

Overview and Background

This is a course for computer system designers and builders, and for people who want to really understand how systems work, especially concurrent, distributed, and fault-tolerant systems.

The course teaches you how to write precise specifications for any kind of computer system, what it means for an implementation to satisfy a specification, and how to prove that it does. It also shows you how to use the same methods less formally, and gives you some suggestions for deciding how much formality is appropriate (less formality means less work, and often a more understandable spec, but also more chance to overlook an important detail).

The course also teaches you a lot about the topics in computer systems that we think are the most important: persistent storage, concurrency, naming, networks, distributed systems, transactions, fault tolerance, and caching. The emphasis is on

- careful specifications of subtle and sometimes complicated things,
- the important ideas behind good implementations, and
- how to understand what makes them actually work.

We spend most of our time on specific topics, but we use the general techniques throughout. We emphasize the ideas that different kinds of computer system have in common, even when they have different names.

The course uses a formal language called Spec for writing specs and implementations; you can think of it as a very high level programming language. There is a good deal of written introductory material on Spec (explanations and finger exercises) as well as a reference manual and a formal semantics. We introduce Spec ideas in class as we use them, but we do not devote class time to teaching Spec per se; we expect you to learn it on your own from the handouts.

Because we write specs and do proofs, you need to know something about logic. Since many people don't, there is a concise treatment of the logic you will need at the end of this handout.

This is not a course in computer architecture, networks, operating systems, or databases. We will not talk in detail about how to implement pipelines, memory interconnects, multiprocessors, routers, data link protocols, network management, virtual memory, scheduling, resource allocation, SQL, relational integrity, or TP monitors, although we will deal with many of the ideas that underlie these mechanisms.

Topics

General

Specifications as state machines.

The Spec language for describing state machines (writing specs and implementations).

What it means to implement a spec.

Using abstraction functions and invariants to prove that a program implements a spec.

What it means to have a crash.

What every system builder needs to know about performance.

Specific

Disks and file systems.

Practical concurrency using mutexes (locks) and condition variables; deadlock.

Hard concurrency: models, specs, proofs, and examples.

Transactions: simple, cached, concurrent, distributed.

Naming: principles, specs, and examples.

Distributed systems: communication, fault-tolerance, and autonomy.

Networking: links, switches, reliable messages and connections.

Remote procedure call and network objects.

Fault-tolerance, availability, consensus and replication.

Caching and distributed shared memory.

System management.

Previous editions of the course have also covered security (authentication, authorization, encryption, trust), but this year we are omitting this topic in order to spend more time on concurrency and semantics.

Prerequisites

There are no formal prerequisites for the course. However, we assume some knowledge both of computer systems and of mathematics. If you have taken 6.033 and 6.042, you should be in good shape. If you are missing some of this knowledge you can pick it up as we go, but if you are missing a lot of it you can expect to have serious trouble. It's also important to have a certain amount of maturity: enough experience with systems and mathematics to feel comfortable with the basic notions and to have some reliable intuition.

If you know the meaning of the following words, you have the necessary background. If a lot of them are unfamiliar, this course is probably not for you.

Systems

Cache, virtual memory, page table, pipeline

Process, scheduler, address space, priority

Thread, mutual exclusion (locking), semaphore, producer-consumer, deadlock

Transaction, commit, availability, relational data base, query, join

File system, directory, path name, striping, RAID
LAN, switch, routing, connection, flow control, congestion
Capability, access control list, principal (subject)

If you have not already studied Lampson's paper on hints for system design, you should do so as background for this course. It is Butler Lampson, Hints for computer system design, *Proceedings of the Ninth ACM Symposium on Operating Systems Principles*, October 1983, pp 33-48. There is a pointer to it on the course Web page.

Programming

Invariant, precondition, weakest precondition, fixed point
Procedure, recursion, stack
Data type, sub-type, type-checking, abstraction, representation
Object, method, inheritance
Data structures: hash table, binary search, B-tree, graph

Mathematics

Function, relation, set, transitive closure
Logic: proof, induction, de Morgan's laws, implication, predicate, quantifier
Probability: independent events, sampling, Poisson distribution
State machine, context-free grammar
Computational complexity, unsolvable problem

If you haven't been exposed to formal logic, you should study the summary at the end of this handout.

References

These are places to look when you want more information about some topic covered or alluded to in the course, or when you want to follow current research. You might also wish to consult Prof. Saltzer's bibliography for 6.033, which you can find on the course web page.

