

# Mathematics

for Future Computing  
and Communications



EDITED BY  
Liao Heng and Bill McColl

# CAMBRIDGE UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

One Liberty Plaza, 20th Floor, New York, NY 10006, USA

477 Williamstown Road, Port Melbourne, VIC 3207, Australia

314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi – 110025, India

103 Penang Road, #05–06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9781316513583](http://www.cambridge.org/9781316513583)

DOI: 10.1017/9781009070218

© Cambridge University Press 2022

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2022

Printed in the United Kingdom by TJ Books Limited, Padstow Cornwall

*A catalogue record for this publication is available from the British Library.*

ISBN 978-1-316-51358-3 Hardback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

# Contents

*Preface*

page ix

<b>Part I Computing</b>	<b>1</b>
<i>Introduction to Part I</i>	<b>3</b>
<b>1 Mathematics, Models and Architectures</b>	
<i>Bill McColl</i>	<b>6</b>
1.1 Introduction	<b>6</b>
1.2 Moving Beyond von Neumann	<b>8</b>
1.3 Programming-Oriented Models	<b>10</b>
1.4 An Algorithm-Oriented Model	<b>15</b>
1.5 A Bridging Model	<b>17</b>
1.6 Parallel Algorithms and Complexity	<b>25</b>
1.7 Networks and Communications	<b>32</b>
1.8 Resilience-Oriented Models	<b>35</b>
1.9 New Research Directions	<b>48</b>
<b>2 Mathematics and Software Verification</b>	
<i>Chen Haibo and Gao Xin</i>	<b>54</b>
2.1 Introduction	<b>54</b>
2.2 Basic Theories of Formal Methods	<b>55</b>
2.3 Spectrum of Formal Methods	<b>60</b>
2.4 Applications of Formal Methods	<b>61</b>
2.5 Challenges of Formal Verification in Software Systems	<b>66</b>
2.6 Towards Well-Engineered Formal Verification	<b>67</b>
2.7 Conclusion	<b>70</b>
<b>3 Mathematics for Quantum Computing</b>	
<i>Kong Yunchuan</i>	<b>74</b>
3.1 Introduction	<b>75</b>
3.2 Quantum Algorithms	<b>78</b>
3.3 Quantum Error Correction	<b>83</b>
3.4 Quantum Control	<b>87</b>

<b>4</b>	<b><u>Mathematics for AI: Categories, Toposes, Types</u></b>	
	<i><u>Daniel Bennequin and Jean-Claude Belfiore</u></i>	98
	<u>4.1 Introduction</u>	99
	<u>4.2 History</u>	100
	<u>4.3 Categories, Topos, Types and Stacks</u>	103
	<u>4.4 Topos of Deep Neural Networks</u>	115
	<u>4.5 Information Theories</u>	122
	<u>4.6 Higher Categories and Homotopy Types</u>	124
	<u>4.7 Categories and Toposes in Computer Science</u>	126
	<b>Part II Communications</b>	133
	<i><u>Introduction to Part II</u></i>	135
<b>5</b>	<b><u>Mathematics and Compressed Sensing</u></b>	
	<i><u>Zhang Rui and Long Zichao</u></i>	138
	<u>5.1 Introduction</u>	138
	<u>5.2 Sampling Theory and Data Recovery</u>	139
	<u>5.3 Main Theory and Breakthroughs</u>	140
	<u>5.4 Algorithms</u>	145
	<u>5.5 General Compressed Sensing</u>	147
	<u>5.6 Applications in Industry</u>	149
	<u>5.7 Open Questions</u>	150
<b>6</b>	<b><u>Mathematics, Information Theory, and Statistical Physics</u></b>	
	<i><u>Mérouane Debbah</u></i>	153
	<u>6.1 Mathematics of Propagation: Maximum Entropy Principle</u>	153
	<u>6.2 Mathematics of Matrices: Statistical Physics</u>	165
	<u>6.3 Mathematics of Communications: Information Theory</u>	175
	<u>6.4 Conclusion</u>	183
<b>7</b>	<b><u>Mathematics of Data Networking</u></b>	
	<i><u>Li Zongpeng, Miao Lihua and Tang Siyu</u></i>	187
	<u>7.1 Introduction</u>	187
	<u>7.2 System Capacity Region</u>	188
	<u>7.3 Theory and Algorithms of Network Optimization</u>	188
	<u>7.4 The Theory of Network Coding</u>	194
	<u>7.5 Mathematics for Internet Quality of Service (QoS)</u>	200
	<u>7.6 Conclusion</u>	206
<b>8</b>	<b><u>Mathematics and Network Science</u></b>	
	<i><u>Sun Jie</u></i>	211
	<u>8.1 Introduction</u>	212
	<u>8.2 Characterizations of Real Networks</u>	213
	<u>8.3 Structural Models of Complex Networks</u>	215



8.4	Community Detection and Network Partition	218
8.5	Network Dynamics: Synchronization, Control, and Optimization	220
8.6	<a href="#">Data-Driven Analysis: Causal Inference, Automated Modeling</a>	<a href="#">222</a>
8.7	<a href="#">Conclusion</a>	<a href="#">223</a>
<b>Part III</b>	<b>Artificial Intelligence</b>	<b>225</b>
	<i>Introduction to Part III</i>	<a href="#">227</a>
<b>9</b>	<b><a href="#">Mathematics, Information and Learning</a></b>	
	<i><a href="#">Tong Wen and Ge Yiqun</a></i>	<a href="#">230</a>
	<a href="#">9.1 Introduction</a>	<a href="#">230</a>
	<a href="#">9.2 Definition of Information</a>	<a href="#">231</a>
	<a href="#">9.3 Neural Network Information</a>	<a href="#">251</a>
	9.4 Learnability	276
	<a href="#">9.5 Conclusion</a>	<a href="#">283</a>
<b>10</b>	<b><a href="#">Mathematics and Bayesian Inference</a></b>	
	<i><a href="#">Guo Kaiyang, Lv Wenlong and Zhang Jianfeng</a></i>	<a href="#">285</a>
	<a href="#">10.1 Introduction</a>	<a href="#">285</a>
	<a href="#">10.2 Bayesian Inference</a>	<a href="#">287</a>
	<a href="#">10.3 Exact Inference in Bayesian Linear Regression</a>	<a href="#">290</a>
	<a href="#">10.4 Approximate Inference</a>	<a href="#">293</a>
	<a href="#">10.5 Distributed Inference</a>	<a href="#">300</a>
	<a href="#">10.6 Bayesian Optimization</a>	<a href="#">301</a>
	<a href="#">10.7 Bayesian Transfer Learning</a>	<a href="#">304</a>
	10.8 Designing a Prior	305
	10.9 Duality between Control and Inference	306
<b>11</b>	<b><a href="#">Mathematics, Optimization and Machine Learning</a></b>	
	<i><a href="#">Jui Shang-Ling</a></i>	<a href="#">309</a>
	<a href="#">11.1 Introduction</a>	<a href="#">309</a>
	<a href="#">11.2 Stochastic Convex Optimization</a>	<a href="#">311</a>
	<a href="#">11.3 Direct Methods for Non-Convex Optimization</a>	<a href="#">315</a>
	<a href="#">11.4 Optimization for Deep Learning</a>	<a href="#">319</a>
	<a href="#">11.5 Open Problems</a>	<a href="#">324</a>
<b>12</b>	<b><a href="#">Mathematics of Reinforcement Learning</a></b>	
	<i><a href="#">Wu Shuang and Wang Jun</a></i>	<a href="#">329</a>
	12.1 Introduction	329
	12.2 Bayesian Decision Principle	330
	12.3 Markov Decision Process	330
	12.4 Algorithmic Development	340
	<a href="#">12.5 Theoretical Foundations</a>	<a href="#">355</a>
	12.6 Challenges	366

<b>Part IV</b>	<b>Future</b>	<b>375</b>
<b>13</b>	<b>Mathematics and Prospects for Future Breakthroughs</b>	
	<i><b>Dang Wenshuan</b></i>	<b>377</b>
13.1	Future AI: From Perception to Cognition	377
13.2	Future Discovery: From Digital Twin to Quantum Twin	378
13.3	Future Unified Computing Architectures	379
13.4	Future Wireless Systems	380
13.5	Future IP Networks	381
13.6	Future Optical Technologies	381
13.7	Future Autonomous Driving Networks	382
13.8	Future Mathematics: The Analytical Approach	383
	<b>Editors and Contributing Authors</b>	<b>384</b>

# Preface

LIAO Heng

The vitality of the computing and communications industry is remarkable. Accomplished experts occasionally fall into the trap of thinking that all major inventions in their field have been made, and that what remains to be discovered are mere refinements to known methods. Advances in the field of computing and communications have proven experts wrong time and time again. Current notable examples include 5G communication networks and AI computing technologies. Such major technological advances are deeply rooted in mathematical science.

How can we ensure the industry preserves the vitality necessary to generate such advances? With this sizeable challenge in mind, we assembled mathematical scientists and engineers from the computing and communications industry to contribute to this volume. The authors are devoted to providing answers to the question presented above. The chapters are intended to provide insights from a variety of different perspectives.

Mathematics is a beautifully self-coherent body of interconnected concepts, but it was not developed in isolation – it has been accompanied all the way by natural science, especially physics. Bertrand Russell once said, “Physics is mathematical not because we know so much about the physical world, but because we know so little; it is only its mathematical properties that we can discover”. Mathematics has been regarded as the universal language for natural science, and likewise, physics has been described as a rich source of inspiration and insight in mathematics. Cross-disciplinary work has led to some of the greatest discoveries of all time. For example, Newton’s pursuit of classical mechanics resulted in the invention of calculus. David Hilbert, best known for setting much of the agenda for twentieth-century mathematics with his famous 23 problems, defined the differential equations of gravity that gave mathematical formulation to Einstein’s theory of general relativity. Eugene Wigner went so far as to describe the intimacy between mathematics and physics as “a miracle”, and his experiences consistently proved him right. Leading research institutes around the world have achieved great success in both mathematics and physics by promoting intimate exchanges across the two fields of study.

In contrast, the relationship between modern computing, communications and mathematics has been slightly less direct. The history can be traced back to Turing, von Neumann and Shannon, the three mathematicians who have come to be seen as the founders of computing and communications. Their work followed a familiar pattern: first, being drawn to a practical challenge (such as sending information across a noisy channel, or performing complex computational calculations) with a particular set of tools available (such as a transistor capable of processing bits); second, formulating the challenge as a mathematical problem; and last, developing mathematical solutions to prove that the theory addresses practical difficulties. In more recent times, Hinton’s work on backpropagation generated new excitement in the study of AI and neural network processing quickly became a major focus of the computing industry. Can the drive to obtain

knowledge from massive amounts of data, to simulate complex phenomena accurately, dealing with intrinsic uncertainty, and communicating at the semantic level – rather than at the bit level – become the inspiration for a new generation of mathematicians? We believe this volume will persuade many others to build fruitful relationships between computing, communications and mathematics.

Many industry visionaries have long realized that the key to business success is to convert scientific methods into technologies, subsequently applying them to product design through the process of research and development (R&D). The bulk of research is often conducted well before the R&D process. Such a realization has led to confusion among leaders about how to support research in mathematical science. We have produced this volume to assuage this confusion by addressing the major challenges we face in the industry, and in doing so, to open up mathematical problems to more fruitful cross-disciplinary discussion. This will help strengthen the confidence and commitment from industry leaders to provide sustained support for mathematical studies, and to direct resources in the most effective directions.

