 & (0 \vee 1) \Rightarrow 1 & (1 \vee 0) \Rightarrow 1 & (1 \vee 1) \Rightarrow 1 \\ (0 \wedge 0) \Rightarrow 0 & (0 \wedge 1) \Rightarrow 0 & (1 \wedge 0) \Rightarrow 0 & (1 \wedge 1) \Rightarrow 1 \\ (\neg 0) \Rightarrow 1 & (\neg 1) \Rightarrow 0 & (0) \Rightarrow 0 & (1) \Rightarrow 1 \end{array}$$

These transformations can be produced by local systematic string replacements.

Step 3: Check whether the resulting literal equals 0 or 1 .

Note that after a single cycle through the replacements in step 2 the depth of the involved propositional expression has been decreased by at least 1. If the depth of the propositional kernel of the given instance was  $k$ , then after the transformations of step 1 the depth of the intermediate formula is  $k + n$  ( $n$  being the number of quantifiers eliminated) . Therefore the number of iterations in step 2 is polynomial.

It remains to show how to perform the transformations in step 1 . Clearly it suffices to locate and read the innermost quantifier and to form the conjunction or disjunction of two copies of the propositional kernel provided the quantified variable has been replaced by 0 and 1 respectively in these copies. But since our *EDITRAM* program allows only literal strings as arguments in systematic replacement instructions we must program the later substitutions. We design therefore a subroutine which copies the string of characters representing variable  $x_i$  into a special purpose text-file (this encoding will include some binary representation of its index  $i$ ), and next subjects all occurrences of variables in the propositional kernel  $P$  to a treatment of systematic replacements which will turn all occurrences of  $x_i$  into a special pattern, and which will leave all other variables undisturbed. Then by substituting 0 or 1 for the special pattern, the required substitutions are obtained. For details of this subroutine see [123]. This completes the proof of the second inclusion.

□

In order to prove the first inclusion for the case of the vector machine and the *MRAM* a similar recursive procedure can be defined which evaluates the contents of a given bit of a given vector or a given bit of a given arithmetic register at some given time. In the vector machine the size of a vector grows exponentially but not worse, whereas for the *MRAM* all arithmetic registers may grow exponentially in length. Due to the presence of carries the simulation of a multiplication becomes as complicated as the case of a string replacement in the *EDITRAM*. Divisions have been reduced to multiplications inside the *MRAM* model itself at an earlier stage of the proof. Also the length of register addresses

remains bounded by the standard trick enabling the machine to use consecutive registers in its memory. In general these simulations achieve the required space-bound at the price of a huge consumption of time; the same values are computed over and over again by the recursion.

The proofs of the corresponding second inclusion for the cases of vector machines and *MRAM*'s invoke a direct simulation of the transitive closure algorithm by subroutines which build the matrix  $A(x, M, S)$  into a register and which compute its transitive closure by iterated squarings. A main ingredient is the programming of a routine which builds a bit-string consisting of the  $2^K$  bit-strings representing the first  $2^K$  integers, separated by markers, and of bit-strings to be used as masks for extracting a given bit-position from these integers in parallel in a single instruction. The idea of simplifying these proofs by invoking *QBF* as a *PSPACE*-complete problem was also used in [9].

### 3.4 Machines with true parallelism

In this section we consider models which provide observable parallelism by having multiple processors operate on shared data and/or shared channels. In these models computation can proceed either synchronously (all processors perform a step of the computation at the same time, driven by a local clock) or asynchronously (each processor computes at its own speed).

There exists a number of possible strategies for resolving the write conflicts which arise when several processors attempt to write in the same location of shared memory. In the *priority write* strategy the processor with the lowest index will succeed in writing in such a shared location and the other values will be lost. Other strategies which have been investigated are *exclusive write* (no two processors can write in the same global register at all), *common write* (if two processors try to write different values at the same time in the same register then the computation jams but writing the same value is permitted) , and *arbitrary write* (one of the writers becomes the winner but it is nondeterministically determined which one). For yet another approach to multiple writes see [25]. The computational power of the parallel *RAM* models based on these resolution strategies has been compared in [28] [29] [68] under the unrealistic assumption that the individual processors have arbitrary computing power. In these investigations the priority model has established itself as the most powerful one.

There are several methods of controlling the creation of parallel processors. Some models have a large or even infinite collection of identical processors which operate in parallel. In other models processors by their own action can create a finite number of new processors running in parallel, and in this way an arbitrary large tree of active processors can be activated as time proceeds.

Alltogether there are a large number of possible models and variations thereof. In this chapter I will discuss only a few of these models which indeed can be shown to obey the parallel computation thesis. As will become clear in chapter 4 some minor modifications of the models suffice to increase the power of these machines.

Curiously enough there exists no proposal for a parallel version of the *SMM*. It seems that this possibility has never been investigated from a complexity point of view. A parallel version of the *KUM* has been introduced in the Soviet literature already 25 years ago; this *Kolmogorov-Barzdin' machine* has been described by Barzdin' [4] [5] [6] but mainly for



its computational properties like the existence of a universal machine and the problem of giving a precise formulation of the involved parallel transformation of a graph (predating the literature on graph grammars). In two papers with Kalnin'sh [7] [8] which appeared 10 years later these computational investigations were continued. In this collection of papers there is one paper which in particular deals with complexity issues: in [6] it is established that a higher dimensional grid cannot be simulated on a lower dimensional one with constant factor overhead in space. As such the construction of a parallel version of the *SMM* which behaves like a second machine class member remains an interesting but doable exercise.

### 3.4.1 The *SIMDAG* model

In the *SIMDAG* model [39] (*Single Instruction, Multiple Data AGregate*) there exists a single global processor which can broadcast instructions to a potentially infinite sequence of local processors, in such a way that only a finite number of them are activated. The mechanism to keep the number of processors activated in a single step finite uses the signature of the local processors. Each local processor contains a read only register, called signature, containing a number which uniquely identifies this local processor. The global processor, in broadcasting an instruction includes a threshold value, and any local processor with a signature less than the threshold value transmitted performs the instruction while the others remain inactive.

Since the global processor can at most double the value of its threshold during a single step it follows that the number of sub-processors activated is bounded by an exponent in the running time of the *SIMDAG* computation.

The local processors operate both on local memory and on the shared memory of the global processor, where write-conflicts are resolved by priority; the local processor with the lower index becomes the winner in case of a write-conflict.

For a device like the *SIMDAG* it is necessary to restrict the power of the arithmetic instructions involved. Otherwise—as we will see in the sequel—the machine may become too powerful.

In the *SIMDAG* model the instruction repertoire for local and global processors involves additive arithmetic and parallel Boolean operations, combined with restricted shifting (division by 2). The writing of data in global storage by local processors is made conditional by stipulating that writing the value 0 is suppressed. In this way the *or* of a list of Boolean values computed by the local processors can be computed in a single write instruction by letting each local processor write a 1 in a fixed register in global memory if its bit equals 1, whereas the value 0 is not transmitted to the global memory.

Based on the above incomplete description I can sketch why the parallel computation thesis is true for the *SIMDAG* model.

**Theorem 16** *The SIMDAG model is a member of the second machine class*

**Proof:** The inclusion  $SIMDAG-NPTIME \subseteq PSPACE$  is shown by an argument similar to that used for the corresponding inclusion in case of the *EDITRAM*.

The trace of a nondeterministic *SIMDAG* computation can be guessed and be written down on a work-tape in polynomial space. Next one defines a pair of recursive procedures  $global(time, register)$  and  $local(time, register, signature)$  which evaluate to the value stored

at the given time in the given register of the global and given local processor respectively. The arguments of these recursive procedures can be written down in polynomial space (due to the restrictions on the arithmetics of the *SIMDAG*) and the recursion depth is bounded by the running time of the *SIMDAG* computation. Using these recursive procedures the guessed trace of the *SIMDAG* computation can be certified to be a correct accepting computation. As before the time needed for the simulation is very large due to the re-computation of intermediate results.

The converse inclusion  $PSPACE \subseteq SIMDAG-PTIME$  is shown by presenting an implementation of the transitive closure algorithm which runs in polynomial time.

This algorithm first loads the  $K \times K$  matrix  $A(x, M, S)$  which was introduced in chapter 3.1 in global memory. For convenience assume that  $K$  is a power of 2. The matrix is constructed by letting processor  $i + K \cdot j$  evaluate the entry  $A(x, M, S)[i, j]$ . By inspecting its signature the processor (using the Boolean operations and the division by 2 as a shift operator) can determine first the values of  $i$  and  $j$  and next unravel these values as bit-patterns in order to see whether the two encoded Turing machine configurations are equal or are connected by a single step. By a global write the matrix is loaded into global memory.

After formation of the matrix the transitive closure is computed by iterated squaring. Each squaring is computed by letting local processor  $i + K \cdot j + K \cdot K \cdot k$  read the values of  $A(x, M, S)[i, k]$  and  $A(x, M, S)[k, j]$ , form the *and* of these two values and write the result (conditionally) in  $A(x, M, S)[i, j]$ . This requires a constant number of steps.