Books

Some of these are fat books better suited for reference than for reading cover to cover, especially Cormen, Leiserson, and Rivest, Jain, Mullender, Hennessy and Patterson, and Gray and Reuter. But the last two are pretty easy to read in spite of their encyclopedic character.

Systems programming: Greg Nelson, ed., *Systems Programming with Modula-3*, Prentice-Hall, 1991. Describes the language, which has all the useful features of C++ but is much simpler and less error-prone, and also shows how to use it for concurrency (a version of chapter 4 is a handout in this course), an efficiently customizable i/o streams package, and a window system.

Performance: Jon Bentley, *Writing Efficient Programs*, Prentice-Hall, 1982. Short, concrete, and practical. Raj Jain, *The Art of Computer Systems Performance Analysis*, Wiley, 1991. Tells you much more than you need to know about this subject, but does have a lot of realistic examples.

Algorithms and data structures: Robert Sedgwick, *Algorithms*, Addison-Wesley, 1983. Short, and usually tells you what you need to know. Tom Cormen, Charles Leiserson, and Ron Rivest, *Introduction to Algorithms*, McGraw-Hill, 1989. Comprehensive, and sometimes valuable for that reason, but usually tells you a lot more than you need to know.

Distributed algorithms: Nancy Lynch, *Distributed Algorithms*, Morgan Kaufmann, 1996. The bible for distributed algorithms. Comprehensive, but a much more formal treatment than in this course.

Computer architecture: John Hennessy and David Patterson, *Computer Architecture: A Quantitative Approach*, 2nd edition, Morgan Kaufmann, 1995. The bible for computer architecture. The second edition has lots of interesting new material, especially on multiprocessor memory systems and interconnection networks. There's also a good appendix on computer arithmetic; it's useful to know where to find this information, though it has nothing to do with this course.

Transactions, data bases, and fault-tolerance: Jim Gray and Andreas Reuter, *Transaction Processing: Concepts and Techniques*, Morgan Kaufmann, 1993. The bible for transaction processing, with much good material on data bases as well; it includes a lot of practical information that doesn't appear elsewhere in the literature.

Networks: Radia Perlman, *Interconnections: Bridges and Routers*, Addison-Wesley, 1992. Not exactly the bible for networking, but tells you nearly everything you might want to know about how packets are actually switched in computer networks.

Distributed systems: Sape Mullender, ed., *Distributed Systems*, 2nd ed., Addison-Wesley, 1993. A compendium by many authors that covers the field fairly well. Some chapters are much more theoretical than this course. Chapters 10 and 11 are handouts in this course. Chapters 1, 2, 8, and 12 are also recommended. Chapters 16 and 17 are the best you can do to learn about real-time computing; unfortunately, that is not saying much.

User interfaces: Alan Cooper, *About Face*, IDG Books, 1995. Principles, lots of examples, and opinionated advice, much of it good, from the original designer of Visual Basic.

Journals

You can find all of these in the LCS reading room. The cryptic strings in brackets are call numbers there.

For the current literature, the best sources are the proceedings of the following conferences. 'Sig' is short for "Special Interest Group", a subdivision of the ACM that deals with one field of computing. The relevant ones for systems are SigArch for computer architecture, SigPlan for programming languages, SigOps for operating systems, SigComm for communications, and SigMod for data bases.

Symposium on Operating Systems Principles (SOSP; published as special issues of ACM *SigOps Operating Systems Review*; fall of odd-numbered years) [P4.35.06]

Operating Systems Design and Implementation (OSDI; Usenix Association, now published as special issues of ACM *SigOps Review*; fall of even-numbered years) [P4.35.U71]

Architectural Support for Programming Languages and Operating Systems (ASPLOS; published as special issues of ACM SigOps *Operating Systems Review* or SigArch *Computer Architecture News* or SigPlan *Notices*; fall of even-numbered years) [P6.29.A7]

Applications, Technologies, Architecture, and Protocols for Computer Communication, (SigComm conference; published as special issues of ACM SigComm *Computer Communication Review*; annual) [P6.24.D31]

Principles of Distributed Computing (PODC; ACM; annual) [P4.32.D57]

Very Large Data Bases (VLDB; Morgan Kaufmann; annual) [P4.33.V4]

International Symposium on Computer Architecture (ISCA; published as special issues of ACM SigArch *Computer Architecture News*; annual) [P6.20.C6]

Less up to date, but more selective, are the journals. Often papers in these journals are revised versions of papers from the conferences listed above.