Researchers and developers in the communications and computing industries are mostly trained in highly specialized domains. They often lack an up-to-date knowledge and awareness of mathematical science beyond their own niche. As a result, many opportunities for applying new mathematical methods to solve problems in engineering are lost, simply because they are developed in other fields. Many engineers also lack the training to be able to convert engineering problems into mathematical models, and so must seek help from mathematicians. This volume aims to serve as a map for engineers, to help them navigate the boundary between engineering and mathematics.

Many areas of mathematical science are highly practical. Although some mathematicians primarily focus on proving theorems, others create and apply models to solve real-life problems. It is not uncommon for mathematicians to underestimate the impact of their work. Equally, many mathematicians are not fully aware of the key mathematical problems in any given applied domain. One major purpose of this volume, therefore, is to highlight the main mathematical problems currently being tackled within the computing and communications industries, and to encourage more mathematicians to direct their efforts towards solving them. It is hoped that interdisciplinary research can be promoted to maximize both its academic impact and the benefits it brings for society.

It must be stressed that the role of advancing fundamental research in computing and communications cannot be undertaken by one organization alone. The volume is intended as a beacon to generate a new wave of excitement among the international research community. Ideally, it will motivate policymakers and university executives to fund research and build education programs with a clearer purpose. We hope it will inspire a new generation of young mathematicians to join this grand effort.

The chapter authors are drawn from a very broad group of researchers, from academics working on core areas of mathematics to experts who have made outstanding contributions to the foundation of modern communication networks and advanced computing devices. We greatly appreciate and sincerely thank the contributors for their capacity to envision a new era of mathematical science that will pave the way for the creation of new machines that can perceive, learn, communicate, think and create.

# Introduction to Part I

The field of Theoretical Computer Science (TCS) covers many areas of computer science and modern mathematics. The topics studied in TCS include algorithms, data structures, computational complexity, combinatorial search and optimization, parallel and distributed computation, probabilistic computation, cryptography, program semantics and verification, as well as computational aspects of logic, geometry, number theory, and algebra.

New research directions in classical areas such as automata theory and information theory are also part of TCS, as are various areas of AI such as machine learning, computational learning theory, theorem proving, and constraint programming. In recent years, whole new areas have been added to TCS such as quantum computing, computational biology, computational economics and algorithmic game theory. Today, the field is still rapidly expanding as new algorithmic or computational theories of natural and artificial phenomena emerge. For example, recent advances in neuroscience are opening up exciting new directions for research in developing theories of cognitive computation, of brain functionality, and of consciousness.

The field of TCS has many major mathematical open questions at its core. Perhaps the most famous is the  $P =? NP$  problem in the area of computational complexity, which is widely regarded as one of the most important open problems in mathematics today. In this theme we introduce just a few of the areas of TCS that have a major significance for innovation in computing and communications technologies.

## Mathematics, Models and Architectures

In the first chapter the relationships between mathematics, models of computation, and the design, analysis and optimization of hardware and software architectures are explored. Mathematics and models of computation have been at the heart of computer science since the earliest research on computing by Turing, von Neumann and other pioneers. Today, models guide how we design and analyze algorithms, how we design and compare architectures, and how we design software that can be automatically adapted to run efficiently on different architectures. For the past 50 years a major goal of computer science research has been to develop a universal “post-von-Neumann” parallel model for algorithms, software and architectures. Some of the history of attempts to produce

such a new universal parallel model are described. also Mathematics continues to be a central element in all of this research.

## **Mathematics and Software Verification**

Mathematical proofs and computer programs can be modeled as the same types of mathematical objects. This makes it possible to use mathematical theories and logics to prove properties of computer programs. In the second chapter, this area of theoretical computer science is explored. Key mathematical concepts such as the lambda calculus, type systems, Hoare logic, and weakest preconditions are introduced, and software techniques such as program analysis, model checking, and automated and interactive theorem proving, are described. Recent research results are given on specifying and verifying operating system kernels, on security-critical components and protocols, and on correctness of synchronization primitives under weak memory models. The chapter closes with a discussion of some of the practical challenges of applying formal methods in industry, and of the benefits that can be obtained from formal methods and formal verification.

## **Mathematics for Quantum Computing**

The concept of computation as a mathematically precise notion was formalized in the 1930s by Emil Post, Alonzo Church, and perhaps most prominently, by Alan Turing. Today, more than 80 years later, it remains the case that what we regard as computable is precisely that which corresponds to the Church–Turing Thesis. A problem is computable, or computationally solvable, if and only if it can be solved on a Turing machine. Our analysis of the computational complexity of problems also assumes that our machines correspond to these classical models. In recent years, the concept of quantum computing has emerged as an interesting new alternative model of computation. Research on the theory and practice of quantum computing is proceeding on many fronts. The third chapter explores this area. Due to the radically different way in which information is represented and manipulated in quantum systems, there is the potential to search for solutions in a space exponentially larger than its classical counterparts, with the same number of bits. With this potential in mind, quantum algorithms have been developed in fields ranging from number theory and algebra, simulation, optimization, and machine learning. However, as the chapter describes, many significant theoretical and practical problems remain to be solved in order to fully achieve this potential.

## **Mathematics for AI : Categories, Toposes, Types**

The fourth chapter reviews possible connections between theoretical computer science and artificial intelligence. These connections are essentially based on notions from

contemporary mathematics, and are heavily influenced by the revolutionary and visionary ideas of Alexander Grothendieck. The generalizing capability of the notions Grothendieck developed (all of them widely using the category language) encompasses computer science and machine learning. After a reminder about category theory, the history of this contemporary mathematical field is described: Grothendieck's discovery of toposes as a category equivalent to the category of sheaves of sets on a site, the existence, in each topos, of a subobject classifier which shows the deep connection with logic.

A description of neural networks as a Grothendieck site is then presented with an explanation of which Grothendieck topology is required to ensure functionality. Many characteristics of neural networks (for example, weights, training datasets and *a priori* knowledge) can then be seen as objects in the topos of sheaves corresponding to the base site of the neural network architecture. We can even take into account the invariance to some symmetries (like CNNs for the group of spatial translations) by generalizing toposes to stacks (fibered in groupoids).

Type intuitionist languages can be interpreted using toposes. In this case, a type is an object of the topos. Martin-Löf's dependent type theory constituted a significant step forward. A further advancement was the construction of new high-level functional programming, which provided effective methods of computation with a large mathematical scope. In turn, this facilitated proof assistance by computers. This path has led us today to programming with homotopy type theory. The definition of infinite-categories/infinite-toposes has paved the way for notions that give all known models of the theories of types of Martin-Löf and the univalent theories of Voevodsky.

# 1 Mathematics, Models and Architectures

---

Bill MCCOLL

## Overview

Models of computation are at the heart of computer science. The notion of what it means to be computable can be precisely and mathematically defined in terms of the Turing machine model, first defined by Alan Turing in 1936. The architecture of sequential computers is based on the von Neumann model, first proposed by John von Neumann in 1945. Turing and von Neumann are today regarded as the founders of our computing industry, but they are also widely recognized as two of the greatest mathematicians of the twentieth century.

Models and mathematics play fundamental roles in computing. They guide how we design and analyze algorithms, how we design and compare architectures, and how we design software that can be automatically adapted to run efficiently on different architectures. For the past 50 years it has been recognized that a post-von-Neumann model will need to be a parallel model, in order to guide the design of parallel algorithms, software and architectures.

In this chapter we will explain some of the history of attempts to produce a new universal parallel model that could potentially supersede and replace the sequential von-Neumann model. We will also show how mathematics continues to be a central element in all of this research. At the end of the chapter we will suggest a number of new areas for fundamental mathematical research that can help guide and influence future work on parallel algorithms, software and architectures.

## 1.1 Introduction

In the nineteenth century and early twentieth century, “computation” was something done by people who were performing tasks such as counting and analyzing census data, or analyzing and predicting the motion of planets.

The concept of computation as a mathematically precise notion was formalized by several researchers in the 1930s – Emil Post, Alonzo Church, and perhaps most prominently, by Alan Turing. Their formalizations, or models, varied significantly – Post’s Machine, Church’s Lambda Calculus (Church 1941), and Turing Machines (Turing 1937). However, it was soon realized that all of these models defined essentially the same notion

<sup>a</sup> From *Mathematics of Future Computing and Communications*, edited by Liao Heng and Bill McColl © 2022 Cambridge University Press.



of what it means to be “computable”. Today, it remains the case that what we regard as computable is precisely that which corresponds to the Church–Turing thesis. We regard a problem to be computable, or computationally solvable, if and only if it can be solved on a Turing machine.

Besides defining the model, Turing also showed other remarkable mathematical results and insights obtained from using the model. For example, he showed that a single “universal” Turing machine could be designed that could simulate an arbitrary Turing machine on arbitrary input. This indicated that it was possible, in principle, to design and build “general-purpose” computers. This mathematical result was the key insight that would later provide the foundation for the launch of the computing industry and its massive global growth over the past 70 years.

Another spectacular mathematical result published by Turing in 1936 showed that many natural problems were not computable. For example, he showed that it was not possible to write a general program that would take a description of any other program and its input and correctly determine whether or not the program would halt on that input or continue running forever.

The halting problem was one of the first problems to be shown to be non-computable, or undecidable. Also in 1936, Church independently showed that a problem in the lambda calculus was undecidable. Since then, using these mathematically precise models of computation, many other natural and important problems have been shown to be undecidable.

These negative mathematical results, and their methods of proof, have massive ramifications for our computing industry. For example, using essentially the same mathematical methods as those used for the halting problem, we can show that it is not possible to write software that will tell if two given programs will always produce the same results. In fact, an even stronger result is Rice’s theorem from 1951 which says essentially that any non-trivial property of a Turing-complete formalism is undecidable. Informally, what this means is that if you have a programming model or computational model  $M$  that is as general and powerful as a Turing machine, then no non-trivial questions about the behavior of programs of  $M$  can be answered computationally!

Without a precise mathematical model of computation, none of these major theoretical achievements would have been possible. Today, theory and models remain at the heart of modern computer science as we shall see. They guide how we design and analyze algorithms, how we design and compare architectures, and how we design software that can be automatically adapted to run efficiently on different architectures.

Following Turing’s groundbreaking theoretical work in the 1930s on universal Turing machines, John von Neumann in 1945 proposed a practical design for a computer architecture (von Neumann 1945). Von Neumann’s model consisted of:

- Processing unit with arithmetic logic unit and registers.
- Control unit with instruction register and program counter.
- Memory for data and instructions.
- External mass storage.
- Input and output mechanism.

## 1.3 Programming-Oriented Models

The von Neumann model is an example of an architectural model that is neutral in terms of the particular styles of program that can be efficiently run on it. A von Neumann sequential machine is effective at handling not only simple numerical and scientific computations, but also complex data structures and all of the many types of applications that are required across the global computing industry. In contrast, a number of the early examples of models proposed as post-von-Neumann took a more programming-oriented approach – essentially first define the programming model and then try to design an architectural model to support it. In many cases, the resulting architectures have proved to be unable to deliver many of the other characteristics required of a model. We will later see that an alternative approach, that of defining a “bridging model” between the many styles of programming and the many types of architecture, has offered a more convenient way of achieving the many simultaneous goals. In this section we consider some of these programming-oriented models.

### 1.3.1 Dataflow

Around 1974, Jack Dennis (1974) and others suggested that programs should not be tied to a specific sequential temporal execution based on a program counter, as in the von Neumann model. Instead, they argued, the dataflow in a computation was the fundamental aspect, and operations should instead be scheduled and executed as soon as their required input data was available. In this dataflow model, instead of von Neumann assignment to variables, computations would be represented as graphs of single-assignment operations. Such graphs could be static or dynamic.