After these squarings the existence of an accepting computation is determined by the global processor by inspecting the proper matrix entry of  $A(x, M, S)$ .  $\square$

### 3.4.2 The array processing machine

Our next model, the *array processing machine* (*APM*) was proposed by van Leeuwen and Wiedermann [134]. It has been inspired by the contemporary vectorized supercomputers. This machine has the storage structure of an ordinary *RAM* but contains besides the traditional accumulator also a vector accumulator which consists of a potentially unbounded linear array of standard accumulators.

The array processing machine combines the instruction set of a standard *RAM* with a new repertoire of vector instructions which operate on the vector accumulator. These instructions allow for reading, writing, transfer of data and arithmetic on vectors of matching size which consist of consecutive locations in storage and/or an initial segment of the vector accumulator.

Each operation on the vector accumulator destroys its previous content. Conditional control on a vector operation is possible by the use of a mask which consists of an array of Boolean values (0 or 1) of the same size as the vector operands; the vector instruction now is performed only at those locations corresponding to occurrences of 1 in the mask. A complete address for a vector operation therefore may consist of four integers: lower and upper bounds of the vector argument and the mask respectively.

The power of parallelism is provided to the model by the time measure used: uniform time or logarithmic time, where every vector instruction is charged according to its most expensive scalar component. So in a vector *LOAD* the logarithmic time complexity is

proportional to the logarithm of the upper bound of the operand and/or mask plus the logarithm of the largest value loaded into the vector accumulator.

In their paper [134] the authors prove that the above array processor is a member of the second machine class by providing mutual simulations with respect to the *SIMDAG*, where it turns out that the simulations require polynomial (more specifically  $n^4$ ) overhead. In both directions the simulations require non trivial programming techniques, and a parallel  $O(\log^2(n))$  implementation of Batchers's sort is an essential element of the simulation.

Inspection of the model shows that a proof that the *APM* obeys the parallel computation thesis can be easily obtained by the techniques used for other devices in this chapter.

**Theorem 17** *The APM model is a member of the second machine class*

**Proof:** To prove  $PSPACE \subseteq APM-PTIME$  one can show that *QBF* can be solved in polynomial time on an *APM* ; here the main ingredient is the construction of  $n$  vectors in storage of length  $2^n$  where vector  $j$  contains the value of bit  $j$  in the binary representation of the numbers  $0 \dots 2^n - 1$ . Using those vectors one can evaluate a given propositional kernel in linear time using the vector instructions , and the resulting vector can be folded together according to the quantifiers in order to provide the final result.

For the converse inclusion:  $APM - NPTIME \subseteq PSPACE$  the usual technique of writing down a computation trace and certifying it by means of a recursive procedure will work. The detour via *PSPACE* implies the existence of mutual polynomial time overhead simulations of *SIMDAG* and *APM* but it does not provide explicit simulation overheads as indicated above.  $\square$

### 3.4.3 Models with recursive parallelism

As an example of a parallel machine of a different character I mention the recursive Turing machine introduced by Savitch [106]. In this model every copy of the device can spawn off new copies which start computing in their own environment of work-tapes, and which communicate with their originator by means of channels shared by two copies of the machine.

A similar model, based on the *RAM* is the  $k-PRAM$  described by Savitch and Stimson [108]. In this model a *RAM*-like device can create up to  $k$  copies of itself, which start computing in their local environment, while their creator is computing onwards. These copies themselves can create new offspring as well. Data are transmitted from parent to child at generation time by loading parameters in the registers of the offspring. Upon termination a child can return a result to its parent by writing a value in a special register of the parent which is read-only for the parent. The parent can inquire about this register whether his child has already written into it; if he reads the register before the child has assigned a value to it, the parent's computation is suspended. Using these features a parent can activate a number of children and consume the first value which is returned, aborting the computations of the remaining children which have not yet terminated.

Clearly models based on such local communication channels are much slower in broadcasting information to a collection of sub-processors and in gathering the answers. However, the validity of the parallel computation thesis is not disturbed by this delay due to the fact that the slowdown is at worst polynomial: in order to activate an exponential

number of processors by spawning off sub-processors polynomial time suffices in case a complete binary tree is formed; the time consumed for activating this tree usually is polynomial in terms of the time consumed by writing down the data to be processed by the sub-processors.

It is not difficult to see that the above models support an implementation of the recursive version of the transitive closure algorithm from chapter 3.1 which runs in polynomial time. This explains why *PSPACE* is contained in the  $\| PTIME$  class for these models. The time overhead which used to be  $n^2$  before becomes now  $n^3$  due to the fact that the recursive algorithm contains also a loop over all intermediate nodes. This loop would require exponential time if performed sequentially, and therefore it is replaced by a recursive expansion where a call at order  $w$  produces an exponential number of calls at order  $w - 1$ . The recursion depth of the procedure now becomes  $\log^2(n)$  instead of  $\log(n)$ .

The same recursive expansion trick is used in order to show that for these models  $\| NPTIME \subseteq \| PTIME$ . A single nondeterministic computation of length  $T$  is replaced by  $2^T$  deterministic computations relative to a choice string which is different for every copy. Transforming this idea into a complete proof exposes a number of complications. The final proof for this inclusion yields an  $(n^6)$  time overhead for elimination of nondeterminism. For the details I refer to the papers [106] [108].

A crucial difference between the *SIMDAG* model and the recursive models is that in the *SIMDAG* model the local processors, if they are active at all at some time, all execute the same instruction on data which may be different. As a consequence it is possible to write down the trace of executed instructions of a *SIMDAG* computation in polynomial space. In the recursive models each processor, once being activated, performs its own program except for the impact of communication with its parent or its offspring. It becomes therefore impossible to write down the complete computation trace in polynomial space if an exponential number of processors is activated.

This has consequences for the proof of the inclusion  $\| PTIME \subseteq PSPACE$  for these models. Rather than writing down the entire trace of the computation and certifying it by a recursive procedure, the entire genealogical tree of machine incarnations is searched by a recursive procedure. Recursion depth is bounded by the polynomial bound on the running time of the parallel machine. A polynomial bound on the size of the recursive stack frames is obtained from the fact that the machine has no powerful arithmetic. It is crucial at this place that the communication between parent and offspring is entirely local; if offspring can write in global memory the proof breaks down.

A further complication results from the fact that a parent can abort some of its children before they have terminated; evaluating the result produced by such a child may lead to a divergent computation. Therefore all devices are extended with a counter keeping track of global synchronous time. Such counters are needed anyhow in order to let a parent know which one of its two children was the first to terminate. Again I refer for more details to the full papers [106] [108].

One obtains the feeling that the models in this section are less "stable" than the powerful sequential models or the synchronous *SIMDAG*. This is also illustrated by the results in the final chapter of this report.

## Chapter 4

# Parallel machine models outside the second machine class

In this chapter I discuss several models which can not be classified as being either first or second class devices. First I discuss a parallel model which is much weaker than the models in the previous chapter. Next I will discuss the models which are more powerful than the second machine class members, including a model claimed to perform all computations in constant time.

### 4.1 A weak parallel machine

The *parallel Turing machine*, (abbreviated *PTM*), introduced by Wiedermann [140] should not be confused with the devices introduced by Savitch under the name recursive Turing machines [106]. In both cases one considers a Turing machine with a nondeterministic program where a choice of possible successor states leads to the creation of several devices, each continuing in one of the possible configurations. Where in Savitch's model the entire configuration is multiplied, it is the case that in Wiedermann's model only the finite control and the heads are multiplied, thus leading to a proliferation of Turing automata all operating on the same collection of tapes.

Wiedermann's device consists of a finite control with  $k$   $d$ -dimensional work-tapes, the first of which contains the input at the start of the computation. Each control has one head on every tape. The program of the device is a standard nondeterministic Turing program for a machine with  $k$   $d$ -dimensional single-head tapes; however, instead of choosing a next state when facing a nondeterministic move the machine creates new copies of its control and heads, which go on computing on the same tapes. There are no read conflicts; write conflicts are resolved by the common write strategy: if two heads try to write different symbols at the same square the computation aborts and rejects; if two heads try to write the same symbol his symbol is written and the heads move on.

An accepting computation is a computation where every finite control that is created during the computation halts in an accepting state.

The crucial observation which makes this model weaker than the true members of the second machine class which we have met in the previous chapter is related to the achievable degree of parallelism. Although the machine can activate an exponential number of copies

of itself in polynomial time, these copies operate on the same tapes and therefore only a polynomial number of essentially different copies will be active at any moment in time. If two finite controls are in the same internal state and have their heads positioned on the same tape squares, their behavior will be equal from that time onwards; hence these two controls actually merge into a single unit. This leads to an upper bound of  $q \cdot S^k$  on the number of different controls being active at the same time, where  $q$  denotes the number of states in the program and  $S$  denotes the space used by the device. Since space is bounded by time this leads to a polynomial bound on the number of different copies.