ACM Transactions on Computer Systems

ACM Transactions on Database Systems

ACM Transactions on Programming Languages and Systems

There are often good survey articles in the less technical IEEE journals:

IEEE Computer, Networks, Communication, Software

Rudiments of Logic

Propositional Logic

The basic type is `Bool`, which contains two elements `true` and `false`. Expressions in these operators (and the other ones introduced later) are called ‘propositions’.

Basic operators. These are \wedge (and), \vee (or), and \sim (not).¹ The meaning of these operators can be conveniently given by a ‘truth table’ which lists the value of $a \text{ op } b$ for each possible combination of values of a and b (the operators on the right are discussed later) along with some popular names for certain expressions and their operands.

		negation	conjunction	disjunction	equality		implication
		not	and	or			implies
a	b	$\sim a$	$a \wedge b$	$a \vee b$	$a = b$	$a \# b$	$a \Rightarrow b$
T	T	F	T	T	T	F	T
T	F		F	T	F	T	F
F	T	T	F	T	F	T	T
F	F		F	F	T	F	T
name of a			conjunct	disjunct			antecedent
name of b			conjunct	disjunct			consequent

Note: In Spec we write `==>` instead of the \Rightarrow that mathematicians use for implication. Logicians write \supset for implication, which looks different but is shaped like the $>$ part of \Rightarrow .

In case you have an expression that you can’t simplify, you can always work out its truth value by exhaustively enumerating the cases in truth table style. Since the table has only four rows, there are only 16 Boolean operators, one for each possible arrangement of T and F in a column. Most of the ones not listed don’t have common names, though ‘not and’ is called ‘nand’ and ‘not or’ is called ‘nor’ by logic designers.

The \wedge and \vee operators are
 commutative and
 associative and
 distribute over each other.

That is, they are just like $*$ (times) and $+$ (plus) on integers, except that $+$ doesn’t distribute over $*$:

$$a + (b * c) \neq (a + b) * (a + c),$$

but \vee does distribute over \wedge :

$$a \vee (b \wedge c) = (a \vee b) \wedge (a \vee c).$$

An operator that distributes over \wedge is called ‘conjunctive’; one that distributes over \vee is called ‘disjunctive’. So both \wedge and \vee are both conjunctive and disjunctive. This takes some getting used to.

¹ It’s possible to write all three in terms of the single operator ‘nor’, but our goal is clarity, not minimality.

The relation between these operators and \sim is given by DeMorgan's laws (sometimes called the "bubble rule" by logic designers), which say that you can push \sim inside \wedge or \vee by flipping from one to the other:

$$\sim (a \wedge b) = \sim a \vee \sim b$$

$$\sim (a \vee b) = \sim a \wedge \sim b$$

Because `Bool` is the result type of relations like $=$, we can write expressions that mix up relations with other operators in ways that are impossible for any other type. Notably

$$(a = b) = ((a \wedge b) \vee (\sim a \wedge \sim b))$$

Some people feel that the outer $=$ in this expression is somehow different from the inner one, and write it \equiv . Experience suggests, however, that this is often a harmful distinction to make.

Implication. We can define an ordering on `Bool` with `false > true`, that is, `false` is greater than `true`. The non-strict version of this ordering is called 'implication' and written \Rightarrow (rather than \geq as we do with other types; logicians write it \supset , which also looks like an ordering symbol). So `(true \Rightarrow false) = false` (think of this as: "true is greater than or equal to false" is false) but all other combinations are `true`. The expression `a \Rightarrow b` is pronounced "a implies b", or "if a then b".²

There are lots of rules for manipulating expressions containing \Rightarrow ; the most useful ones are given below. If you remember that \Rightarrow is an ordering you'll find it easy to remember most of the rules, but if you forget the rules or get confused, you can turn the \Rightarrow into \vee by the rule

$$(a \Rightarrow b) = \sim a \vee b$$

and then just use the simpler rules for \wedge , \vee , and \sim . So remember this even if you forget everything else.

The point of implication is that it tells you when one proposition is stronger than another, in the sense that if the first one is true, the second is also true (because if both `a` and `a \Rightarrow b` are `true`, then `b` must be `true` since it can't be `false`).³ So we use implication all the time when reasoning from premises to conclusions. Two more ways to pronounce `a \Rightarrow b` are "a is stronger than b" and "b follows from a". The second pronunciation suggests that it's sometimes useful to write the operands in the other order, as `b \Leftarrow a`, which can also be pronounced "b is weaker than a" or "b only if a"; this should be no surprise, since we do it with other orderings.