The dataflow or task parallelism approach has been tried many times over the past 40 years, most recently in the TensorFlow model for machine learning and AI computation. At modest scale the dataflow model can certainly be useful. However, the optimization of dataflow computations at large scale remains a major challenge. Also, as it is a relatively low-level programming model, the productivity of programmers is low, especially when trying to achieve high performance at scale with dataflow models.

As it has been proposed quite recently, it is perhaps important to elaborate further on the TensorFlow model (Dean et al. 2015), developed by Google. TensorFlow computations are expressed as stateful dataflow graphs consisting mainly of operations on tensors (multidimensional data arrays). TensorFlow is not intended to be a model for general purpose parallel computing, or even a model for a broad class of AI computations. It is only intended to support a class of tensor-based machine learning applications. However, these computations (particularly training) are often very computationally intensive, communication-intensive, and highly iterative. They therefore require a model that can offer high performance at scale.

Viewing TensorFlow as a parallel computing model, it is essentially the same as the dataflow models and frameworks that were proposed originally in the 1970s. As a special purpose computing model, TensorFlow provides much of what is needed to build deep learning applications on small-scale parallel systems. As the scale increases and there

is a need to carry out the computation on large-scale distributed memory architectures, the standard problems of data distribution and communication minimization become the dominant challenges, as they always do.

In the case of deep learning at small scale, batch splitting for data parallelism can be used. However, at larger scale, this approach suffers from all of the standard problems in parallel computing: memory constraints, high latency, and inefficiency due to fine-grain small batch sizes. To address these issues we need to shift to more of a “model parallelism” approach, using the kind of standard parallel algorithm, software and architectural techniques that have been used in BSP, MPI and other frameworks for over 25 years, such as high performance point-to-point communications and optimized collective communications such as AllReduce implemented efficiently at scale on powerful networks. Today this trend is already underway in systems such as Mesh-TensorFlow (Shazeer et al. 2018) that are starting to use some basic high performance computing concepts.

In summary, TensorFlow provides a convenient special purpose framework for a certain class of tensor dataflow computations. Such special purpose frameworks are not intended to support a broad class of computations. As such they are like scripting languages that can be useful in certain areas, but do not replace more general purpose programming languages. In terms of scalability and performance, it is not clear whether such a dataflow approach can offer convenience in terms of simplicity and software productivity, while at the same time delivering high performance at scale.

### 1.3.2 Functional Programming

In 1977, John Backus, who had designed Fortran in the 1950s, won the Turing Award. His Turing Award lecture (Backus 1978) was entitled “Can Programming Be Liberated from the von Neumann Style? A Functional Style and Its Algebra of Programs” In it he argued that the world should move on from the temporal imperative programming style of von Neumann and Fortran to a new era in which programming would be much more mathematical and abstract, one where computation would be functional.

Functional models of computation were not new, even in the 1970s. They have been around for a very long time. As noted earlier, the lambda calculus, a functional model of computation, was first introduced by Alonzo Church (1941), around the same time as Turing (1937) introduced the Turing Machine model. Haskell Curry’s work on combinatory logic (Curry & Feys 1958) was another early example of a functional model of computation, closely related to the work of Church. Although even here, Curry was not the first. The first work on combinatory logic was done by Schönfinkel in 1920 in Göttingen.

Not only are functional models of computation an old idea, so are functional programming languages! In 1958, just one year after the development of Fortran, John McCarthy and his colleagues developed the functional programming language Lisp which is based on the lambda calculus. Remarkably, even today, Fortran and Lisp are still both in widespread use around the world, more than 60 years after they were first developed. Since then many other functional programming languages have been developed, includ-

ing Scheme, ML and Haskell. Today, many popular programming languages such as Python, Lua and Scala have incorporated important elements of functional programming.

Functional models of computation offer significant potential in terms of improving software productivity, as they provide a flexible high-level abstraction that makes programming easier using concepts such as first class and higher-order functions. The main challenge in this area is to provide the necessary software automation to ensure that the required load balancing, communication-efficiency, synchronization-efficiency, etc., can be achieved without the programmer having to specify how it should be done at scale on complex heterogeneous parallel architectures. To date, this remains an unsolved problem.

### 1.3.3 Data Parallelism

Data parallelism is perhaps the simplest model of parallel computation, where the control flow is sequential, as in the von Neumann model, but the data is distributed so that in each step the processors can operate independently on separate parts of the data, with no need for coordination, synchronization, communication or any of the other aspects that make parallel programming challenging.

SIMD machines, vector machines, GPUs and other architectures support various styles of data parallelism. Data parallelism is attractive in terms of programmer productivity, as it is a very simple programming model for simply structured parallel computations, however it lacks the flexibility and power to handle a broad class of computations in an efficient manner. Besides low-level GPU libraries, three recent examples of programming models that are predominantly data parallel are MapReduce, Spark and GraphBLAS. We will say a little about each of these.

#### MapReduce

MapReduce (Dean & Ghemawat 2004) is one of the simplest models of computation ever proposed. It takes two of the most basic functions map and reduce and builds a whole system platform around that basic functionality. Originally developed by Google for simple big data computations, it gained prominence for a brief period in the form of Apache Hadoop. Today, Google and others have moved on beyond the limitations of this primitive model.

Perhaps the most worthwhile element of this minimalist model and platform was that it provided automatic support for redundancy and fault tolerance on cheap commodity clusters. However, this functionality was only achieved at a huge cost in terms of efficiency, since the computations on these MapReduce and Hadoop platforms often involved rescheduling and re-executing large-scale tasks and achieved data resilience using slow and inefficient distributed file systems.

For a very small class of parallel computations with very simple and light communication requirements, MapReduce could be used, but for any computations that required any kind of iteration, the MapReduce model would be hopelessly inefficient due to the requirement to write to, and read from, distributed persistent storage in order to

communicate. Since most large-scale computations today such as graph computing, analytics and machine learning involve many rounds of iterative computation, this model is no longer relevant. Moreover, many modern parallel computations are also highly communication-intensive. Since most MapReduce implementations write communications to distributed persistent storage for fault recovery, this also severely limits the applicability of these MapReduce platforms.

### Spark

Recognizing the severe limitations of the MapReduce model for any iterative parallel computation, the Spark model (Zaharia et al. 2010) and platform were developed to address some of these shortcomings. The Spark model supports a class of data parallel computations and, like MapReduce, provides a degree of fault tolerance.

Spark's core idea is the Resilient Distributed Dataset (RDD) which can be viewed as a working set for certain types of parallel computation. Spark provides a form of distributed shared memory. The main benefit is that data management and communication is handled automatically, but of course this is also the main limitation, in that the algorithm designer or programmer has insufficient control over data distribution and communication to achieve high performance in many cases. Like MapReduce, Spark requires a cluster manager and a distributed storage system to handle the various automated elements.

The core programming model for Spark is functional programming on RDDs, using for example the Scala programming language. While this provides a nice high level programming abstraction, the Spark model also has some imperative programming features such as accumulators, and some shared variable mechanisms such as broadcast variables.

The Spark model was an important development in the area of models of parallel computation, as it provided a simple data parallel programming abstraction for a range of big data computations in which one is applying the same operation to all elements of a dataset. However, for large-scale communication-intensive parallel computations, the convenience of automatic management of memory and communications in a model such as Spark comes with a significant price in terms of performance.

### GraphBLAS

The Basic Linear Algebra Subprograms (BLAS) have, for many years, provided a convenient set of algorithmic software tools for high performance computing. The objective of work on GraphBLAS (Buluç et al. 2017) is to provide a similar framework of building blocks for graph computing. The starting motivation for this work is the simple observation that any weighted graph can be represented by its adjacency matrix, and with such a representation, we can easily transform linear algebra on matrices and vectors into graph computations such as breadth first search by simply changing the scalar functions or operators. For example, changing the standard operator pair  $(*,+)$  into the semiring  $(+,\min)$  we can compute various functions on paths in graphs.

GraphBLAS provides a simple high level functional programming model for computations on sparse and dense graphs, with the potential to improve software productivity. Programs operate directly on vectors and matrices, both of which can be dense or

sparse matrices. Implementations of GraphBLAS can take advantage of this sparsity to achieve high performance. They can also exploit the freedom that is provided by such a high level programming model, in order to optimize data distribution and the associated parallel communication and synchronization in distributed memory architectures. The various standard higher order functions are provided as data-centric primitives or building blocks:

- Matrix-vector product.
- Matrix-matrix product.
- Inner product.
- Data parallel elementwise addition and multiplication.
- Fold left and fold right.

It should be noted that the last two correspond closely to map and reduce, and can be generalized in that direction. Similarly, the basic mechanisms of BSP data-centric “Think Like a Vertex” graph frameworks such as Pregel, which we will discuss later, can also be modeled using GraphBLAS primitives. So, the GraphBLAS model may be extensible in ways that offer more flexibility than one would expect at first sight.

### 1.3.4 Message Passing

In 1977, Tony Hoare introduced the communicating sequential processes (CSP) model of concurrent computation (Hoare 1978). The CSP model is based on the fundamental idea that synchronization and communication should be tightly coupled. CSP computations use synchronous communication via channels for sending and receiving. Other models of concurrent computation were introduced around the same time, such as Robin Milner’s Calculus of Communicating Systems (CCS) and other process algebras.

In the 1980s the message passing model was used in a number of early parallel computer architectures, including the Inmos Transputer with its associated programming language Occam, which was heavily influenced by CSP. This message passing approach came to be more generally referred to as “distributed memory parallel computation” in contrast to the kind of shared memory models that we will consider next.

In the 1990s there was a rapid proliferation of distributed memory message passing architectures and programming tools. In response to this, a committee was formed – the MPI Forum – with the goal of developing a standard Message Passing Interface (MPI) for point-to-point and collective communications in a distributed memory parallel architecture. The first version of the MPI standard was designed in 1993 consisting of a library of functions usable with standard sequential programming languages such as C. Since then a series of further updated versions have been defined, and some of them have been implemented.

Hoare’s original CSP model can be seen as addressing the need for a sound mathematical framework for low-level communications. As such, it succeeds, in that complex low-level protocols can be analyzed and verified using the CSP model and its associated tools. In the case of the Transputer architecture and the Occam language, this very low-level message passing is fundamental to the hardware architecture. In the case of

accelerators, then the numbers of such cores within a single device may be quite high, but it remains the case that for systems in which the number of chips themselves is large, shared memory is not a viable model.

## 1.5 A Bridging Model

In 1990, Leslie Valiant published the first description of a new model of computation, the Bulk Synchronous Parallel (BSP) model (Valiant 1990a), based on the idea of computing in rounds, or supersteps. In the initial description, a strong emphasis was placed on demonstrating that the idealized PRAM shared memory model could be realized by simulation on certain forms of BSP machine with appropriate properties. This “Automatic-Mode BSP” was based on the ideas of using hashing to randomize the distribution of data, and of using “parallel slackness” to hide or tolerate latency to non-local memory. More importantly, Valiant showed that BSP had a simple cost model that could be used to capture the important costs of communication and synchronization in a parallel computation in a convenient way, enabling powerful analysis, comparison and optimization of parallel algorithms.

These two initial achievements of the BSP superstep model in demonstrating PRAM simulation via automatic memory management, and accurate cost modeling of parallel algorithms, were merely the first steps in a long line of BSP innovations that followed through the 1990s. Most of the subsequent research, rather than focusing on Automatic-Mode BSP, PRAM simulation, strobing barriers etc. instead pursued the alternative direction of BSP as an architectural and/or programming model. In this “Direct-Mode BSP”, data distribution and the associated communications and synchronization would be controlled by the algorithm designer or programmer, with the system only automating the low-level scheduling of those communications to avoid congestion and maximize throughput.