Based on this observation it becomes possible to simulate this parallel machine by a standard deterministic Turing machine with polynomial time overhead: the simulator maintains on some additional work-tape a list of all active finite controls with their head positions, and by maintaining a pair of old and new work-tapes the machine can process the updates of each control in sequence, taking care of the needed multiplication of controls and checking for write conflicts.

It can not be inferred from this simulation that the device is a standard first class machine, since it is not clear whether the above simulation can be modified in such a way that the space overhead becomes a constant factor. For the special case of a single tape parallel machine a constant factor space overhead is achieved by storing with each tape cell the set of states achieved by heads scanning this cell; this set (being a subset of the fixed set of states of the machine) can be written down in an amount of space which is independent of the length of the input. But for the case of more tapes or more heads on a single tape it is not sufficient to mark the tape cells by the states in which they are scanned, since one must also know which heads belong together to a single finite control, and this requires the encoding of head positions for each device. Therefore the naive simulation as indicated above requires space  $O(S^k \cdot \log(S))$ .

Wiedermann observes that one can recognize the language of palindromes with a *PTM* with two one-dimensional tapes, where the first tape is a read-only input tape, in space  $O(1)$  if the space of the input tape is not counted. Still the heads on the read-only input tape may multiply, thus storing information on this tape by their positions. It seems therefore unlikely that a constant factor space overhead simulation is possible since palindromes cannot be recognized in constant space by a standard machine. This particular example, however, seems to depend on the particular interpretation of the space on the input tape not being counted, and an example where the input head is common for all copies of the finite control seems to be required for a more convincing separation result.

These simulations show that  $P = PTM-PTIME$  and that  $PSPACE = PTM-PSPACE$ ; the *PTM* model therefore does not obey the parallel computation thesis, unless  $P = PSPACE$ .

The *PTM* model has the interesting property that for several practical problems in  $P$  an impressive speed-up by pipelining can be achieved, even though the device is not a second machine class member. For details I refer to [140].

## 4.2 Beyond the second machine class

There exist in the literature several models which seem to be more powerful than a second class machine. Fortune and Wyllie [31] have described a hybrid of the *SIMDAG* with a parallel machine based on branching, called *P-RAM* in [31] and called *MIMD-RAM* in my

earlier survey [131]. For this model it has been established that  $MIMD-RAM-PTIME = PSPACE$  and  $MIMD-RAM-NPTIME = NEXPTIME$ .

Savitch [107] has defined a model based on an  $MRAM$  which creates parallel copies by branching, called the  $LPRAM$ , and proves that  $LPRAM-PTIME = PSPACE$  and  $LPRAM-NLOGTIME = NP$ .

N. Blum has written a short note arguing against the parallel computation thesis [12]. His claim is based on a pair of models called the  $PRAM$  and the  $WRAM$  respectively. These models resemble the  $SIMDAG$  but instead of the priority strategy for write conflicts he uses the exclusive write strategy in the  $PRAM$  and common write strategy in the  $WRAM$ . For these models he claims the results  $NEXPTIME \subseteq WRAM-PTIME$  and  $EXPTIME \subseteq PRAM-PTIME$ . The correctness of his claims is questionable since his proof requires the presence of powerful arithmetic as well.

Below I give a sketch of these models and the simulations involved.

#### 4.2.1 The $MIMD-RAM$

The  $MIMD-RAM$  which was introduced by Fortune and Wyllie [31] as a hybrid between the  $SIMDAG$  and the  $k-PRAM$  described in chapter 3.4. In the  $MIMD-RAM$  a machine can create offspring by forking, where the offspring will perform the instructions of its own program. Upon creation the offspring machine will start executing at the first instruction in its program. The subsequent course of the computation is influenced by the initial value of the accumulator which has been set by its creator. The machine has standard additive  $RAM$  arithmetic.

Each offspring processor has its own local memory. The machines communicate furthermore by global memory. Global and local memories are standard  $RAM$  memories. Each device can both read and write in its local memory using the standard  $RAM$  instructions. For passing information upwards processors write in global memory. Simultaneous reads from global memory are allowed, but simultaneous writes are prohibited. The machine accepts if the oldest copy accepts by halting with 1 in its accumulator.

In the model a special input convention is used which supports the reading of an input of length  $n$  in time  $O(\log(n))$ , so sublinear running times become meaningful.

The  $MIMD-RAM$  has the property that its deterministic version is a true second machine class device whereas its nondeterministic version provides us with a full exponential speed-up in time. This is expressed by the following result

**Theorem 18** *The  $MIMD-RAM$  model obeys the following equations*

1.  $MIMD-RAM-PTIME = PSPACE$  and
2.  $MIMD-RAM-NPTIME = NEXPTIME$ .

Before starting the proof I first introduce an  $NP$ -complete problem which we will use in order to establish one of the required inclusions. This is the problem  $BOUNDED TILING$  [103], called  $SQUARE TILING (GP13)$  by Garey and Johnson [37]

**BOUNDED TILING:**

**INSTANCE:** a finite set  $W$  of tiles (squares with colors given on their edges), and an  $N \times N$  square  $V$  with a given coloring on the  $4 \cdot N$  unit edge segments

on the border of  $V$

*QUESTION*: is it possible to tile the square  $V$  with copies of the tiles in  $W$  (without rotations or reflections) such that each pair of adjacent tiles have matching colors on their common edge, and such that the tiles adjacent to the border of the square  $V$  have colors matching the given coloring of the border on their exterior edges ?

**Proof**: The proofs for these equalities can be sketched as follows:

*PSPACE*  $\subseteq$  *MIMD-RAM-PTIME*: one proof consists of simulation of a 2-*PRAM* where the channels have been replaced by global registers. Another proof consists of a version of the transitive closure algorithm for the computation graph. Simultaneous writes are prevented by designing a fan-in algorithm for evaluation of the *and* or *or* of  $O(S)$  bits. In both cases the resulting program is a deterministic *MIMD-RAM* program.

*MIMD-RAM-PTIME*  $\subseteq$  *PSPACE*: Due to the restrictions on the arithmetic in time  $T$  at most  $2^T$  machines are created, each operating on at most  $2^T$  registers and operate on values bounded by  $2^T$ . This implies that all arguments for a recursive procedure  $val(i, j, t)$ , which yields the value of register  $j$  in device  $i$  at time  $t$ , can be written down in polynomial space. A similar procedure can be obtained for the contents of the global memory. However, it is no longer possible to write down in polynomial space the entire trace of the computation since each device performs its own instructions. Instead the trace of the computation is built together with the tree of recursive calls; a certain value is present in some register as a result of the execution of previous instructions which have been executed because of the presence of certain values in other registers at some earlier time etc. It is not hard to see that an alternating *RAM* can evaluate this recursion guessing and certifying the trace at the same time. From the fact that the *MIMD-RAM* involved is deterministic it follows that only the true values and instructions performed can be certified. It also follows that this alternating *RAM* consumes time polynomially bounded by  $T$ . Hence the above inclusion follows.

*MIMD-RAM-NPTIME*  $\subseteq$  *NEXPTIME* and *MIMD-RAM-NLOGTIME*  $\subseteq$  *NP*: These inclusions can be shown by brute force simulations: the overhead in simulating a nondeterministic *MIMD-RAM* computation by writing down an entire record of its computation is at most exponential in time.

*NP*  $\subseteq$  *MIMD-RAM-NLOGTIME* and *NEXPTIME*  $\subseteq$  *MIMD-RAM-NPTIME*: These inclusions are shown by designing a nondeterministic *MIMD-RAM* algorithm which recognizes an *NP*-complete problem in logarithmic time. This yields the first inclusion and the second will follow by a similar argument or by a padding argument. For the *NP*-complete problem we select the problem *BOUNDED TILING*.

An instance of *BOUNDED TILING* can be solved by a nondeterministic *MIMD-RAM* which operates as follows: first it creates  $N^2$  copies, each covering a single unit square inside  $V$ . This requires time  $O(\log(N))$ . Next in a nondeterministic move each processor guesses the tile to be placed on its square. This is the unique nondeterministic move in the entire program. Subsequently, using the global memory, each processor exchanges its tile with the processors representing its neighbors, and certifies whether the choices match. The result is obtained by a standard fan-in communication of the bits computed in this way. Evaluation of the match requires constant time; the fan-in procedure requires  $O(\log(N))$  steps.  $\square$



### 4.2.2 The *LPRAM*

The *LPRAM* is a model introduced by Savitch [107]. It is a hybrid between the  $k$ -*PRAM* and the *MRAM*, since it combines the recursive parallelism of the  $k$ -*PRAM* with the powerful arithmetic of the *MRAM*. Communication is by channels only as in the  $k$ -*PRAM*.

**Theorem 19** *The LPRAM satisfies the equations*

1.  $LPRAM-PTIME = PSPACE$  and
2.  $LPRAM-NPTIME = NEXPTIME$ .

**Proof:**  $PSPACE \subseteq LPRAM-PTIME$  : trivial; don't use the vector instructions.

$LPRAM-PTIME \subseteq PSPACE$  : by a standard guess and verify method. The same alternating verifier can be used as in chapter 4.2.1 , but since the vector instructions allow for exponential growth of values the certification of values should proceed at the bit level rather than at the register level. This is the same trick as is used for the inclusion  $MRAM-NPTIME \subseteq PSPACE$  in the papers by Pratt and Stockmeyer [94] and Hartmanis and Simon [44] [45].