Of course, implication has the properties we expect of an ordering:

² It sometimes seems odd that `false` implies `b` regardless of what `b` is, but the "if ... then" form makes it clearer what is going on: if `false` is `true` you can conclude anything, but of course it isn't. A proposition that implies `false` is called 'inconsistent' because it implies anything. Obviously it's bad to think that an inconsistent proposition is true. The most likely way to get into this hole is to think that each of a collection of innocent looking propositions is true when their conjunction turns out to be inconsistent.

³ It may also seem odd that `false > true` rather than the other way around, since `true` seems better and so should be bigger. But in fact if we want to conclude lots of things, being close to `false` is better because if `false` is true we can conclude anything, but knowing that `true` is true doesn't help at all. Strong propositions are as close to `false` as possible; this is logical brinkmanship. For example, `a \wedge b` is closer to `false` than `a` (there are more values of the variables `a` and `b` that make it `false`), and clearly we can conclude more things from it than from `a` alone.

Transitive: If $a \Rightarrow b$ and $b \Rightarrow c$ then $a \Rightarrow c$.⁴

Reflexive: $a \Rightarrow a$.

Anti-symmetric: If $a \Rightarrow b$ and $b \Rightarrow a$ then $a = b$.⁵

Furthermore, \sim reverses the sense of implication (this is called the ‘contrapositive’):

$$(a \Rightarrow b) = (\sim b \Rightarrow \sim a)$$

More generally, you can move a disjunct on the right to a conjunct on the left by negating it.

Thus

$$(a \Rightarrow b \vee c) = (a \wedge \sim b \Rightarrow c)$$

As special cases in addition to the contrapositive we have

$$(a \Rightarrow b) = (a \wedge \sim b \Rightarrow \text{false}) = \sim(a \wedge \sim b) \vee \text{false} = \sim a \vee b$$

$$(a \Rightarrow b) = (\text{true} \Rightarrow \sim a \vee b) = \text{false} \vee \sim a \vee b = \sim a \vee b$$

since *false* and *true* are the identities for \vee and \wedge .

We say that an operator *op* is ‘monotonic’ in an operand if replacing that operand with a stronger (or weaker) one makes the result stronger (or weaker). Precisely, “*op* is monotonic in its first operand” means that if $a \Rightarrow b$ then $(a \text{ op } c) \Rightarrow (b \text{ op } c)$. Both \wedge and \vee are monotonic; in fact, any conjunctive operator is monotonic, because if $a \Rightarrow b$ then $a = (a \wedge b)$, so $a \text{ op } c = (a \wedge b) \text{ op } c = a \text{ op } c \wedge b \text{ op } c \Rightarrow b \text{ op } c$.

If you know what a lattice is, you will find it useful to know that the set of propositions forms a lattice with \Rightarrow as its ordering and (remember, think of \Rightarrow as “greater than or equal”):

top = *false*

bottom = *true*

meet = \wedge

join = \vee

least upper bound, so $(a \wedge b) \Rightarrow a$ and $(a \wedge b) \Rightarrow b$

greatest lower bound, so $a \Rightarrow (a \vee b)$ and $b \Rightarrow (a \vee b)$

This suggests two more expressions that are equivalent to $a \Rightarrow b$:

$(a \Rightarrow b) = (a = (a \wedge b))$ ‘and’ing a weaker term makes no difference,
because $a \Rightarrow b$ iff $a =$ least upper bound(a, b).

$(a \Rightarrow b) = (b = (a \vee b))$ ‘or’ing a stronger term makes no difference,
because $a \Rightarrow b$ iff $b =$ greatest lower bound(a, b).

Predicate logic

Propositions that have free variables, like $x < 3$ or $x < 3 \Rightarrow x < 5$, demand a little more machinery. You can turn such a proposition into one without a free variable by substituting some value for the variable. Thus if $P(x)$ is $x < 3$ then $P(5)$ is $5 < 3 = \text{false}$. To get rid of the free variable without substituting a value for it, you can take the ‘and’ or ‘or’ of the proposition for all the possible values of the free variable. These have special names and notation⁶:

⁴ We can also write this $((a \Rightarrow b) \wedge (b \Rightarrow c)) \Rightarrow (a \Rightarrow c)$.

⁵ Thus $(a = b) = (a \Rightarrow b \wedge b \Rightarrow a)$, which is why $a = b$ is sometimes pronounced “a if and only if b” and written “a iff b”.