Between 1990 and 1992, Leslie Valiant and I worked on ideas for a distributed memory BSP programming model. Between 1992 and 1997, I led a large research team at Oxford that developed various BSP programming libraries, languages and tools, and also numerous massively parallel BSP algorithms (McColl 1995). With interest and momentum growing, I then led an international group that developed and published the BSPlib Standard for BSP programming (Hill et al. 1998) in 1996. Since then, there have been hundreds of papers on BSP research, and BSP has become a standard model for parallel computing (Skillicorn et al. 1997). It is now widely used in research and in industry in many areas, including data analytics, graph computing, machine learning and HPC.

### 1.5.1 BSP Architectures

A BSP computer consists of a set of processor-memory pairs, a global communications network, and a mechanism for the efficient barrier synchronization of the processors. A BSP computer operates in the following way. A computation consists of a sequence

of parallel supersteps, where each superstep consists of a sequence of steps, followed by a barrier synchronization at which point all data communications will be completed. During a superstep, each processor can perform a number of computation steps on values held locally at the start of the superstep, send and receive a number of messages, and handle various remote read/get and write/put requests.

The BSP computer is a two-level memory model, i.e. each processor has its own physically local memory module; all other memory is non-local, and is accessible in a uniformly efficient way. By uniformly efficient, we mean that the time taken for a processor to read from, or write to, a non-local memory element in another processor-memory pair should be independent of which physical memory module the value is held in. The algorithm designer and the programmer should not be aware of any hierarchical memory organization based on network locality corresponding to the particular structure of the communications network.

BSP computations have both a horizontal (spatial) structure and a vertical (temporal) structure. The horizontal structure arises from concurrency, and consists of a fixed number of virtual threads. These are each associated, at run-time, with a physical processor. The vertical structure arises from the progress of a computation through time. For BSP, this is a sequential composition of global supersteps, which conceptually occupy the full width of the executing architecture. Each superstep is further subdivided into three ordered phases consisting of:

- Computation locally in each thread, using only values stored in the local memory of each processor.
- Communication actions, e.g. put and get, amongst the threads, involving movement of data between processors.
- Barrier synchronization, which waits for all of the communication actions to complete, and which then makes the data that was moved available in the local memories of the destination processors.

A superstep is shown in Figure 1.1.

If the target parallel computer has fewer processors than the virtual parallelism (parallel slackness), then a simple transformation can be used to convert the BSP program into a slimmer version. For example, if number of virtual processes is twenty times larger than the number of physical processors then we can map twenty virtual processes to each processor. Moreover, this can be done in a way that maximizes the number of communications that become internal operations.

### 1.5.2 BSP Cost Modeling

If we define a time step to be the time required for a single local operation, i.e. a basic operation (such as addition or multiplication) on locally held data values, then the performance of any BSP computer can be characterized by three parameters:

- $p$  = Number of processors.



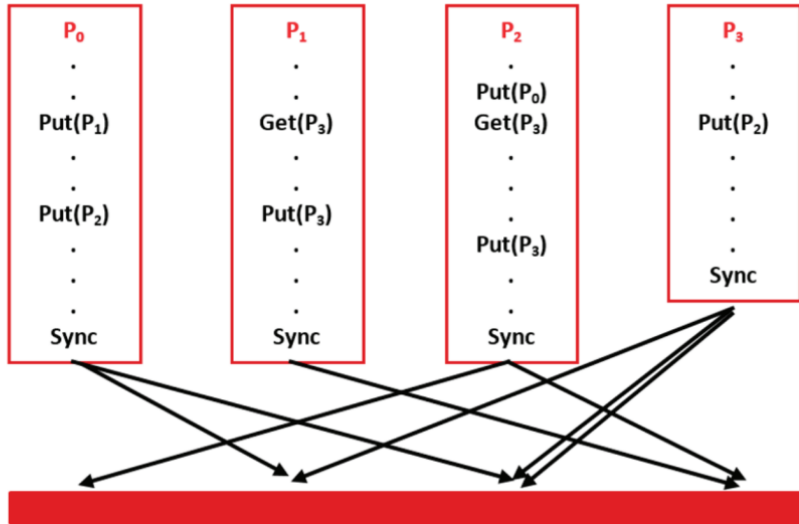


Figure 1.1 Superstep

- $g = (\text{Total number of local operations performed by all processors in one second}) / (\text{Total number of words delivered by the communications network in one second, in a situation of continuous traffic})$ .
- $L = \text{Number of time steps for barrier synchronization}$ .

There is also, of course, a fourth parameter  $s$ , the number of time steps per second. However, since the other parameters are normalized with respect to that one, it can be ignored in the design of algorithms and programs. The parameter  $g$  corresponds to the frequency with which non-local memory accesses can be made; in a machine with a higher value of  $g$  one must make non-local memory accesses less frequently. More formally,  $g$  is related to the time required to realize  $h$ -relations in a situation of continuous message traffic. An  $h$ -relation is a communication pattern in which each processor sends up to  $h$  messages and receives up to  $h$  messages. The parameter  $g$  is the value such that any  $h$ -relation can be performed in  $g * h$  time steps.

Any parallel computing system can be regarded as a BSP computer, and can be benchmarked accordingly to determine its BSP parameters  $L$  and  $g$ . The BSP model is therefore not prescriptive in terms of the physical architectures to which it applies. Every general purpose parallel architecture can be viewed by an algorithm designer or programmer as simply a point  $(p, g, L)$  in the space of all BSP machines.

In theoretical terms, the BSP model can be regarded as a generalization of the PRAM model which permits the frequency of barrier synchronization, and hence the demands on the network, to be controlled. If a BSP architecture has a very small value of  $g$ , for example  $g = 1$ , then it can be regarded as a PRAM and we can use hashing to automatically achieve efficient memory management. The value of  $L$  will determine the

degree of parallel slackness required to achieve optimal efficiency. The case  $g = L = 1$  corresponds to the idealized PRAM, where no parallel slackness is required.

The time for a superstep is determined as follows. Let the work  $w$  be the maximum number of local computation steps executed by any processor during the superstep, and let  $h_s$  be the maximum number of messages sent by any processor, and  $h_r$  be the maximum number of messages received by any processor during the superstep. The time for the superstep is then at most  $w + g * \max\{h_s, h_r\} + L$  steps. The total time required for a BSP computation is easily obtained by adding the times for each superstep. Analyzing and predicting the cost of a BSP program is, therefore, no more difficult than analyzing and predicting the cost of a sequential program.

By adding the time for each of  $S$  supersteps we obtain an expression of the form  $W + g * H + L * S$  where  $W$ ,  $H$ ,  $S$  will typically be functions of  $n$  and  $p$ . In designing an efficient BSP algorithm or program for a problem which can be solved sequentially in time  $T(n)$  our goal will, in general, be to produce an algorithm requiring total time  $W + g * H + L * S$  where  $W(n, p) = T(n)/p$ ,  $H(n, p)$  and  $S(n, p)$  are as small as possible, and the range of values for  $p$  is as large as possible. In many cases, this will require that we carefully arrange the data distribution so as to minimize the frequency of remote memory references. Another property of interest in BSP algorithm design is the space (or memory) efficiency of the computation. We use  $M(n, p)$  to denote the maximum number of values which any one processor has to store at any point during the computation.

### 1.5.3 Why BSP?

#### **Universal. General purpose.**

The essence of the BSP approach to parallel programming is the notion of the superstep, in which communication and synchronization are completely decoupled. A “BSP program” is simply one which proceeds in phases, with the necessary global communications taking place between the phases. This approach to parallel programming is applicable to all kinds of parallel architectures and parallel algorithms. It provides a consistent, and very general, framework within which to develop portable parallel software for scalable parallel architectures.

#### **Simple.**

BSP programs are much the same as sequential programs. Only a bare minimum of extra information needs to be supplied to describe the use of parallelism.

#### **Easy to debug. No deadlock.**

Since communication and synchronization are decoupled in a BSP program, the programmer does not have to worry about problems such as deadlock, which can occur with synchronous message passing. Debugging a BSP program is also made much easier by this decoupling. The barrier at the end of a superstep provides an appropriate breakpoint at which the global state of the parallel computation is well defined and can be interrogated. Debugging and reasoning about the correctness of a BSP program are, therefore, not much more difficult than for a sequential program.

#### **Portable. Independent of target architectures.**

Unlike many parallel programming models, BSP is designed to be architecture-independ-

ent, so that programs run unchanged when they are moved from one architecture to another. Thus BSP programs are portable in a strong sense.

**Performance of a program on a given architecture is predictable.**

The execution time of a BSP program can be computed from the text of the program and a few simple parameters of the target architecture. This makes design space exploration possible, since the effect of a decision on performance can be easily determined.

**Analyzable.**

The BSP cost model makes it clear what strategies should be adopted to write efficient BSP programs. We should balance and minimize the computation between threads, since  $W$  is a *maximum* over computation times. We should balance and minimize the communication between threads, since  $h$  is a *maximum* over fan-in and fan-out of data. We should minimize the number of supersteps, since this determines the number of times  $L$  appears in the final cost. The cost model also shows how it can be used to predict performance on a particular target architecture. The values of  $W$  and  $h$  for each superstep, and the number of supersteps can be determined by inspection of the program code. Values of  $p$ ,  $g$ , and  $L$  can then be inserted to give execution time before the program is executed.

**Easy to checkpoint for resilience.**

The fact that there we can easily capture a well-defined state at each barrier also means that we have a simple way in which to perform checkpointing and recovery from hardware failures. In contrast, capturing and checkpointing the state of a parallel computation in a message passing system such as MPI is incredibly difficult. In a later section, on the BSP cloned computing model, we will see that BSP superstep semantics not only allows us to handle hardware failures, but also to efficiently handle long tail latencies, which is a much more challenging problem.

**Communications can be optimized for global exchange.**

In large-scale parallel computing it is almost always the case that the parallelism in the software vastly exceeds the parallelism in the hardware. This is a simple consequence of the power of modern computing hardware. A single core is capable of performing over ten billion operations per second. In the example mentioned above, where a graph algorithm has one trillion parallel subtasks, these subtasks will be mapped onto a physical machine which has a much, much smaller degree of parallelism. For example, it might be mapped onto 200 16-core machines, in which case each core will be handling around 300 million small parallel subtasks. The same phenomenon of parallel slackness occurs in every area of large-scale parallel computing – sparse matrix computations, machine learning, modeling, optimization etc.

A major consequence of parallel slackness is that it will normally be the case that in any large-scale parallel computation, each processor will be communicating with every other processor as the computation proceeds. So, to achieve peak performance we need to design parallel software systems that are optimized for this global pattern of communication in which each processor is sending to all the others, and is receiving from all the others. This global pattern is called total exchange. BSP software systems are optimized for total exchange, rather than for single, individual pairwise communications, as in MPI

systems have appeared. Today, data-centric BSP has become the standard way to perform large-scale graph computations in databases and other analytics systems.

As noted above, the basic idea behind data-centric BSP systems such as Pregel is to “Think Like a Vertex”. What this means is that the data structure, for example a graph with one billion vertices should be mapped to a virtual parallel computation in which vertices become local computations and edges become communications. As in all BSP systems, the computation proceeds in rounds or supersteps. In each round, each vertex will perform local computation, send messages to neighbors in the graph, and receive messages from neighbors, as shown in Figure 1.2. This iterative computation continues until some given termination condition is met.



Figure 1.2 Data-centric BSP

As in all BSP systems, superstep semantics ensures that we have no concurrency issues such as race conditions or deadlock. The degree of fine-grain virtual or logical parallelism in such a computation can of course be huge, perhaps one billion or one trillion. Practical parallel machines, on the other hand, will typically have only a few hundred or a few thousand cores. So, as noted previously, this massive scale virtual computation can be easily mapped onto a much smaller scale physical parallel machine in a straightforward way, and in a way that takes account of the structure of locality within the graph to optimize communications.