$LPRAM-NPTIME \subseteq NEXPTIME$  and  $LPRAM-NLOGTIME \subseteq NP$  : These inclusions again are obtained by a brute force simulation.

$NP \subseteq LPRAM-NLOGTIME$  and  $NEXPTIME \subseteq LPRAM-NPTIME$ : These inclusions are again proved by implementing a nondeterministic *LPRAM*-algorithm for *BOUNDED TILING* which runs in logarithmic time: the machine first creates a tree of  $N^2$  offspring processors which guess a tile for some unit square inside  $V$ . These guesses are communicated to the oldest processor by a fan-in procedure, where along the way the exponential growing information is stored using the vector instructions. Next the same information is distributed once more over  $N^2$  processors in such a way that each processor obtains a cell with its four neighbors. These  $N^2$  processors then check whether the tiles match and finally the results are communicated to the oldest processor by a bounded fan-in procedure.  $\square$

### 4.2.3 Extending the *SIMDAG* with powerful arithmetic

In this section we consider a hybrid of a *SIMDAG* with the *MRAM*. To my knowledge such a model has never been formally introduced in the literature, but implicitly it seems to have been considered.

Starting point here is a more refined analysis of the running time of the transitive closure algorithm in chapter 3.3.1. Such an analysis shows that the total running time consists of three contributions:

1.  $O(\log(K)) = O(S)$  for evaluation of the matrix size  $K$
2.  $O(S)$  for unraveling the configurations and computing  $A(x, M, S)$
3.  $O(\log(T)) = O(S)$  for computing the transitive closure of  $A(x, M, S)$

This more refined analysis explains why the power of the arithmetic in the model is a crucial factor in establishing whether the model is a second machine device or not. Assume that we can use multiplication in unit time, the first contribution is reduced to  $O(\log \log(K)) = O(\log(S))$ ; given more powerful parallel Boolean instructions the unraveling can be distributed over  $O(S)$  processors and therefore the contribution for step 2 becomes  $O(\log(S))$  as well. As a consequence for the resulting *PSIMDAG* (for powerful *SIMDAG*) model one obtains  $NEXPTIME \subseteq PSIMDAG-PTIME$ . Hence it is unlikely that such a model obeys the parallel computation thesis, unless  $PSPACE = NEXPTIME$ .

The above ideas form the basis for the objections against the parallel computation thesis as put forward by N. Blum [12]. However, in his paper he only considers the third phase of the above algorithm: the computation of the transitive closure itself in  $\log(T)$  steps. It should be clear that in a model where addition is the only arithmetic available steps 1 and 2 still will require time  $O(S)$ . But given a suitable set of powerful arithmetical instructions the required speed-up of steps 1 and 2 from  $O(S)$  to  $O(\log(S))$  become possible. It seems therefore that his model requires such strong arithmetical instructions, but this makes his claims much weaker, given the fact that models showing this kind of behavior had been proposed before by Fortune and Wyllie and by Savitch.

Remember that Blum's models resemble the *SIMDAG* but instead of the priority strategy for write conflicts the exclusive write strategy is used in the *PRAM* and common write strategy in the *WRAM*. For these models Blum claims the following results:

**Theorem 20** *When equipped with powerful arithmetic the Blum models satisfy:*

1.  $NEXPTIME \subseteq WRAM - PTIME$
2.  $EXPTIME \subseteq PRAM - PTIME$ .

**Proof:** Taking the computation of the transition matrix for granted, given the powerful arithmetic, the first inclusion claimed by Blum:  $NEXPTIME \subseteq WRAM-PTIME$ , follows from the fact that in the transitive closure algorithm only 0's and 1's need to be written, and the writing of a 0 can be suppressed. So in the case of a multiple write the writes will be consistent. Since a bound  $T$  on the running time of the simulated machine is given we know that instead of the transitive closure the  $T$ -th power of the transition matrix can be computed.

The second inclusion:  $EXPTIME \subseteq PRAM-PTIME$ , follows by realizing if the machine is deterministic all powers of the transition matrix have the property that every row contains at most a single non-zero entry. From this it follows that during a squaring every matrix element  $A(x, M, S)[i, j]$  will be written by at most one processor. Note that also step 2 of the algorithm must be modified in order to have the required fan-in of  $O(S)$  bits of local information into a single matrix element: this can be obtained at the price of an additional time of  $O(\log(S))$ . So Blum's construction is valid on the *PRAM* for deterministic computations.  $\square$

#### 4.2.4 Arbitrary computations in constant time

Intuition in computation theory states that, whatever model of a computing device one selects, it should always remain true that running time on this model behaves like a

complexity measure [10]. Hence in particular it should never be possible to perform arbitrary complex computations in constant time.

A model which violates this very mild restriction has been described in the literature. We find the model in the context of an analysis of the correct interpretation of the parallel computation thesis by Ian Parberry. These ideas, which were originally presented in his Ph.D. Thesis [83], have been published recently in SIGACT News [84]; for similar results with slightly improved resource bounds see [85].

In his discussion Parberry considers a parallel *RAM* model which runs fully synchronously. The number of active processors is determined at the start of the computation. Each processor has beside the standard additive arithmetic some shifts or multiplication like instructions which enable processors to extract bits from bit-strings at arbitrary locations. The processors communicate via a global memory. Each processor has a read-only register initialized at its index.

Resources considered are time  $T$ , storage  $S$  (counted in *RAM* words), word-size  $W$  (both in local and global storage), and number of processors  $n$ :  $T(n)$ ,  $S(n)$ ,  $W(n)$  and  $P(n)$ .

The model resembles the *SIMDAG* with two major differences: in the first place, rather than having all processors being activated by a single central processor, all processors with index  $\leq P(n)$  are activated at time  $t = 0$  and remain so till all have halted. The result of the computation can be found by inspecting register 0 in global memory. The second difference is that the model is non-uniform: some quantities like  $P(n)$  and derived quantities are given to all processors in advance.

Having done so we have given the machine sufficient power to perform the transitive closure algorithm from chapter 3.3.1 in constant time. Remember the three contributions of the running time analysis introduced in chapter 4.2.3 above; here  $T'(n)$  and  $S'(n)$  denote the time and space bounds on the Turing machine computations on input  $x$  with  $|x| = n$ , and  $K(n)$ , the size of the transition matrix  $A(x, M, S'(n))$  equals  $2c' \cdot S'(n)$  for a suitable constant  $c'$ .

1.  $O(\log(K(n))) = O(S'(n))$  for evaluation of the matrix size  $K(n)$
2.  $O(S'(n))$  for unraveling the configurations and computing  $A(x, M, S'(n))$
3.  $O(\log(T'(n))) = O(S'(n))$  for computing the transitive closure of  $A(x, M, S'(n))$

Step 1 becomes  $O(1)$  due to the non-uniformity involved. Quantities like the space  $S'(n)$  and time  $T'(n)$  for the Turing machine to be simulated are given in advance, as are the number of processors  $P(n)$  which will be invoked.

Step 2 and step 3 are combined. In the simulation by Blum a team of  $O(S'(n))$  processors is called upon in order to decompose the index of processor  $i \cdot K(n) + j$  into bits and pieces in order to determinate whether  $A(x, M, S'(n))[i, j]$  equals 0 or 1. So the index of a processor is analyzed whether it represents a legal transition between two successive transitions of the machine which is simulated. In Parberry's simulation this idea is carried to the extreme: the index of a processor is analyzed whether it represents the coding of an entire accepting computation. The length of a string coding of such a computation becomes  $O(S'(n) \cdot T'(n))$ , hence the total number of indices analyzed in this way becomes  $2^c \cdot T'(n) \cdot S'(n)$  for some constant  $c$ . The analysis of one such index can be performed by

breaking it into bits and pieces: given a coding like the one used in time-space diagrams of Turing machines, a team of  $T'(n) \cdot S'(n)$  processors can perform the comparisons of the coded configurations at time  $t$  and time  $t + 1$  in time  $O(1)$ ; as Parberry indicates  $T$  processors in a single team already suffice. One must check furthermore that the initial configuration at time  $t = 0$  represents a proper configuration on the real input but that can be done in a similar way.

The result obtained in this way is that the simulation time becomes  $O(1)$  at the price of using  $P(n) = 2^{c \cdot T'(n) \cdot S'(n)} \cdot T'(n)$  processors and word-size  $W(n) = O(T'(n) \cdot S'(n))$ .

By a similar proof Parberry shows that if you can simulate with word-size  $W(n)$  a  $B(n)$  time bounded deterministic Turing machine in time  $B'(n)$  on his model one can use this simulation as a tool for speeding-up a  $T'(n)$  time-bounded deterministic Turing machine computation with word-size  $O(W(n) + B(n) + \log(T(n)))$  to time  $O(T(n)/B(n) + B'(n))$ ; the second term in the time analysis of the simulation is the time required to set up a transition matrix structure for computing  $B(n)$  steps in time  $O(1)$  by table look-up, whereas the first term measures the time needed for simulation of  $T(n)$  steps in blocks of  $B(n)$  steps each.