⁶ There is no agreement on what symbol should separate the $\forall x$ or $\exists x$ from the $P(x)$. We use ‘|’ here as Spec does, but other people use ‘.’ or ‘:’ or just a space, or write $(\forall x)$ and $(\exists x)$. Logicians traditionally write $(\forall x)$ and $(\exists x)$.

$\forall x \mid P(x) = P(x_1) \wedge P(x_2) \wedge \dots$ for all $x, P(x)$. In Spec,
 (ALL $x \mid P(x)$) or $\wedge : \{x \mid P(x)\}$
 $\exists x \mid P(x) = P(x_1) \vee P(x_2) \vee \dots$ there exists an x such that $P(x)$. In Spec,
 (EXISTS $x \mid P(x)$) or $\vee : \{x \mid P(x)\}$

Here the x_i range over all the possible values of the free variables.⁷ The first is called ‘universal quantification’; as you can see, it corresponds to conjunction. The second is called ‘existential quantification’ and corresponds to disjunction. If you remember this you can easily figure out what the quantifiers do with respect to the other operators.

In particular, DeMorgan’s laws generalize to quantifiers:

$$\sim (\forall x \mid P(x)) = (\exists x \mid \sim P(x))$$

$$\sim (\exists x \mid P(x)) = (\forall x \mid \sim P(x))$$

Also, because \wedge and \vee are conjunctive and therefore monotonic, \forall and \exists are conjunctive and therefore monotonic.

It is not true that you can reverse the order of \forall and \exists , but it’s sometimes useful to know that having \exists first is stronger:

$$\exists y \mid \forall x \mid P(x, y) \Rightarrow \forall x \mid \exists y \mid P(x, y)$$

Intuitively this is clear: a y that works for every x can surely do the job for each particular x .

If we think of P as a relation, the consequent in this formula says that P is total (relates every x to some y). It doesn’t tell us anything about how to find a y that is related to x . As computer scientists, we like to be able to compute things, so we prefer to have a function that computes y , or the set of y ’s, from x . This is called a ‘Skolem function’; in Spec you write $P.\text{func}$ (or $P.\text{setF}$ for the set). $P.\text{func}$ is total if P is total. Or, to turn this around, if we have a total function f such that $\forall x \mid P(x, f(x))$, then certainly $\forall x \mid \exists y \mid P(x, y)$; in fact, $y = f(x)$ will do. Amazing.

⁷ In general this might not be a countable set, so the conjunction and disjunction are written in a somewhat misleading way, but this complication won’t make any difference to us.

Summary of logic

The \wedge and \vee operators are commutative and associative and distribute over each other.

DeMorgan's laws: $\sim (a \wedge b) = \sim a \vee \sim b$
 $\sim (a \vee b) = \sim a \wedge \sim b$

Implication: $(a \Rightarrow b) = \sim a \vee b$

Implication is the ordering in a lattice (a partially ordered set in which every subset has a least upper and a greatest lower bound) with

top	= false	so false \Rightarrow true
bottom	= true	
meet	= \wedge	least upper bound, so $(a \wedge b) \Rightarrow a$
join	= \vee	greatest lower bound, so $a \Rightarrow (a \vee b)$

For all x , $P(x)$:

$$\forall x \mid P(x) = P(x_1) \wedge P(x_2) \wedge \dots$$

There exists an x such that $P(x)$:

$$\exists x \mid P(x) = P(x_1) \vee P(x_2) \vee \dots$$

Index for logic

\sim , 6	implication, 6, 7
\Rightarrow , 6	join, 8
\Rightarrow , 6	lattice, 8
ALL, 9	least upper bound, 8
and, 6	meet, 8
antecedent, 6	monotonic, 8
Anti-symmetric, 8	negation, 6
associative, 6	not, 6
bottom, 8	only if, 7
commutative, 6	operators, 6
conjunction, 6	or, 6
conjunctive, 6	ordering on Bool, 7
consequent, 6	predicate logic, 8
contrapositive, 8	propositions, 6
DeMorgan's laws, 7, 9	quantifiers, 9
disjunction, 6	reflexive, 8
disjunctive, 6	Skolem function, 9
distribute, 6	stronger than, 7
existential quantification, 9	top, 8
EXISTS, 9	transitive, 8
follows from, 7	truth table, 6
free variables, 8	universal quantification, 9
greatest lower bound, 8	weaker than, 7
if a then b, 7	