Data-centric BSP has been extensively used in recent years by many companies around the world for practical computations on graphs with billions of vertices and trillions of edges. Besides delivering massive scalability and extremely high performance, the model is also very easy to learn and to use effectively. For example, Google’s PageRank algorithm, which is central to their search methods can be described in just 15-20 lines of code, since the algorithm can be formulated in data-centric BSP as the following simple iterative computation:

- Read neighbor ranks.
- Update local rank.
- Send new rank to neighbor.

The database company TigerGraph recently announced that their data-centric BSP graph

database was now the world's fastest transactional graph database in production, handling huge volumes of ecommerce and financial transactions.

Although we have emphasized graph analytics in discussing the data-centric BSP model, it is also applicable to many other classes of parallel computation. For example, it can be used in scientific and engineering computations such as stencil computations,  $n$ -body problems, finite elements, and circuit modeling. Also in machine learning and other areas of AI.

## 1.6 Parallel Algorithms and Complexity

The simplicity of the BSP cost model enables “cost-driven design” of parallel algorithms, parallel software and parallel architectures. As noted above, in designing an efficient BSP algorithm or program for a problem which can be solved sequentially in time  $T(n)$ , our goal will, in general, be to produce an algorithm requiring total time  $W + g * H + L * S$  where  $W(n, p) = T(n)/p$ ,  $H(n, p)$  and  $S(n, p)$  are as small as possible, and the range of values for  $p$  is as large as possible.

### 1.6.1 BSP Algorithms for Common Parallel Patterns

Many static computations can be conveniently modeled by directed acyclic graphs (DAGs), where each node corresponds to some simple operation, and the arcs correspond to inputs and outputs. These DAGs often have a simple structure or pattern. In this section we will look at some of those patterns. In what follows, in the interests of simplicity and clarity, we will often ignore small constant factors.

#### Tree Computations

An important pattern in parallel computation is the tree pattern. Figure 1.3 shows a ternary tree.

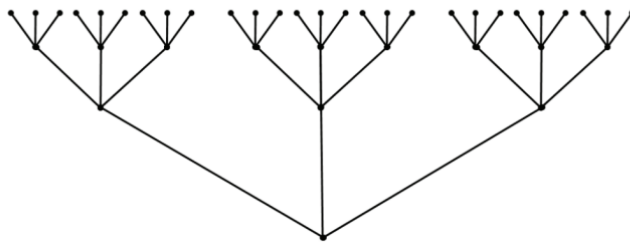


Figure 1.3 Ternary tree pattern

Consider the simple problem of broadcasting a single value from one processor to all other  $p - 1$  processors. We can use a balanced  $k$ -ary tree of height  $s$ , where  $s = \log_k p$ . This corresponds to a BSP algorithm with  $s$  supersteps, each costing  $g * k + L$ . The total cost is therefore  $s * (g * k + L)$ , so if we choose  $k = L/g$ , then the cost will be

$2 * \log_k p * L$ . If we choose  $k = 2$ , then the cost will be  $\log_2 p(g * 2 + L)$ . If we chose  $k = p$ , then we simply broadcast in a single round in which case the cost will be  $g * p + L$ . So, even for a very simple problem such as broadcasting a single value, there are a wide range of choices, and the cost model allows these to be easily analyzed and compared.

Consider now the related problem of broadcasting an array of  $n$  values from one processor to all other  $p - 1$  processors, where  $n$  is much larger than  $p$ . For this problem we can use a different approach. In the first superstep, the source processor sends each of the other processors a distinct sub-array of size  $n/p$ . The cost of this step is  $g * n + L$ . In the second superstep, each processor sends its sub-array to all of the other processors. The cost of this step is again  $g * n + L$ , so the total cost is only  $2 * (g * n + L)$ .

So, given a particular problem of broadcasting an array of  $n$  values on a  $p$  processor machine with BSP parameters  $g$  and  $L$ , for some particular values of  $n, p, g$  and  $L$ , this kind of cost-driven design can be used to obtain the best possible solution. The same kind of analysis can be applied to find optimal algorithms for related tree-structured problems such as Reduce (single tree) and Prefix Sums (double tree).

### Butterfly Computations

Another important pattern in parallel computation is the butterfly pattern. Figure 1.4 shows an 8-point butterfly.

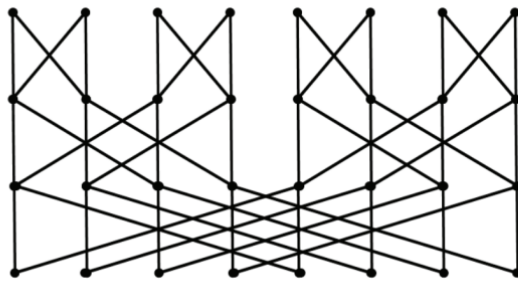


Figure 1.4 8-point butterfly

Consider the problem of computing the Fast Fourier Transform (FFT) of  $n$  points. It is well known that this algorithm has sequential complexity  $O(n \log n)$  and that it can be represented as a  $(\log n)$ -level butterfly graph with a bit reversal permutation applied to the inputs. It can be implemented on a  $p$ -processor BSP machine in  $(\log n)/(\log(n/p))$  supersteps, in each of which each processor sequentially computes the next  $\log(n/p)$  levels of the butterfly graph on its  $n/p$  points in time  $(n/p) * \log(n/p)$ . The cost of each superstep is  $(n/p) * \log(n/p) + g * (n/p) + L$ , and therefore the total time required is at most  $((n \log n)/p) * (1 + g/\log(n/p) + L/((n/p) \log(n/p)))$ . In most practical situations, the number of processors  $p$  will be no more than  $n^{1/2}$ . In such cases, the number of supersteps would be at most two, and the total time required would be  $(n \log n)/p + g * (n/p) + L$ .



### 2D Computations

Another important pattern in parallel computation is the two-dimensional grid DAG. Figure 1.5 shows a  $9 \times 9$  2D grid in which all of the 144 arcs are directed downwards.

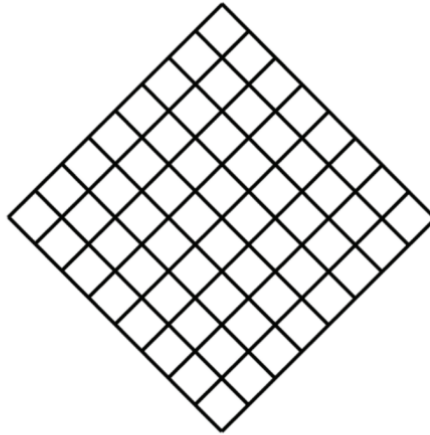


Figure 1.5 2D grid DAG

Let  $D_n$  denote the 2D Grid DAG which has  $n^2$  nodes  $v_{i,j}$ ,  $0 \leq i, j < n$ , and arcs from  $v_{i,j}$  to  $v_{i+1,j}$  and  $v_{i,j+1}$  where those nodes exist. In McColl (1995) it is shown that the graph  $D_n$  can be efficiently scheduled for a  $p$  processor BSP computer by partitioning  $D_n$  into  $p^2$  subgraphs, each of which is isomorphic to  $D_{n/p}$ . The resulting schedule shows that any computation that can be mapped directly onto  $D_n$  has a BSP implementation with cost at most  $n^2/p + g * n + L * p$ . Note that for such a BSP algorithm, in some cases, using more processors will actually *increase* the runtime. For example, consider a BSP architecture based on a simple ring. Such an architecture will have parameters  $g = L = p$ . For such a machine the runtime will be  $n^2/p + n * p + p^2$ . This runtime is minimized when  $p = n^{1/2}$ . Increasing the number of processors beyond this value will increase the runtime.

In McColl (1995) it is shown that the problem of solving a triangular  $n \times n$  linear system can be mapped onto  $D_n$ , so we obtain the above upper bound for that problem. Many other problems can be similarly mapped onto  $D_n$ . Dynamic programming algorithms are often used in combinatorial search problems such as string comparison and computing string edit distance. These dynamic programming algorithms use tabulation to successively compute the entries of a 2D cost matrix. Such computations map directly onto  $D_n$ .

### 3D Computations

The 2D Grid pattern can be easily extended to higher dimensions. Figure 1.6 shows a  $6 \times 6 \times 6$  3D Grid DAG in which all of the 540 arcs are directed downwards.

Let  $C_n$  denote the 3D Grid DAG which has  $n^3$  nodes  $v_{i,j,k}$ ,  $0 \leq i, j, k < n$ , and arcs from  $v_{i,j,k}$  to  $v_{i+1,j,k}$ ,  $v_{i,j+1,k}$  and  $v_{i,j,k+1}$  where those nodes exist. In McColl (1995) it

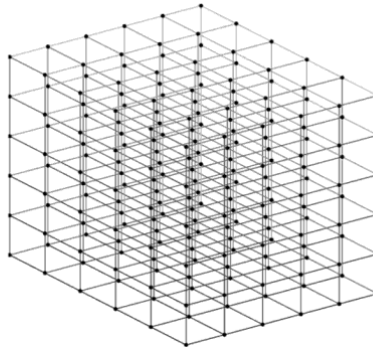


Figure 1.6 3D grid DAG

is shown that the graph  $C_n$  can be efficiently scheduled for a  $p$ -processor BSP computer by partitioning  $C_n$  into  $p^{3/2}$  subgraphs, each of which is isomorphic to  $C_{n/p^{1/2}}$ . The resulting schedule shows that any computation that can be mapped directly onto  $C_n$  has a BSP implementation with cost at most  $n^3/p + g * (n^2/p^{1/2}) + L * p^{1/2}$ .

In McColl (1995) it is shown that the problem of computing the LU decomposition of an  $n \times n$  linear system can be mapped onto  $C_n$ , so we obtain the above upper bound for that problem. Many other problems can be similarly mapped onto  $C_n$  or onto closely related 3D graphs. One such example is the Algebraic Path Problem (APP) over a closed semi-ring, which includes as special cases the problems of computing shortest paths and computing transitive closures. The standard Floyd-Warshall dynamic programming algorithm for the APP can be mapped directly onto a minor variant of  $C_n$  in which some of the arc directions are modified. With this approach we get the above upper bound for all of the various instances of the APP too.

We have explicitly considered 2D and 3D computations. However, the same general scheduling techniques can be easily extended to similar but higher dimensional directed acyclic graphs.

### 1.6.2 Communication-Optimality and Immortal Algorithms

To further illustrate some of the important issues in BSP algorithm design, we will consider several fundamental computational problems involving vectors and matrices.

#### Dense Matrix-Vector Multiplication

Consider the problem of multiplying an  $n \times n$  matrix  $M$  by an  $n$ -element vector  $v$  on  $p$  processors, where  $M, v$  are both dense. In McColl (1995) it is shown that by using a “block-block” or a “block-grid” distribution of the matrix  $M$  we can produce a simple BSP algorithm with cost  $n^2/p + g * (n/p^{1/2}) + L$ .

To compute  $u = M.v$ , the matrix  $M$  and vectors  $u, v$  are partitioned uniformly across the  $p$  processors. In the first superstep, each processor gets all the vector elements  $v_j$  for which it holds a corresponding matrix element  $m_{i,j}$ . In the second superstep, each



$$\begin{aligned} a_{i,j,k} &= a_{i,j-1,k} \\ b_{i,j,k} &= b_{i-1,j,k} \\ c_{i,j,k} &= c_{i,j,k-1} + (a_{i,j,k} * b_{i,j,k}) \end{aligned}$$

where  $a_{i,0,k} = a_{i,k}$ ,  $b_{0,j,k} = b_{k,j}$  and  $c_{i,j,0} = 0$ . These definitions can be directly translated into a labeled version of the 3D directed acyclic graph  $C_n$ . The BSP cost of Method 3 is therefore at most  $n^3/p + g * (n^2/p^{1/2}) + L * p^{1/2}$  which is the same as Method 2. It also has the same optimal memory cost as Method 2.