This brings Parberry to his analysis of how the above simulations support in fact the parallel computation thesis: a “reasonable” parallel device should satisfy the constraint that  $W(n)$  is polynomially bounded by  $T(n)$  and that (as a consequence)  $P(n)$  is bounded by some expression  $2^{c \cdot T(n)^k}$  for some constant  $k$ . The result above shows that within these bounds arbitrarily large polynomial speed-ups can be obtained with “reasonable” devices, but the speed-up to constant time is possible only at the price of violating this “reasonability” criterion.

To my opinion Parberry’s result shows once more the hideous power of shifts and multiplicative instructions. Also the non-uniformity of the model makes it suspect. If quantities like  $T'(n)$ ,  $S'(n)$  or  $P(n)$  must be computed from the input the  $O(1)$  time bound is invalidated.

# Acknowledgements

This chapter originates from the confluence of two research interests. The idea of introducing the first machine class has its origin in a private discussion I had seven years ago with L. Torenvliet, when I tried to convince him of the self evidency of what later would become the invariance thesis. The truth about the state of affairs with respect to the traditional space measure for the *RAM* was uncovered in close cooperation with C. Slot. The exchange of ideas with J.P. Schmidt last summer was most useful for clarifying the issue of the true power of restricted unit time multiplication in the *RAM*. For the case of the *SMM* valuable discussions with A. Schönhage and M.C. Loui and D.R. Luginbuhl are acknowledged for. The difference between the orthodox and the liberal interpretation of the invariance thesis was first pointed out by Wiedermann.

The interest in parallel machines which led to the introduction of the second machine class (which predates the first machine class by about one year) was provoked by the invitation by J. van Leeuwen and P.M.B. Vitányi for presenting an overview on models of parallelism at a local seminar in Utrecht [131]. This work underwent subsequent revisions during presentations at the Banach seminar on foundations of computation theory in 1985 in Warsaw [132].

Valuable comments on preliminary versions of the chapter were presented by L. Stockmeyer, L. Torenvliet, J. van Leeuwen, P.M.B. Vitányi, J. Wiedermann, and the students of the machine models class which has been organized during the fall season of the past three years. I apologize for not yet having written down the second machine class version of the parallel *SMM* which these students have invented. I am grateful to R. Freivalds for providing me the references to the Kolmogorov Barzdin' machine. The inclusion among the Physics section of the translation volumes of the Doklady may have contributed to the lack of attention these papers have received in the West. Finally the appearance of the Wagner and Wechsung book saved me from having to put together the encyclopaedic survey I envisioned when I started drafting this chapter. Looking back I know that I could never have completed it along that road within the 60 pages allowed to me.

Finalizing this chapter would have been impossible without the real-time assistance of L. Torenvliet during the process of having the Mac Write document transformed into  $\text{\LaTeX}$  and having the result being processed on our machine.

# Bibliography

- [1] Aho, A.V., Hopcroft, J.E. and Ullman, J.D., *The Design and Analysis of Computer Algorithms*, Addison Wesley (1974)
- [2] Andreev, A.E., *On a method for obtaining lower bounds for the complexity of individual monotone functions*, Soviet Math. Dokl. 31 (1985) 530–534
- [3] Balcázar, J.L., Díaz, J. and Gabarró, J., *Structural complexity I*, EATCS Monographs on Theoretical Computer Science 11, Springer Verlag (1988)
- [4] Barzdin', Ya. M., *Universal pulsing elements*, Soviet Physics-Doklady 9 (1965) 523–525
- [5] Barzdin', Ya. M., *Universality problems in the theory of growing automata*, Soviet Physics-Doklady 9 (1965) 535–537
- [6] Barzdin', Ya. M., *Capacity of the medium and behavior of automata*, Soviet Physics-Doklady 10 (1966) 8–11
- [7] Barzdin' Ya. M. and Kalnin'sh, Ya. Ya., *On a language for the transformation of graphs intended for the specification of automata*, Automatika i Vychislitel'naya Teknika 7(5) (1973) 22–28
- [8] Barzdin' Ya. M. and Kalnin'sh, Ya. Ya., *a Universal automaton with variable structure*, Automatika i Vychislitel'naya Teknika 8(2) (1974) 9–17.
- [9] Bertoni, A., Mauri, G. and Sabadini, N., *Simulations among classes of random access machines and equivalence among numbers succinctly represented*, Ann. Discr. Math. 25 (1985) 65–90
- [10] Blum, M., *A machine-independent theory of the complexity of recursive functions*, J. Assoc. Comput. Mach. 14 (1967) 322–336
- [11] Blum, N., *A Boolean function requiring  $3n$  network size*, Theor. Comp. Sci. 28 (1984) 337–345
- [12] Blum, N., *A note on the "Parallel Computation Thesis"*, Inf. Proc. Letters 17 (1983) 203–205
- [13] Böhling, K.H. and von Braunmühl, B., *Komplexität bei Turingmaschinen*, B.I. Wissenschaftsverlag, Reihe Informatik 14, 1974

- [14] Borodin, A., *On relating time and space to size and depth*, SIAM J. Comput. 6 (1977) 733–744
- [15] Borodin, A., Fisher, M., Kirkpatrick, D. G., Lynch, N.A. and Tompa, M., *A time-space tradeoff for sorting on non-oblivious machines*, J. Comput. System Sci. 22 (1981) 351–364
- [16] Chandra, A.K., Kozen, D.C. and Stockmeyer, L.J., *Alternation*, J. Assoc. Comput. Mach. 28 (1981) 114–133
- [17] Chazelle, B. and Monier, L., *Unbounded hardware is equivalent to deterministic Turing machines* J. Theor. Comp. Sci. 24 (1983) 123–120
- [18] Cook, S.A., *The complexity of theorem proving procedures*, Proc. ACM Symposium Theory of computing 3 (1971), 151–158
- [19] Cook, S.A., *Towards a complexity theory of synchronous parallel computation*, Enseign. Math. 27 (1980) 99–124
- [20] Cook, S.A., *The classification of problems which have fast parallel algorithms*, in M. Karpinski (ed.), Proc. Fundamentals of Computation Theory 1983, Borgholm Sweden, Springer Lecture notes in Computer Science 158 (1983), 78–93
- [21] Cook, S.A. and Reckhow, R.A., *Time bounded random access machines*, J. Comput. Syst. Sci. 7 (1973) 354–375
- [22] Davis, M., *Computability and Unsolvability*, Mc Graw Hill, 1958
- [23] Davis, M., *Why Gödel didn't have Church's thesis*, Inf. and Control 54 (1982) 3–24
- [24] Dymond, P.W. and Cook, S.A., *Hardware complexity and parallel computation*, Proc. IEEE Foundations of computer science 21 (1980), 360–372
- [25] Dymond, P.W. and Ruzzo, W.L., *Parallel RAMs with owned global memory and deterministic context-free language recognition*, Proc. Int. Colloquium on Automata Languages and Programming 13 , Springer Lecture notes in Computer Science 226 (1986), 95–104
- [26] Edmonds, J., *Paths, trees, and flowers*, Canad. J. Math. 17 (1965), 449–467
- [27] Even, S., and Tarjan, R.E., *A combinatorial problem which is complete in polynomial space*, J. Assoc. Comput. Mach. 23 (1976) 710–719
- [28] Fich, F.E., Radge, P.L. and Widgerson, A., *Relations between concurrent-write models of parallel computation*, SIAM J. Comput. 17 (1988) 606–627
- [29] Fich, F.E., Meyer auf der Heide, F., Radge, P.L. and Widgerson, A., *One, two, three ... infinity: lower bounds for parallel computation*, Proc. ACM Symposium Theory of computing 17 (1985), 48–58
- [30] Fischer, P.C., Meyer, A.R. and Rosenberg, A.L., *Real-time simulation of multihead tape units*, J. Assoc. Comput. Mach. 19 (1972) 590–607

