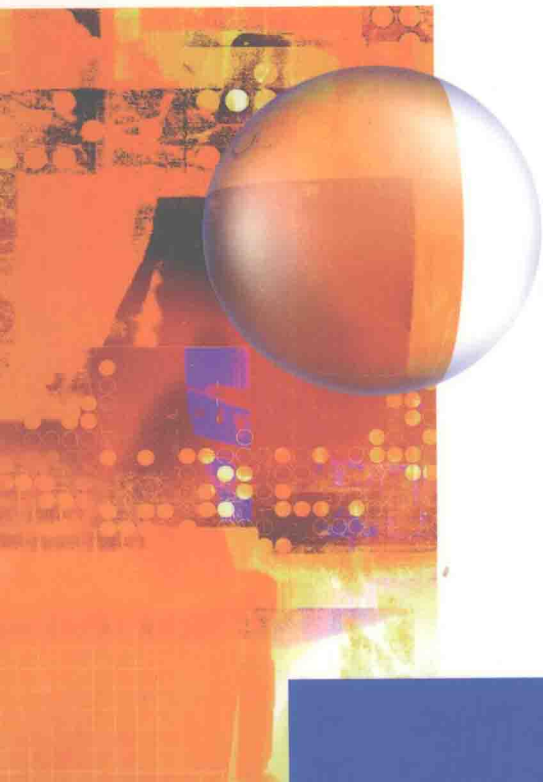# FIBRE CHANNEL for SANs

Step-by-step guidance through all 5 FC layers

Bandwidth-saving solutions for multimedia

Troubleshooting FC networks

## ALAN F. BENNER

# Fibre Channel
# for SANs

Alan F. Benner

# List of Figures

# Preface

In the preface to the first edition of this book I wrote about the difficulty of keeping up with the flow of information in a book format, since the state of the technology changes so much during the time between when a book is started and when it actually appears. Since then, the pace of change in Fibre Channel, Storage Area Networks (SANs), and server networking has, if anything, accelerated. At least a dozen new standards have been approved since the first edition, and several more are either under development or have been superseded by better work.

Fortunately, however, most of the current progress is up in the network management and data management layers, so that the basic operation of the network, which is the primary subject of this book, has pretty well stabilized. It's possible, therefore, to expect that the subject matter here will stay stable for a reasonable period of time.

The most fundamental change in Fibre Channel since the first edition of this book, however, has been in where it gets used. When Fibre Channel was first created, it was not clear exactly which problems it would be most effective at solving, since it was such a widely applicable technology. Since then, Fibre Channel has become the de facto standard interconnect technology for storage networking. In fact, the term Storage Area Network, newly defined since that time, has currently become essentially synonymous with Fibre Channel.

Since Fibre Channel has become so tightly linked with storage networking, much of the progress in Fibre Channel has been towards solving specific storage networking problems. In this book, however, I've tried to keep a general focus on the whole of Fibre Channel technology, since restricting the discussion to only those aspects important in current-day SANs runs the risk of ignoring parts of the technology that will be extremely important in the future.

The intention of this book is two-fold. First, it's intended to be an overview guide to the concepts, the structures, and the goals of the Fibre Channel architecture at a fairly detailed level. A dedicated reader should be able to use it to understand most of the details of Fibre Channel before referring to the ANSI materials for authoritative information. Second, and perhaps more importantly, it's an attempt to show the reasons why the network is designed the way it is. Networking technology is consistently becoming more and more important, and new technologies are being developed regularly. As new networking technologies are invented, I hope that this book will help people exploit the goods features in the already-existing networks.

A few words on notation. Determining a consistent notation in this type of work is not trivial, since the subject matter bridges both computer and communications arenas, which have traditionally used slightly different

notations. For example, communications data rates are generally measured in megabits per second, where "mega" means $10^6$, while computer data is measured in megabytes, where "mega" means $2^{20}$. The text mixes both somewhat, using "b" to represent bits, as in Gbps, and "B" to represent bytes, as in MBps. In recognition of the communications-oriented nature of the subject, the prefixes "mega (M)" and "giga (G)" will mean $10^6$ and $10^9$ here, rather than $2^{20}$ and $2^{30}$. All numbers in the book are written in binary (b'0110 0101'), hexadecimal (x'FF FFFD'), or decimal (65,532) formats. Single bits are written as 1 or 0.

A number of common words, such Sequence, Exchange, and Connection, have specific meanings in Fibre Channel that are quite distinct from their common usage. In this book, words with specific Fibre Channel meanings are generally capitalized to distinguish them from the common usage. This capitalization generally matches the format used in the ANSI standards documents. Information provided here is in the public domain, through generally available books, articles, ANSI documents, or other reference material.

Several terms used in this book, such as ATM and HIPPI, are taken from other architectures. Any trademarks used are properties of their rightful owners. Ethernet is a trademark of the Xerox Corporation. ESCON, FICON, and SBCON are trademarks or registered trademark of the IBM Corporation in the United States or other countries or both. InfiniBand is a registered trademark of the InfiniBand Trade Association.

The book is organized as follows. The first few chapters give an overview of the features and goals for the Fibre Channel architecture, along with an example of how data is transmitted under a Fibre Channel network.

The middle chapters cover the concepts and structures of Fibre Channel in a fair amount of detail. These include chapters on all of the Fibre Channel physical components and logical constructs, supported functions, flow control, and error recovery.

The final chapters cover configuration and operation of the Arbitrated Loop topology, mapping of Fibre Channel constructs to upper level protocols such as SCSI and the IP level of TCP/IP. These chapters show how Fibre Channel fits in with currently existing software and operating system levels. In the last chapter I have taken the opportunity to make some predictions for the future of Fibre Channel, SANs, and server networking, as far as I dare.

Thanks are due to far too many people to list here, but I'll try — I apologize in advance to those I may have missed. Many thanks to Carl Zeitler, Ki Won Lee, Mike Yang, Dan Eisenhower, Ron Cash, Roger Weekly, Giles Frazier, Jerry Chapman, Jerry Rouse, Jonathan Thatcher, Bill George, Al Widmer, Tom McConathy, Casey Cannon, Gary Nutt, R. Bryan Cook, Paul Green, Dal Allan, Martin Sachs, Horst Truestedt, Richard Taborek, Schelto Van Doorn, and Roger Cummings, who helped especially during the writing of the first edition. Thanks to Herman Presby, Ivan Kaminow, and Jon Sauer,

# Table of Contents