All of the three methods we have described have the same communication cost  $g * H(n,p) = g * (n^2/p^{1/2})$ . This communication cost can, however, be reduced further if we use the 3D algorithm described in McColl (1995), which Leslie Valiant and I developed for the BSP model, although other similar methods were also developed for other models. The cost of this 3D algorithm (Method 4) is  $n^3/p + g * (n^2/p^{2/3}) + L$ . It is clearly optimal in terms of computation cost and synchronization cost. Moreover, an input-output complexity argument based on the isoperimetric inequality result of Loomis and Whitney (1949) can be used to show that for any BSP implementation of the standard  $n^3$  sequential algorithm, if  $W(n,p) = n^3/p$  then  $H(n,p) \geq n^2/p^{2/3}$ . So Method 4 provides a BSP realization of the standard  $n^3$  matrix multiplication method which simultaneously achieves the optimal values for computation cost  $W(n,p)$ , communication cost  $g * H(n,p)$  and synchronization cost  $L * S(n,p)$ . The memory requirement of this algorithm is, however, slightly inferior to Methods 2 and 3. Its memory complexity  $M(n,p)$  is  $n^2/p^{2/3}$ . In closing we note that there is no single BSP algorithm that is optimal in all four dimensions (computation, communication, synchronization, memory). In Irony et al. (2004) it is shown that if  $M(n,p) = O(n^2/p)$  then  $H(n,p) = \Omega(n^2/p^{1/2})$ . So any optimal memory algorithm must be suboptimal in communication. Depending on whether we wish to optimize memory or communication, we can choose Method 2 or Method 4.

As we have seen, the design, analysis and optimization of parallel algorithms relies heavily on mathematical analysis and theorems such as isoperimetric inequalities. This stems from the fact that viewed through the correct mathematical lens, we see that in the graph structure of parallel computations, the region corresponding to each processor has a surface that corresponds to the communication into and out of that processor, while the volume of the region corresponds to the operations performed by that processor.

### 1.6.3 Recursive Algorithms and Automatic Tradeoffs

For the problem of multiplying two  $n \times n$  matrices  $A$  and  $B$ , we can regard the matrices as each composed of four square  $n/2 \times n/2$  submatrices. In this setting, the  $n \times n$  product matrix  $C = A.B$  can be computed using eight  $n/2 \times n/2$  matrix multiplications and four  $n/2 \times n/2$  matrix additions. These sub-problems can then be further recursively subdivided in the same way. Strassen showed that in such a recursive matrix multiplication algorithm, the number of multiplications can be reduced from 8 to 7 if we allow the number of additions to increase from 4 to 15. The effect of this change is to reduce the asymptotic complexity from  $O(n^3)$  to  $O(n^\alpha)$  where  $\alpha = \log_2 7$ .

Consider the recursive  $O(n^3)$  algorithm. After  $k$  levels of recursion we have  $8^k$  matrix multiplication sub-problems. If we have  $p$  processors then we might stop the recursion

when  $8^k = p$  since the number of sub-problems would match the number of processors. At that stage the size of the square submatrices would be  $n/p^{1/3} \times n/p^{1/3}$ . In McColl & Tiskin (1999) it is shown that the resulting recursive algorithm is essentially Method 4 above. It has optimal computation, communication and synchronization cost.

Now suppose instead that we had continued the recursion beyond that point, and stopped it when  $4^k = p$ . At that stage, the size of the square submatrices would be  $n/p^{1/2} \times n/p^{1/2}$ . In this case, McColl & Tiskin (1999) show that the resulting recursive algorithm is essentially Method 2 above. It can be realized with optimal memory cost.

If we terminate the recursion at some level  $k$  where  $p^{1/3} < 2^k < p^{1/2}$ , then we will have a BSP algorithm that requires less memory than Method 4 and less communication than Method 2. So the results of McColl & Tiskin (1999) provides us with a single unified recursive algorithm for BSP matrix multiplication that we can use to automatically achieve a tradeoff between memory optimality and communication optimality. The results of McColl & Tiskin (1999) also apply to Strassen's sub-cubic matrix multiplication algorithm.

We noted previously that our BSP algorithm for dense matrix-vector multiplication was an "immortal algorithm" in the sense that it simultaneously achieved the optimal computation cost, the optimal communication cost, and the optimal synchronization cost, all to within small constant factors. So, it achieves the best possible parallel performance for all values of  $n$ ,  $p$ ,  $g$  and  $L$ . For matrix multiplication, Method 4 above is also an immortal algorithm in this sense. However, it does require more memory than Method 2. If we redefine the term immortal to also require memory optimality then the result in Irony et al. (2004) shows that there is no such algorithm for matrix multiplication. Given that fact, the recursive algorithm in McColl & Tiskin (1999) is as close to "immortality" as can be achieved in that it can be simply and automatically tuned to deliver optimality in any of the four dimensions.

## 1.7 Networks and Communications

A large amount of work has been done in recent years on the development of efficient routing methods, on the efficient embedding of one network in another, and on the demonstration of work-preserving emulations of one network by another. We will focus our attention here on the routing problem.

We consider the problem of routing  $h$ -relations on a  $p$ -processor network. We are interested in the development of distributed routing methods in which the routing decisions made at a node at some point in time are based only on information concerning the packets that have already passed through the node at that time. In the non-distributed case where global information is available everywhere, the problem of routing is easier and well understood.

### 1.7.1 Oblivious Routing

Let us first consider deterministic methods for distributed routing. We define a routing method to be *oblivious* if the path taken by each packet is entirely determined by its source and destination. It is known (Borodin & Hopcroft 1985) that, for a 1-relation no deterministic oblivious routing method can do better than proportional to  $(p^{1/2}/d)$  time steps, in the worst case, for any degree  $d$  graph. The most obvious examples of deterministic oblivious approaches are greedy methods in which one sends all packets to their destination by a shortest path through the network.

For 1-relations, the performance of greedy routing on a (fixed-degree) butterfly network can be summarized as follows. All 1-relations can be realized in  $O(p^{1/2})$  steps, which, as we have observed, is an optimal worst case bound for any such fixed degree network. A large number of specific 1-relations which arise in practical parallel computation, e.g. the bit-reversal permutation and the transpose permutation, provably require proportional to  $p^{1/2}$  steps.

What about the “average case”? Define a *random 1-mapping* to be the routing problem where each processor has a single packet which is to be sent to a random destination. Greedy routing of a random 1-mapping on a butterfly will terminate in  $O(\log p)$  steps. Moreover, the fraction of all random 1-mappings which do not finish in  $O(\log p)$  steps is incredibly small, despite the fact that most of the 1-relations which seem to arise in practice do not finish in this time. We can probably conclude from these results that “typical” routing problems, in a practical sense, is a rather different concept from “typical” routing problems in a mathematical sense.

The performance of greedy routing on a  $(\log p)$ -degree hypercube is very similar to the case of the butterfly. All 1-relations can be realized in  $O(p^{1/2}/\log p)$  steps, which is an optimal worst case bound for any  $(\log p)$ -degree network. For the average case, where each packet has a random destination, greedy routing will terminate in  $O(\log p)$  steps. In the case of the hypercube, there are exponentially many shortest paths for a greedy method to choose from, but even randomizing among these choices still gives no better than  $O(p^\alpha)$ ,  $\alpha > 0$ , steps for many 1-relations.

### 1.7.2 Randomized Routing

We have seen that for the butterfly and hypercube, the performance of greedy routing on random 1-mappings is much better than on “worst case 1-relations”, such as the bit-reversal permutation in the case of the butterfly. Around 1980, Valiant made the simple and striking observation that one could achieve efficient distributed routing, in terms of worst case performance, if one could reduce a 1-relation to something like the composition of two random 1-mappings. The resulting technique which emerged from this observation has come to be known as two-phase randomized routing (Valiant 1990b).

Using this approach, a 1-relation is realized by initially sending each packet to a random node in the network, using a greedy method. From there it is forwarded to the desired destination, again by a greedy method. Both phases of the routing correspond closely to the realization of a random 1-mapping. Extensive investigation of this method,

in terms of the number of steps required, size of buffers required etc., has shown that it performs extremely well, both in theory and in practice. Using randomized routing one can show that with high probability, every 1-relation can be realized on a  $p$ -processor butterfly, 2D array and hypercube in a number of steps proportional to the diameter of the network. For these fixed-degree networks, this result is essentially optimal. For the  $(\log p)$ -degree hypercube network, the following even stronger result can be obtained. With high probability, every  $(\log p)$ -relation can be realized on a  $p$ -processor hypercube in  $O(\log p)$  steps. Randomized routing can also be used to achieve good worst case performance on other networks such as the shuffle-exchange network, the cube-connected-cycles, and on fat trees.

An interesting theoretical alternative to using randomized routing on a standard, well-defined network such as a butterfly, is to use deterministic routing on a “randomly wired network”. In Leighton & Maggs (1989), Upfal (1992) it is shown that a simple deterministic routing algorithm can be used to realize a 1-relation in  $O(\log p)$  steps on a randomly wired, bounded degree network known as a multi-butterfly. An important feature of multi-butterflies is that they have powerful expansion properties. In addition to permitting fast deterministic routing, such expander graphs also have very strong fault tolerance properties.

### 1.7.3 Networks, Routing and BSP

The use of the parameters  $L$  and  $g$  to characterize the communications performance of a BSP computer contrasts sharply with the way in which communications performance is described for most distributed memory architectures produced and sold today. A major feature of the BSP model is that it lifts considerations of network performance from the local level to the global level. We are thus no longer particularly interested in whether the network is a 2D array, a butterfly or a hypercube, or whether it is implemented in VLSI or in some optical technology. Our interest instead is in the global parameters of the network, such as  $L$  and  $g$ , which describe its ability to support non-local data communications in a uniformly efficient manner.

In the design and implementation of a BSP computer, the values of  $L$  and  $g$  which can be achieved will depend on the capabilities of the available technology and the amount of money that one is willing to spend on the communications network. As the computational performance of machines continues to grow, we will find that to keep  $L$  and  $g$  low it will be necessary to continually increase our investment in the communications hardware as a percentage of the total cost of the machine. In asymptotic terms, the values of  $L$  and  $g$  one might expect for various  $p$ -processor networks are as shown in Table 1.1 (ignoring small constant factors).

These asymptotic mathematical estimates are based on the degree and diameter properties of the corresponding graph, and on the use of a fast routing method such as randomized routing. In a practical setting, the channel capacities, routing methods used, physical implementation etc. would also have a significant impact on the actual values of  $L$  and  $g$  which could be achieved on a given machine. New optical technologies may offer the prospect of further reductions in the values of  $L$  and  $g$  which can be achieved,

**Table 1.1** Network  $L$  and  $g$  values

Network	$L$	$g$
Ring	$p$	$p$
2D Array	$p^{1/2}$	$p^{1/2}$
Butterfly	$\log p$	$\log p$
Hypercube	$\log p$	1
Completely Connected	1	1

by providing a more efficient means of non-local communication than is possible with VLSI.