- [31] Fortune, S. and Wyllie, J., *Parallelism in random access machines*, Proc. ACM Symposium theory of computing 10 (1978), 114–118
- [32] Fraenkel, A.S., Garey, M.R., Johnson, D.S., Schaefer, T. and Yesha, Y., *The complexity of checkers on an  $N \times N$  board, preliminary report*, Proc. IEEE Foundations of computer science 19 (1978), 55–64
- [33] Fraenkel, A.S. and Lichtenstein, D., *Computing a perfect strategy for  $n \times n$  chess requires time exponential in  $n$* , J. Combin. Theory 31 (1981) 199–214
- [34] Fredman, F.L., Komlós, J. and Szemerédi E., *Storing a sparse table with  $O(1)$  worst case access time*, Proc. IEEE Foundations of computer science 23 (1982), 165–169
- [35] Fredman, F.L., Komlós, J. and Szemerédi E., *Storing a sparse table with  $O(1)$  worst case access time*, J. Assoc. Comput. Mach. 31 (1984) 538–544
- [36] Galil, Z. and Paul, W., *An efficient general-purpose parallel computer*, J. Assoc. Comput. Mach. 30 (1983) 360–387
- [37] Garey, M.S. and Johnson, D.S., *Computers and Intractability, a Guide to the Theory of NP-completeness*, Freeman, 1979
- [38] Gilbert, J.R., Lengauer, T. and Tarjan, R.E., *The pebbling problem is complete in polynomial space*, SIAM J. Comput. 9 (1980) 513–524
- [39] Goldschlager, L.M., *A universal interconnection pattern for parallel computers*, J. Assoc. Comput. Mach. 29 (1982) 1073–1086
- [40] Gruska, J., *Systolic Automata - Power, Characterizations, Nonhomogeneity*, Proc. Mathematical foundations of computer science '84, Springer Lecture notes in Computer Science 176 (1984), 32–49
- [41] Halpern, J.Y., Loui, M.C., Meyer, A.R. and Weise, D., *On Time Versus Space III*, Math. Syst. Theory 19 (1986) 13–28
- [42] Harel, D., *Recurring dominos: Making the highly undecidable highly understandable*, Ann. Disc. Math. 24 (1985) 51–72
- [43] Hartmanis, J., *Observations about the development of theoretical computer science*, Ann. Hist. Comp. 3 (1981) 42–51
- [44] Hartmanis, J. and Simon, J., *On the power of multiplication in random access machines*, Proc. IEEE Switching and automata theory 15 (1974), 13–23
- [45] Hartmanis, J. and Simon, J., *On the structure of feasible computations*, in Rubinoff, M. and Yovits, M.C. (eds.), Advances in Computers 14, Acad. Press 1976, pp. 1–43
- [46] Hartmanis, J. and Stearns, R.E., *On the computational complexity of algorithms*, Trans. Amer. Math. Soc. 117 (1966) 285–306
- [47] Hemmerling, A., *On the space complexity of multidimensional Turing automata*, E.-M. Arndt Univ. Greifswald, Preprint 2, 1979



- [48] Hennie, F.C. and Stearns, R.E., *Two-way simulation of multi-tape Turing machines*, J. Assoc. Comput. Mach. 13 (1966) 533–546
- [49] Hopcroft, J., Paul, W. and Valiant, L., *On time versus space and related problems*, Proc. IEEE Foundations of computer science 16 (1975) 57–64
- [50] Hopcroft, J., Paul, W. and Valiant, L., *On time versus space*, J. Assoc. Comput. Mach. 24 (1977) 332–337
- [51] Hopcroft, J.E. and Ullman, J.D., *Introduction to Automata Theory, Languages, and Computation*, Addison-Wesley (1979)
- [52] Ibarra, O.H., and Moran, S., *Some time-space trade-off results concerning single-tape and offline TM's*, SIAM J. Comput. 12 (1983) 388–394
- [53] Jacobs, C.T.M. and van Emde Boas, P., *Two results on tables*, Inf. Proc. Letters 22 (1986) 43–48
- [54] Jerrum, M., *The complexity of finding minimum-length generator sequences (extended abstract)*, Proc. Int. colloquium on automata languages and programming 11, Springer Lecture notes in computer science 172 (1984), 270–280
- [55] Jerrum, M., *The complexity of finding minimum-length generator sequences*, Theor. Comp.Sci. 36 (1985) 265–290
- [56] Johnson, D.S., *A child's garden of complexity classes*, this handbook
- [57] Karp, R.M. and Lipton, R.J., *Some connections between nonuniform and uniform complexity classes*, Proc. ACM Symposium theory of computing 12 (1980) 302–309
- [58] Katajainen, J., Penttonen, M. and van Leeuwen, J., *Fast simulation of Turing machines by random access machines*, SIAM J. Comput. 17 (1988) 77–88
- [59] Kleene, S.C., *Origins of recursive function theory*, Ann. Hist. Comp. 3 (1981) 52–67
- [60] Kolmogorov, A.N. and Uspenskii, V.A., *On the definition of an algorithm*, Uspehi Mat. Nauk, 13 (1958) 3–28 ; AMS Transl. 2nd ser. 29 (1963) 217–245
- [61] Kosaraju, S.R., *Real-time simulation of concatenable double-ended queues* , Proc. ACM Symposium theory of computing 11 (1979) 346–351
- [62] Kosaraju, S.R. and Atallah, M.J., *Optimal simulations between mesh-connected arrays of processors*, J. Assoc. Comput. Mach. 35 (1988) 635–650
- [63] Krapchenko, V.M., *Complexity of the realization of a linear function in the class of P-circuits*J. Math. Notes Acad. Sci. USSR (1971) 21–23
- [64] Leong, B.L., and Seiferas, J.I., *New real-time simulations of multihead tape units*, J. Assoc. Comput. Mach. 28 (1981) 166–180
- [65] Lewis, H.R. and Papadimitriou, C.H., *Elements of the Theory of Computation* , Prentice Hall,1981

- [66] Li, M., Longpré, L. and Vitányi, P.M.B., *The power of the queue*, Proc. Structure in Compl. Theory 1 (1986), Springer Lecture notes in computer science 223 (1986) 219–233
- [67] Li, M. and Vitányi, P.M.B., *Tape versus queue and stacs: the lower bounds*, Inf. and Comp. 78 (1988) 56–85
- [68] Li, M. and Yesha, Y., *Separation and lower bounds for ROM and nondeterministic models of computation*, Inf. and Comp. 73 (1987) 102–128
- [69] Lichtenstein, D. and M. Sipser, *Go is polynomial-space hard*, J. Assoc. Comput. Mach. 27 (1980) 393–401
- [70] Lorys, K. and Liskiewicz, M., *Two applications of Fürer's counter to one-tape non-deterministic TMs*, in Chytil, M.P., Janiga, L. and Koubek, V. (eds.), proc. Mathematical foundations of computer science'88, Springer Lecture notes in computer science 324 (1988) 445–453
- [71] Luginbuhl, D.R. and Loui, M.C., *Hierarchies and space measures for pointer machines*, Report UILU-ENG-88-2245, Univ. of Ill. at Urbana-Champaign, Aug 1988
- [72] Lupanov, O.B., *On the asymptotic bounds of complexities of formulas which realize logic algebra functions*, Dokl. Akad. Nauk SSSR 128 (1959) 464–467; Engl. transl. in Autom. Expr. 2 nr. 6 (1960) 12–14
- [73] Maass, W., *Combinatorial lower bound arguments for deterministic and nondeterministic Turing machines*, Trans. Amer. Math. Soc. 292 (1985) 675–693
- [74] Maass, W. and Schnitger, G., *An optimal lower bound for Turing machines with one work tape and a two-way input tape*, Proc. Structure in complexity theory 1, Springer Lecture notes in computer science 223 (1986) 249–264
- [75] Machtey, M. and P. Young, *An Introduction to the General Theory of Algorithms*, Theory of Computation Series, North Holland, 1978
- [76] McColl, W.F. and Patterson, M.S., *The depth of all Boolean functions*, SIAM J. Comput 6 (1977) 373–380
- [77] Mead, C. and Conway, L., *Introduction to VLSI Systems*, Addison Wesley, 1980
- [78] Mehlhorn, K., *On the program size of perfect universal hash functions*, Proc. IEEE Foundations of computer science 23 (1982) 170–175
- [79] Mehlhorn, K., *Data Structures and Algorithms 1, Sorting and Searching*, EATCS Monographs on Theoretical Computer Science 1, Springer Verlag, 1984
- [80] Meyer auf der Heide, F., *Efficiency of universal parallel computers*, Acta Informatica 19 (1983) 269–296
- [81] Minsky, M., *Computation, Finite and Infinite Machines*, Prentice Hall, 1972
- [82] Neciporuk, È.I., *A Boolean function*, Soviet Math. Dokl. 7 (1966) 999–1000

- [83] Parberry, I., *A Complexity Theory for Parallel Computation*, Ph.D. Thesis, Dept. Computer Science, Univ. of Warwick, May 1984
- [84] Parberry, I., *Parallel speedup of sequential machines: a defense of the parallel computation thesis*, SIGACT News 18 nr.1 (1986) 54–67
- [85] Parberry, I. and Schnitger, G., *Parallel computation with threshold functions (preliminary version)*, Proc. Structure in Complexity Theory 1, Springer Lecture notes in computer science 223 (1986), 272–290
- [86] Paterson, M.S., *Tape bounds for time bounded Turing machines*, J. Comput. Syst. Sci. 6 (1972) 116–124
- [87] Paul, W.J., *On-line simulation of  $k + 1$  tapes by  $k$  tapes requires nonlinear time*, Inf. and Control 53 (1982) 1–8
- [88] Paul, W.J., Pippenger, N., Szemerédi, E and Trotter, W.T., *On determinism versus non-determinism and related problems*, Proc. IEEE Foundations of computer science 24 (1983) 429–438
- [89] Paul, W.J. and Reischuk, R., *On time versus space II*, J. Comput. Syst. Sci. 23 (1981) 108–126
- [90] Pippenger, N., *Probabilistic simulations*, Proc. ACM Symposium theory of computing 14 (1982) 17–26
- [91] Pippenger, N., *On simultaneous resource bounds*, Proc. IEEE Foundations of computer science 20 (1979) 307–311
- [92] Pippenger, N. and Fischer, M.J., *Relations among complexity measures*, J. Assoc. Comput. Mach. 26 (1979) 361–381
- [93] Pratt, V.R., *The effect of basis on size of Boolean expressions*, Proc. IEEE Foundations of computer science 16 (1975), 119–121; Also: Kibern. Sb. Nov. Ser. 17 (1980) 114–123 ( in Russian)
- [94] Pratt, V.R. and Stockmeyer, L.J., *A characterization of the power of vector machines*, J. Comput. Syst. Sci. 12 (1976) 198–221
- [95] Razborov, A.A., *Lower bounds for the monotone complexity of some Boolean functions*, Soviet Mat. Dokl. 31 (1985) 354–357
- [96] Reif, J.H., *Complexity of the mover's problem and generalizations, extended abstract*, Proc. IEEE Foundations of computer science 20 (1979) 421–427
- [97] Reisch, S., *HEX ist PSPACE-vollständig*, Acta Informatica 15 (1981) 167–191
- [98] Robson, J.M.,  *$N$  by  $N$  checkers is exptime complete*, SIAM J. Comput. 13 (1984) 252–267

- [99] Robson, J.M., *The complexity of go*, in R.E.A. Mason, (ed.), Information Processing '83, Proc. Ninth IFIP world computer congress, Paris 1983, North Holland, Amsterdam (1983) 413–418
- [100] Ruzzo, W.L., *An improved characterization of the power of vector machines*, Report Univ. Washington, DCS-TR-78-10-01
- [101] Ruzzo, W.L., *Tree-size bounded alternation*, J. Comput. Syst. Sci. 21 (1980) 218–235
- [102] Savage, J.E., *The Complexity of Computing*, John Wiley, 1976
- [103] Savelsberg, M. and van Emde Boas, P., *BOUNDED TILING: an alternative to SATISFIABILITY?*, in Wechsung, G., (ed.), Proc. 2nd Frege Conference, Akademie Verlag, Berlin German Democratic Republic, (1984) 354–365
- [104] Savitch, W.J., *Relations between deterministic and nondeterministic tape complexities*, J. Comput. Syst. Sci. 4 (1970) 177–192
- [105] Savitch, W.J., *The influence of the machine model on computational complexity*, in Lenstra, J.K., Rinnooy Kan, A.H.G. and van Emde Boas, P., Interfaces between computer science and operations research, Math. Centre Tracts 99 (1978) 1–32
- [106] Savitch, W.J., *Recursive Turing machines*, Inter. J. Comput. Math. 6 (1977) 3–31
- [107] Savitch, W.J., *Parallel random access machines with powerful instruction sets*, Math. Systems Theory 15 (1982) 191–210
- [108] Savitch, W.J., and Stimson, M.J., *Time bounded random access machines with parallel processing*, J. Assoc. Comput. Mach. 26 (1979) 103–118
- [109] Schaefer, T.J., *Complexity of some two-person perfect-information games*, J. Comput. Systems Sci. 16 (1978) 185–225
- [110] Schmidt, J.P. and Siegel, A., *The spatial complexity of oblivious k-probe hash functions*, manuscript Courant Institute, New York, May 1988
- [111] Schnitger, G., *Storage Modification Machines versus Kolmogorov - Uspenskii Machines (an information flow analysis)*, manuscript Penn State University, Oct. 1987
- [112] Schnorr, C.P., *The network complexity and the Turing machine complexity of finite functions*, Acta Informatica 7 (1976) 95–107
- [113] Schönhage, A., *On the power of random access machines*, Proc. Int. colloquium on automata languages and programming 6, Springer Lecture notes in computer science 71 (1979) 520–529
- [114] Schönhage, A., *Storage modification machines*, SIAM J. Comput. 9 (1980) 490–508
- [115] Schönhage, A., *Tapes versus pointers, a study in implementing fast algorithms*, EATCS Bulletin 30 (1986) 23–32

- [116] Schönhage, A., *A nonlinear lower bound for Random-Access Machines under logarithmic cost*, J. Assoc. Comput. Mach. 35 (1988) 748–754
- [117] Schorr, A., *Physical parallel devices are not much faster than sequential ones*, Inf. Proc. Letters 17 (1983) 103–106
- [118] Seiferas, J. and Vitányi, P.M.B., *Counting is easy*, J. Assoc. Comput. Mach. 35 (1988) 985–1000
- [119] Shepherdson, J.C. and Sturgis, H.E., *Computability of recursive functions*, J. Assoc. Comput. Mach. 10 (1963) 217–255
- [120] Sipser, N., *Halting space-bounded computations*, Theor. Comp. Sci. 10 (1980) 335–337
- [121] Slot, C. and van Emde Boas, P., *On tape versus core; an application of space efficient hash functions to the invariance of space*, Proc. ACM Symposium theory of computing 16 (1984) 391–400
- [122] Slot, C. and van Emde Boas, P., *The problem of space invariance for sequential machines*, Inf. and Comp. 77 (1988) 93–122
- [123] Stegwee, R.A., Torenvliet, L. and van Emde Boas, P., *The power of your editor*, Report IBM Research, RJ 4711 (50179) (May 1985)
- [124] Stockmeyer, *The complexity of decision problems in automata theory and logic*, Report MAC-TR-133, MIT 1984
- [125] Stockmeyer, L., *The polynomial time hierarchy*, Theor. Comp. Sci. 3 (1977) 1–22
- [126] Stockmeyer, L., *Classifying the computational complexity of problems*, J. Symb. Logic 52 (1987) 1–43
- [127] Stoss, H.J., *Zwei-Band-Simulation von Turingmaschinen*, Computing 7 (1971) 222–235
- [128] Toffoli, T. and Margolus, N., *Cellular Automata Machines, a New Environment for Modeling*, the MIT Press, 1987
- [129] Torenvliet, L. and van Emde Boas, P., *A note on time and space*, Proc. Computing Science in the Netherlands CSN87 (1987) 225–234
- [130] Turing, A.M., *On computable numbers, with an application to the entscheidungsproblem*, Proc. London Math. Soc. ser. 2, 42 (1936) 230–265
- [131] van Emde Boas, P., *The second machine class, models of parallelism*, in Lenstra, J.K. and van Leeuwen, J., (eds.), *Parallel computers and computations*, CWI syllabus 9 (1985) pp. 133–161
- [132] van Emde Boas, P., *The second machine class 2, an encyclopedic view on the parallel computation thesis*, in H. Rasiowa, (ed.), *Mathematical problems in computation theory*, Banach Center Publications, 21 (1987) pp. 235–256

- [133] van Emde Boas, P., *Space measures for storage modification machines*, Report FVI-UvA-87-16, Nov 1987 (to appear in Inf. Proc. Letters).
- [134] van Leeuwen, J. and Wiedermann, J., *Array processing machines*, BIT 27 (1987) 25–43, ( a preliminary version appeared in L. Budach (ed.), *Fundamentals of Computation Theory 1985*, Cottbus, German Democratic Republic, Springer Lecture notes in computer science 199 (1985), 257–268)
- [135] Vitányi, P.M.B., *An optimal simulation of counter machines*, SIAM J. Comput. 14 (1985) 1–33
- [136] Vitányi, P.M.B., *Non-sequential computation and laws of nature*, Proc. Aegean workshop on computing '86, Springer Springer Lecture notes in computer science 227 (1986), 108–120
- [137] Vitányi, P.M.B., *Locality, communication, and interconnect length in multicomputers*, SIAM J. Comput. 17 (1988) 659–672
- [138] Wagner, K. and Wechsung, G., *Computational Complexity*, Mathematische Monographien 19J, VEB Deutscher Verlag der Wissenschaften, 1986, (Published in the West by Reidel 1986)
- [139] Wiedermann, J., *Deterministic and nondeterministic simulation of the RAM by the Turing machine*, in R.E.A. Mason, (ed.), *Information Processing '83*, proc. Ninth IFIP world computer congress, Paris, North Holland, Amsterdam 1983, 163–168
- [140] Wiedermann, J., *Parallel Turing machines* , Report Rijksuniversiteit Utrecht, RUU-CS-84-11, Nov. 1984
- [141] Wiedermann, J., *Fast simulation of nondeterministic Turing machines with application to the knapsack problem*, Report VUSEI-AR 4/1986, Bratislava, Oct. 1986