Choosing the right values of  $L$  and  $g$  for a network architecture will depend on the class of computations that the parallel machine needs to support. If an application has a BSP cost  $W + g * H + L * S$  then a network architecture where  $g$  is less than  $W/H$  and  $L$  is less than  $W/S$  will ensure that the machine is reasonably balanced for that application, and that the utilization of the available computation capacity will be acceptable. Of course, the lower the values of  $L$  and  $g$ , the higher the overall utilization will be, up to a point. However, attempting to achieve the lowest possible  $L$  and  $g$  will not, in general, be worthwhile, as the cost involved will increase dramatically, and may only provide a minimal increase in performance. Simple mathematical analysis, using the BSP cost model, can easily determine all of these tradeoffs in a precise way.

## 1.8 Resilience-Oriented Models

Resilient, predictable architectures are essential for large-scale parallel computation. As scale increases, the number of faults and long tails increases correspondingly. We therefore need a model that can efficiently handle fault tolerance and tail tolerance at scale. In 2017, a new BSP-based bridging model was developed that, for the first time, solves this problem. This new Cloned Computing model supports high-performance parallel computing with fault tolerance and tail tolerance (McColl 2018).

The model is general purpose – it applies to all forms of large-scale parallel computing, including communication-intensive, highly iterative computations. It provides a new, simple model enabling large-scale parallel software to be run with automatic fault and tail tolerance, with no program changes for load balanced computations, and only a simple Boolean parameter addition for non-load-balanced computations. It also provides an accurate cost model enabling large-scale parallel software to be automatically optimized for any architecture. Finally, the model enables high performance nonstop and realtime parallel computation. Unlike checkpointing-based recovery, the new model enables computations to run continuously without interruption and meeting realtime requirements, with high probability. Below we will provide an overview of this new model.

highly iterative, involving thousands of rounds, and can be very naturally and easily expressed as BSP computations.

The Cloned Computing model extends the BSP model in a major way to support high performance fault tolerance and tail tolerance. An implementation of the model runs parallel programs that compute in rounds using a new execution model which for convenience we will call Nonstop BSP. A Nonstop BSP program has four core features:

- It is a BSP program and is structured to compute in rounds. In each round, each of the processes computes on data in local memory, globally communicates across the network, synchronizes.
- The number of processes in a program is the parallelism parameter  $p$ .
- The number of rounds in a program is the parameter  $R$ , which may be finite or infinite.
- Each process has a BSP synchronization mechanism at the end of each round. This mechanism is parameterized by a (typically Boolean) value indicating whether or not the process is one that can set the minimum time value for the round. This parametric value can be changed dynamically during the execution of a program. If a program has the parameter set so that no process can set the minimum time value in any round then the cloned computing system will execute the program as a normal BSP program without fault tolerance or tail tolerance.

In Figure 1.1 we showed a simple example of a non-load-balanced BSP program with four processes. For that program, we might have an expected time per round of 4 seconds for  $p_0, p_1$  and  $p_2$  but only 100ms for  $p_3$ . In order to turn this into a Nonstop BSP program for execution, the only change required is to parameterize the synchronization primitives in order to show whether or not the associated process is one that can set the minimum time value for the round. So, for this example, we might set the Sync parameter to False for  $p_3$ , and to True for the others, as shown in Figure 1.8.

The following provides a simple high level overview of a cloned architecture with  $P$  processors, running a  $p$ -process parallel program, where  $p \leq P$ .

- Each of the  $P$  processing elements can run one or more of the processes during a single round.
- Multiple instances of a process are referred to as clones.
- Each of the  $p$  processes is run as the first process on at least one of the  $P$  elements.
- The first process to be run on an element is the only one that can possibly set MinTime for the round.
- During each round, if the first process to be run on each element has its synchronization mechanism set indicating that it can set MinTime, then it monitors its elapsed time for the current round (Nonstop BSP) and attempts to write MinTime when it ends. The first such process to write its time sets MinTime for round.
- Each process can have access not only to its own local data and state, but also to other information including
  - A copy of MinTime for the round;
  - Its elapsed time for the round;
  - Which clones of other processes it may need to communicate with;

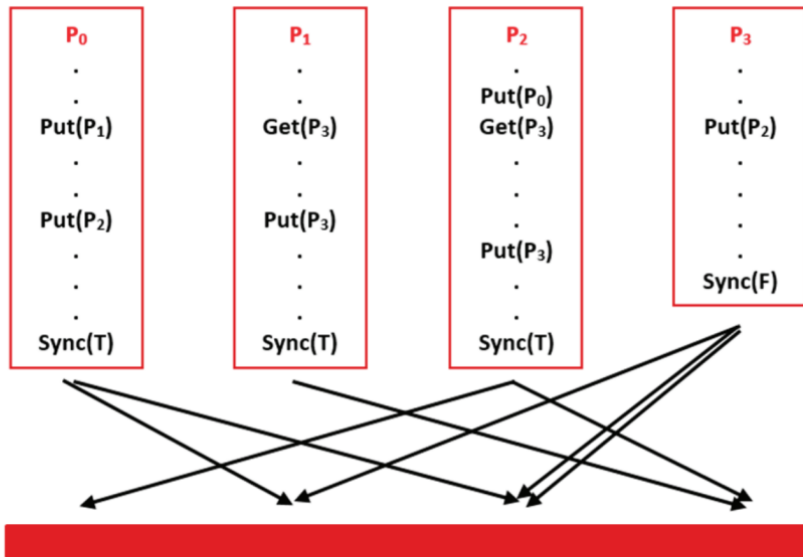


Figure 1.8 Superstep with Sync parameter

- Standby resources available;
- Where other clones of the process are located.
- Cloned computing architectures have a polynomial “Predictability” exponent  $D$ , indicating that an expected fraction  $P/t^D$  of the processing elements will fail to complete any round in less than  $t \cdot \text{MinTime}$ .
- The cloned computation has a TailLimit  $T$ . Any process that fails to complete before  $T \cdot \text{MinTime}$  is marked as a fault/tail, others are marked as live.
- A round is successful if at least one clone of each of the  $p$  processes completes.
- The inter-process communications can be handled in several ways. For example, the first clone of a process to complete can handle the global communications for all the clones of that process. A number of other variations are also possible, depending on other objectives such as balancing communications at endpoints, increasing network latency resilience and other factors.
- Process faults and tails can be handled by transferring state from another live clone of the same process. To improve the speed with which this state transfer and relaunch can be achieved, it may be convenient to maintain a pool of spare containers that are ready to run. This can be done in several ways. For example, by having a static pool of containers directly associated with the various processes, or by having a dynamic pool of containers each of which can be used with any process. The choice between having a static pool, a dynamic pool, or no pool of spare containers can be made based on a tradeoff between speed of relaunch and efficiency of container utilization.



### 1.8.3 Vertical and Horizontal Cloning

To ensure fault tolerance and tail tolerance, processes can be cloned vertically or horizontally, or both. We noted above that in a cloned architecture with  $P$  processing elements running a  $p$ -process parallel program, where  $p \leq P$ .

- Each of the  $P$  processing elements can run one or more of the processes during a single round.
- Multiple instances of a process are referred to as clones.
- Each of the  $p$  processes is run as the first process on at least one of the  $P$  elements.

Let's consider first the case where  $p = P$ , and look at vertical cloning. Suppose  $p = 10$ , and we have processes numbered from 0 to 9 as shown in Figure 1.9.



Figure 1.9 Processes 0 to 9

Let us further assume that the program is perfectly load balanced, so that every process takes exactly the same time as all the others, in every round. Then with no faults or tails we expect that all processes will complete the round at the same time, as shown in Figure 1.10.

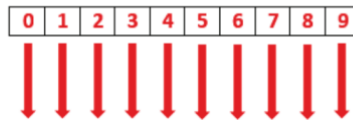


Figure 1.10 Processes 0 to 9 with no faults or tails

If, however, there are faults or tails then the processes may take quite different amounts of time, as shown in Figure 1.11.

To address this challenge we can use vertical cloning, where multiple process instances run consecutively on the same processing element (no extra elements required), as shown in Figure 1.12.

In the example shown, each process appears twice, but vertical cloning can be used with any multiple. The multiples do not even need to be uniform, although this may normally be the case. For example, we might have the cloning shown in Figure 1.13 where each process has three additional clones.

In such uniform cases, we say that the degree of vertical cloning is the parameter  $VC$ . In the example shown in Figure 1.13, we have  $VC = 4$ . The first row of processes are the only ones that can try to set  $\text{MinTime}$  for the round, and can do so only if their synchronization parameter is set accordingly.

Next we introduce the idea of horizontal cloning. Again suppose  $p = 10$ , and we have processes numbered from 0 to 9 as shown in Figure 1.9. As in the case of vertical





Figure 1.11 Processes 0 to 9 with tails

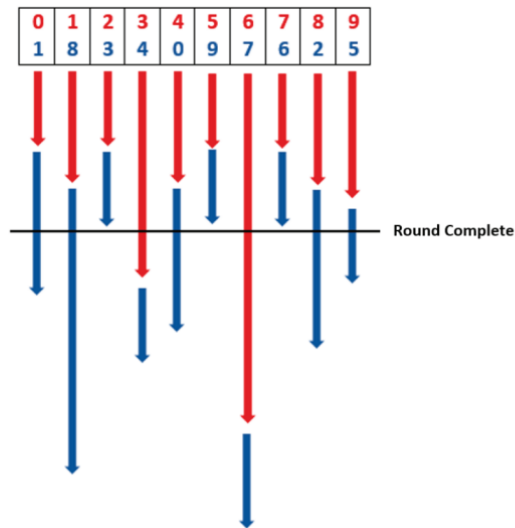


Figure 1.12 Processes 0 to 9 with vertical cloning

0	1	2	3	4	5	6	7	8	9
1	8	3	4	0	9	7	6	2	5
7	5	4	9	1	6	8	0	3	2
2	3	6	1	7	4	0	5	9	8

Figure 1.13 Vertical cloning

cloning, let us further assume that the program is perfectly load balanced, so that every process takes exactly the same time as all the others, in every round. Instead of running multiple process instances on the same element, we can instead run multiple instances

of processes on different elements concurrently, by using more processing elements, as shown in Figure 1.14.

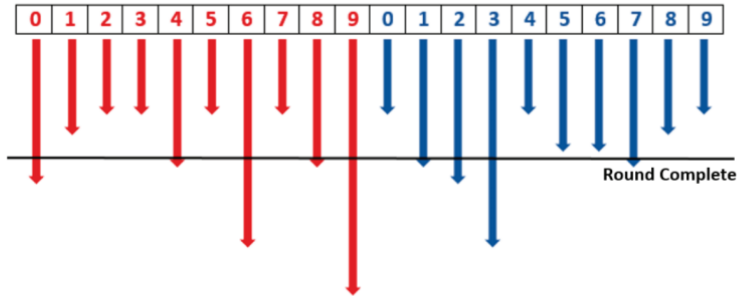


Figure 1.14 Horizontal cloning

In the example shown in Figure 1.14, each process appears twice, but horizontal cloning can be used with any multiple. The multiples do not need to be uniform. For example, for  $p = 6$  and  $P = 21$  we might have the cloning structure shown in Figure 1.15.



Figure 1.15 Horizontal cloning with  $HC=3.5$

We say that the level of horizontal cloning is the parameter  $HC = P/p$ . In the examples shown in Figures 1.14 and 1.15, the values of  $HC$  are 2 and 3.5.

Vertical and horizontal cloning can be combined. For example, we might have  $VC = 3$  and  $HC = 2$ , as shown in Figure 1.16.



Figure 1.16 Horizontal cloning with  $VC=3$  and  $HC=2$

As another example, we might have  $VC = 2$  and  $HC = 1.5$ , as shown in Figure 1.17.



Figure 1.17 Horizontal cloning with  $VC=2$  and  $HC=1.5$

As in the case of simple vertical cloning, the first row of processes are the only ones that can try to set  $MinTime$  for the round, and can do so only if their synchronization parameter is set accordingly.