# Index

- accept 7
- accepting computation tree 47
- accepting configuration 7
- accumulator 24
- alphabet 5
- alternating Turing machine 45
- arbitrary write 52
- array processing machine 54
- asynchronous 8
- basis 38
- Batcher's sort 55
- blank symbol 6
- Blum axioms 25
- bounded nondeterminism 6
- BOUNDED TILING 59
- capacity 21
- cellular automata 41
- center 33
- Church's thesis 3
- circuit 38
- common write 52
- complete basis 38
- complete problem 15
- complexity class 10
- computation 6
- concatenation 24
- configuration 6
- configuration graph 42
- constant factor space overhead simulation 12
- constant factor speed-up 21
- constant-delay simulation 12
- Cook's theorem 15
- counter 19
- crossing-sequence 22
- $\Delta$ -structure 33
- depth 38
- determinism 6
- divergent computation 7
- division 24
- EDITRAM 49
- exclusive write 52
- existential state 46
- final configuration 7
- finite control 5
- first machine class 5
- formula 38
- full computation 7
- fundamental complexity classes 11
- fundamental complexity hierarchy 11
- hard problem 15
- head 18
- history method 21
- indirect address 24
- initial configuration 7
- input 6
- input tape 17
- instruction counter 24
- invariance thesis 4
- Kolmogorov complexity 22
- k-PRAM 55
- KUM 33
- length of input 9
- liberal interpretation 4
- linear-time simulation 12
- local memory 7
- logarithmic measure 25
- logarithmic time measure 9
- logspace reduction 14
- LPRAM 59
- machine 5
- machine model 5
- many-one reduction 14
- master reduction 15
- matrix squaring algorithm 43
- memory 5

memory register 24  
MIMD-RAM 58  
monotone network 38  
MRAM 48  
multiplication 24  
network 38  
nondeterminism 6  
NP-completeness 15  
oblivious 20  
off-line computation 6  
on-line computation 6  
orthodox interpretation 4  
output 6  
output tape 17  
palindrome 22  
paradox of assignment 34  
parallel bitwise operation 24  
parallel computation thesis 4  
parallel machine 7  
parallel Turing machine 57  
parsimonious reduction 15  
pebble game 22  
perfect hashing 31  
pointer machine 34  
polynomial-time reduction 14  
polynomial-time simulation 12  
P-RAM 58  
priority write 52  
program 5  
PSIMDAG 62  
PSPACE-completeness 15  
PTM 57  
QBF 43  
QUANTIFIED BOOLEAN FORMULAS  
    43  
queue 19  
RAM 24  
random access machine 24  
RASP 25  
real-time simulation 12  
recognize 7  
recursive path splitting 43  
recursive Turing machine 55  
reduction 14  
register machine 23  
register machine 25  
rejecting configuration 7  
relative computability 14  
SATISFIABILITY 15  
second machine class 5  
semi-infinite tape 19  
sequential machine 7  
shared memory 7  
SIMDAG 53  
simulation 8  
size 38  
SMM 33  
space 9  
space bounded accepting 10  
space bounded recognition 10  
stack 19  
storage modification machine 33  
structural complexity theory 4  
synchronous 8  
systolic computation 41  
tape cell 18  
terminating computation 7  
time 9  
time bounded accepting 10  
time bounded recognition 10  
transition 6  
transitive closure 41  
Turing machine 17  
unbounded fan-in network 38  
uniform measure 25  
uniform time measure 9  
uniformity 37  
universal state 46  
unused register 29  
vector machine 48  
work tape 17  
write conflict 52



# The ITLI Prepublication Series

## 1986

- 86-01 The Institute of Language, Logic and Information  
86-02 Peter van Emde Boas A Semantical Model for Integration and Modularization of Rules  
86-03 Johan van Benthem Categorical Grammar and Lambda Calculus  
86-04 Reinhard Muskens A Relational Formulation of the Theory of Types  
86-05 Kenneth A. Bowen, Dick de Jongh Some Complete Logics for Branched Time, Part I  
Well-founded Time, Forward looking Operators  
86-06 Johan van Benthem Logical Syntax

## 1987

- 87-01 Jeroen Groenendijk, Martin Stokhof Type shifting Rules and the Semantics of Interrogatives  
87-02 Renate Bartsch Frame Representations and Discourse Representations  
87-03 Jan Willem Klop, Roel de Vrijer Unique Normal Forms for Lambda Calculus with Surjective Pairing  
87-04 Johan van Benthem Polyadic quantifiers  
87-05 Víctor Sánchez Valencia Traditional Logicians and de Morgan's Example  
87-06 Eleonore Oversteegen Temporal Adverbials in the Two Track Theory of Time  
87-07 Johan van Benthem Categorical Grammar and Type Theory  
87-08 Renate Bartsch The Construction of Properties under Perspectives  
87-09 Herman Hendriks Type Change in Semantics:  
The Scope of Quantification and Coordination

## 1988

### *Logic, Semantics and Philosophy of Language:*

- LP-88-01 Michiel van Lambalgen Algorithmic Information Theory  
LP-88-02 Yde Venema Expressiveness and Completeness of an Interval Tense Logic  
LP-88-03 Year Report 1987  
LP-88-04 Reinhard Muskens Going partial in Montague Grammar  
LP-88-05 Johan van Benthem Logical Constants across Varying Types  
LP-88-06 Johan van Benthem Semantic Parallels in Natural Language and Computation  
LP-88-07 Renate Bartsch Tenses, Aspects, and their Scopes in Discourse  
LP-88-08 Jeroen Groenendijk, Martin Stokhof Context and Information in Dynamic Semantics  
LP-88-09 Theo M.V. Janssen A mathematical model for the CAT framework of Eurotra  
LP-88-10 Anneke Kleppe A Blissymbolics Translation Program

### *Mathematical Logic and Foundations:*

- ML-88-01 Jaap van Oosten Lifschitz' Realizability  
ML-88-02 M.D.G. Swaen The Arithmetical Fragment of Martin Löf's Type Theories with weak  $\Sigma$ -elimination  
ML-88-03 Dick de Jongh, Frank Veltman Provability Logics for Relative Interpretability  
ML-88-04 A.S. Troelstra On the Early History of Intuitionistic Logic  
ML-88-05 A.S. Troelstra Remarks on Intuitionism and the Philosophy of Mathematics

### *Computation and Complexity Theory:*

- CT-88-01 Ming Li, Paul M.B. Vitanyi Two Decades of Applied Kolmogorov Complexity  
CT-88-02 Michiel H.M. Smid General Lower Bounds for the Partitioning of Range Trees  
CT-88-03 Michiel H.M. Smid, Mark H. Overmars, Leen Torenvliet, Peter van Emde Boas Maintaining Multiple Representations of Dynamic Data Structures  
CT-88-04 Dick de Jongh, Lex Hendriks, Gerard R. Renardel de Lavalette Computations in Fragments of Intuitionistic Propositional Logic  
CT-88-05 Peter van Emde Boas Machine Models and Simulations (revised version)  
CT-88-06 Michiel H.M. Smid A Data Structure for the Union-find Problem having good Single-Operation Complexity  
CT-88-07 Johan van Benthem Time, Logic and Computation  
CT-88-08 Michiel H.M. Smid, Mark H. Overmars, Leen Torenvliet, Peter van Emde Boas Multiple Representations of Dynamic Data Structures  
CT-88-09 Theo M.V. Janssen Towards a Universal Parsing Algorithm for Functional Grammar  
CT-88-10 Edith Spaan, Leen Torenvliet, Peter van Emde Boas Nondeterminism, Fairness and a Fundamental Analogy  
CT-88-11 Sieger van Denneheuvel, Peter van Emde Boas Towards implementing RL

### *Other prepublications:*

- X-88-01 Marc Jumelet On Solovay's Completeness Theorem

## 1989

### *Logic, Semantics and Philosophy of Language:*

- LP-89-01 Johan van Benthem The Fine-Structure of Categorical Semantics

### *Computation and Complexity Theory:*

- CT-89-01 Michiel H.M. Smid Dynamic Deferred Data Structures  
CT-89-02 Peter van Emde Boas Machine Models and Simulations


### *Other prepublications:*

- X-89-01 Marianne Kalsbeek An Orey Sentence for Predicative Arithmetic

**HET GROTE SPOED B.V.P.**

Faculteit Wiskunde & Informatica	OPDRACHTBON VOOR DRUKWERK, OFFSET EN KOPIEERWERK
Datum	
Opdrachtgever <u>De J. de Jongh</u>	Omschrijving <u>CT-89-02</u>
Vakgroep <u>Logica en theor. Informatica</u>	
Code vakgroep <u>30</u> (22=verkoop syllabi)	
Aantal originelen <u>81</u>	Aantal masters _____
Aantal afdrucken per origineel	
Aantal sets <u>180 K</u>	Bandbreedte _____ mm.
<del>Enkel</del> zijdig/dubbelzijdig	Rugomslag ja/nee
Vergaren ja/nee	Kaft dubbel ja/nee
Niet ja/nee	Kaft enkel ja/nee
Vouwen ja/nee	Kaftkleuren <u>(z</u> uurblauw, karmijnrood, steenrood, goudgeel
Offset	
Stadsdrukkerij	
Fotokopiëren offset	
Fotokopiëren portier	
Fotokopiëren Océ	
Fotokopiëren RX	
Fotokopiëren Canon	

Handtekening budgetbeheerder



Let op: laatste bladzijde op binnen-zijde v/d blauwe achterkaft