- Cloned Computing Parameters:
  - Predictability  $D$ ;
  - TailLimit  $T$ ;
  - VerticalCloneLevel  $VC$ ;
  - HorizontalCloneLevel  $HC$ ;
  - StandbyLevel  $S$ .

Given  $D$ ,  $p$ ,  $R$  we can automatically compute the most appropriate values for  $T$ ,  $VC$ ,  $HC$ ,  $S$  to optimize performance and ensure nonstop resilience. The cost is at most  $T * HC * \text{Perfect}$ , where Perfect is the cost for an idealized architecture where  $D$  is infinite (no failures or tails ever). We can place this new cost model at the top level of a hierarchy of cost models for different models of parallel computing, including PRAM (idealized shared memory), MapReduce, and standard BSP with no fault tolerance or tail tolerance.

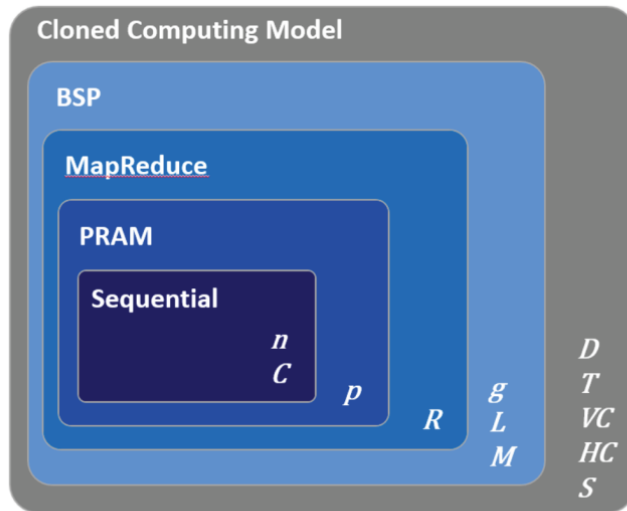


Figure 1.22 Cost model hierarchy

The model parameters in the cost model hierarchy shown in Figure 1.22 are as follows:

- $n$  = Problem Size.
- $C$  = (Sequential) Complexity.
- $p$  = Parallelism.
- $R$  = Rounds.
- $g$  = Network Throughput.
- $L$  = Network Latency.
- $M$  = Local Memory.
- $D$  = Predictability.
- $T$  = TailLimit.

- $VC$  = VerticalClone Level.
- $HC$  = HorizontalClone Level.
- $S$  = StandbyLevel.

For the cloned computing model:

- The parameters  $p$ ,  $R$  and data distribution are decided by the programmer.
- The parameters  $g$ ,  $L$ ,  $M$ ,  $D$  are parameters determined by the infrastructure (hardware+software).
- Using the cost model, optimal  $T$ ,  $VC$ ,  $HC$ ,  $S$  parameters can be automatically calculated to guarantee a given Quality of Service Level (performance and resilience).

So, the model, in addition to providing high-performance general-purpose parallel computing with fault tolerance and tail tolerance, is also easy to use. Any large-scale parallel software can be automatically run as a cloned computation. Just add a Boolean to sync, and make it always true if you have a load balanced program. Also, any large-scale parallel software can be automatically optimized for cloned execution. Given program Parallelism  $p$  and Rounds  $R$ , and Predictability  $D$ , we can automatically generate the optimal values of TailLimit, VerticalClone level, HorizontalClone level, StandbyLevel.

### 1.8.6 Why Cloned Computing?

#### **General purpose.**

Applies to all forms of large-scale parallel computing, including communication-intensive, highly iterative computations. Previous approaches only covered a very small class of MapReduce parallel computations.

#### **Automated.**

Provides a new, simple programming model (Nonstop BSP) enabling large-scale parallel software to be automatically run on the fault tolerant and tail tolerant system with no changes for load balanced computations, and only a simple Boolean parameter addition for non-load-balanced computations. The new model provides an accurate cost model enabling large-scale parallel software to be automatically optimized for any infrastructure.

#### **High Performance. Nonstop. Predictable.**

The new model enables high performance, nonstop, predictable, and realtime parallel computation. Unlike checkpointing-based recovery, it enables computations to run continuously without interruption and meeting realtime requirements, with high probability.

The simple diagrams in Figure 1.23 show the potential impact of this new model in reducing the time taken to perform five rounds of a load balanced parallel computation.

### 1.8.7 Coded Computing

As we have seen, cloned computing provides a solution to the resilience and predictability problem that can be applied to any parallel computation. This goes far beyond previous

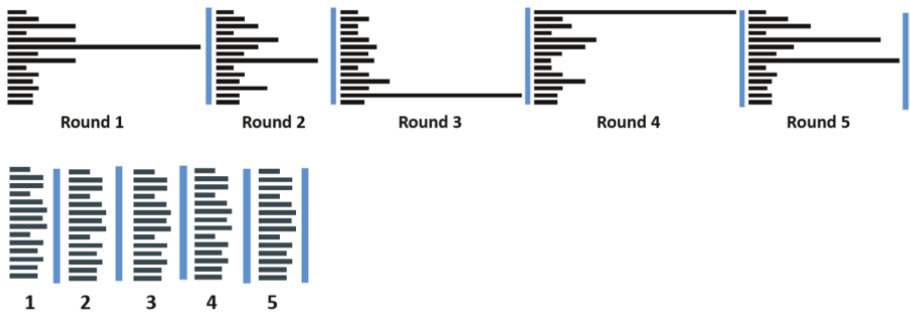


Figure 1.23 Impact of cloned computing

results that only applied to simple tree-structured parallel computations such as MapReduce and other related forms of master-worker distributed computing. Another approach to the resilience and predictability problem that has been explored recently is “Coded Computing” (Yu et al. 2019). This is aimed at the more limited range of tree-structured parallel computations, but targets not only the tail latency (straggler) problem but also issues of communication, security and data privacy. It can be seen as a complementary approach to cloned computing, addressing a narrower scope of computations but a broader set of issues.

As noted in the discussion of cloned computing, the performance of a modern parallel and distributed systems is significantly affected by anomalous system behavior and bottlenecks – a form of “system noise”. Given the individually unpredictable nature of the nodes in these systems, we are faced with the challenge of achieving fast results in the face of massive uncertainty.

One way of tackling this challenge is using coding theoretic techniques. The role of codes in providing resiliency against noise has been studied for decades in several other engineering contexts, and is part of our everyday infrastructure (smartphones, laptops, WiFi and cellular systems etc.). Coding is also now routinely used to transform the storage layer of distributed systems in modern datacenters supporting regenerating with locally repairable codes for distributed storage.

The basic idea of coded computing is to be able to compute in a redundant way so that the computation can be completed even without collecting the results from the stragglers. The redundancy in coded computation can be used not only to address the straggler problem in master-worker computations, but also communication bottlenecks. Coded computing has been used in a number of tree-structured master-worker parallel computations, including FFT, Matrix Multiplication, Machine Learning, and MapReduce. The method is based on the use of Lagrange interpolation polynomials and related mathematical techniques.

## 1.9 New Research Directions

As we have seen, despite more than 40 years of effort, we still do not have a universal model of computation that can handle all of the standard algorithmic structures and patterns that it needs to support, and all of the other requirements. No single model has yet demonstrated that it is satisfactory in all dimensions. We are still looking for a model that can address the spatial and temporal decomposition of computations on current and future heterogeneous massively parallel architectures.

### 1.9.1 Research Challenges for Current Models

The dataflow approach (such as TensorFlow) has been tried many times over the past 40 years, and at modest scale can be useful. However, the optimization of dataflow computations at large scale remains a major challenge. Also, as it is a relatively low-level programming model, the productivity of programmers is low, especially when trying to achieve high performance at scale with dataflow models.

The functional /algebraic approach (such as GraphBLAS) offers considerable potential in terms of improving software productivity, as it provides a flexible high-level functional abstraction that makes programming easier via higher-order functions. The main challenge in this area is to provide the necessary software automation to ensure that the required load balancing, communication-efficiency, synchronization-efficiency, etc., can be achieved without the programmer having to specify how it should be done at scale on complex heterogeneous parallel architectures.

The BSP model has clearly demonstrated massive scalability and very high performance at scale. It has also provided a foundation for other models that have come later, such as Data-Centric BSP (Pregel and many others) and nonstop resilient BSP (Cloned Computing). Data-centric BSP is now routinely used today at many companies for data analytics on graphs with billions of nodes and trillions of edges. It is also used in world-leading commercial graph database systems. Being a lower-level abstraction than the functional approach means that the programmer has potentially more flexibility, but as always this comes at the price of increased programming complexity.

As research on BSP-based models has shown, the challenges in developing a model and associated framework are not only concerned with the programming model or style. Indeed, as the programming model becomes simpler and more abstract in order to improve software productivity, these other aspects become more and more important. For example, it is easy to propose an appealing model where the programmer is presented with a simple algebraic programming abstraction based on higher-order functions. The real challenge in such a case is that any viable implementation of the model needs to offer not only a scalable high performance low-level execution engine, but also a range of powerful new technologies for cost-model driven software automation that together enable the highest levels of scalability and performance to be achieved. We need to provide not only a run-time layer for scalable heterogeneous parallel computing, but also the parallelization and optimization layer.

Can a new model be developed that can combine the software productivity advantages

of simple functional models, with the scalability, performance and resilience benefits of lower-level BSP-based models? Or is there something else that is even better that we have just not thought of yet? That remains a research challenge for the future.

It is remarkable to reflect that the early work of von Neumann in the 1940s and, before that, Turing in the 1930s, provided us with a universal model of sequential computation that has endured to this day. The stability that it provided has been essential to the remarkable growth of the global software and hardware industries over the past 60-70 years. Despite the desire for a post-von Neumann model, we are not there yet.

### 1.9.2 Scale Simplifies

Before looking at specific new research directions, it is perhaps appropriate to consider a very simple and fundamental question. Does increasing scale increase the complexity of the challenge of producing a model that can achieve all the various aspects required – Universality, Scalability, Performance, Portability, Predictability, Analyzability, Productivity etc.? Or does scale simplify the problem?

It is natural to think that increasing scale will make problems more complex and challenging, but paradoxically, as scale increases, fewer and fewer of the approaches that work at small scale can be used. An expert programmer can attempt to explicitly define every aspect of a system with a few cores. However, when the number of cores is one million or more, this becomes an unrealistic task, unless a large part of the process is uniform and automated. So, at massive scale, we can expect that few models will work, and those that do will have to be simple, uniform and highly automated. So where should we be looking for new directions?

### 1.9.3 Communication-Light Models

As we have seen, many computational problems are communication-light. In BSP terms, these are problems with a cost  $W + g * H + L * S$  where  $H$  and  $S$  are much smaller than  $W$ . For example, there are many important computational problems that are “embarrassingly parallel”, i.e. where  $H$  and  $S$  are constant for any problem size, and in some cases may be zero. A simple example of such a communication-light problem would be data-parallel video encoding.

There are also many important problems where  $S$  is constant or small, and where  $H$  grows much more slowly than  $W$ . For example, we have many “high arithmetic intensity” parallel matrix computations where  $W = O(n^3/p)$  but  $H$  is only  $O(n^2/p^{1/2})$  or less. In these communication-light scenarios, a number of the models we have considered will provide good scalability and performance. We can certainly use BSP or message passing, but other less flexible and more automated models such as MapReduce and Spark can also be used in many such cases. Using simpler models in these cases also allows us to more easily provide resilience with automated fault tolerance.

For communication-light parallel computation, the most extreme model is perhaps the disaggregated “serverless” architecture. Serverless computing (Jonas et al. 2019) has