

# Principles for Computer System Design

**10 years ago:** *Hints for Computer System Design*

**Not that much learned since then—disappointing**

*Instead of standing on each other's shoulders, we stand on each other's toes.* (Hamming)

**One new thing: How to build systems more precisely**

*If you think systems are expensive, try chaos.*

# Collaborators

**Bob Taylor**

**Chuck Thacker**

Workstations: Alto, Dorado, Firefly  
Networks: AN1, AN2

**Charles Simonyi**

Bravo WYSIWYG editor

**Nancy Lynch**

Reliable messages

**Howard Sturgis**

Transactions

**Martin Abadi**

Security

**Mike Burrows**

**Morrie Gasser**

**Andy Goldstein**

**Charlie Kaufman**

**Ted Wobber**

# From Interfaces to Specifications

## Make modularity precise

*Divide and conquer (Roman motto)*

*Design*

*Correctness*

*Documentation*

## Do it recursively

*Any idea is better when made recursive (Randell)*

*Refinement:* One man's implementation is another man's spec.  
*(adapted from Perlis)*

*Composition:* Use actions from one spec in another.

# Specifying a System with State

**A *safety* property: nothing bad ever happens**

**Defined by a *state machine*:**

*state*: a set of values, usually divided into named *variables*

*actions*: named changes in the state

**A *liveness* property: something good eventually happens**

**These define *behavior*: all the possible sequence of actions**

**Examples of systems with state:**

Data abstractions

Concurrent systems

Distributed systems

You can't observe the actual state of the system from outside.  
All you can see is the results of actions.

# Editable Formatted Text

**State**     *text*: sequence of (Char, Property)     **H**e l l o

**Actions**     *get*(2) returns ('e', (Times-Roman, ...))

*replace*(3, 5, 2, 3, (a p p l e))     **H**e l **p**

H e l l o

*look*(0, 5, *italic* := true)

***H*** e l l o

This interface was used in the Bravo editor.  
The implementation was about 20k lines of code.

# How to Write a Spec

## Figure out what the state is

Choose it to make the spec clear, not to match the code.

## Describe the actions

What they do to the state

What they return

## Helpful hints

Notation is important; it helps you to think about what's going on.

Invent a suitable vocabulary.

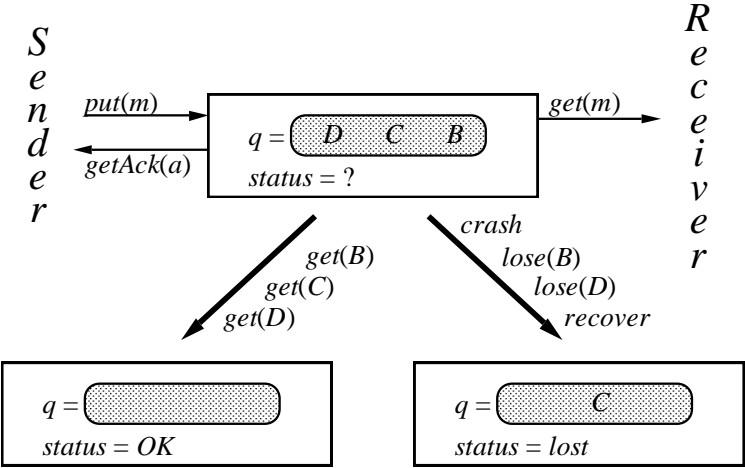
Fewer actions are better.

*Less is more.*

More non-determinism is better; it allows more implementations.

*I'm sorry I wrote you such a long letter; I didn't have time to write a short one.*  
(Pascal)

# Reliable Messages



# Spec for Reliable Messages

*q* : sequence[*M*] :=  $\langle \rangle$   
*status* : {*OK*, *lost*, ?} := *lost*  
*rec<sub>s/r</sub>* : Boolean := *false* (short for ‘recovering’)

Name	Guard	Effect	Name	Guard	Effect
<b>**put(<i>m</i>)</b>		append <i>m</i> to <i>q</i> , <i>status</i> := ?	<b>*get(<i>m</i>)</b>	<i>m</i> first on <i>q</i>	remove head of <i>q</i> , if <i>q</i> = $\langle \rangle$ , <i>status</i> = ? then <i>status</i> := <i>OK</i>
<b>*getAck(<i>a</i>)</b>	<i>status</i> = <i>a</i>	<i>status</i> := <i>lost</i>			
<i>lose</i>	<i>rec<sub>s</sub></i> or <i>rec<sub>r</sub></i>	delete some element from <i>q</i> ; if it’s the last then <i>status</i> := <i>lost</i> , or <i>status</i> := <i>lost</i>			



# What “Implements” Means?

**Divide actions into *external* and *internal*.**

**Y implements X if**

every external behavior of Y is an external behavior of X, and  
Y’s liveness property implies X’s liveness property.

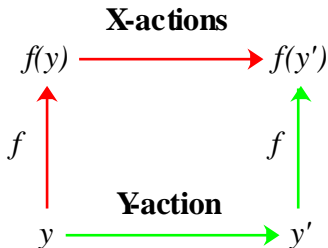
This expresses the idea that Y implements X if  
you can’t tell Y apart from X by looking only at the external actions.

# Proving that Y implements X

Define an *abstraction function*  $f$  from the state of Y to the state of X.

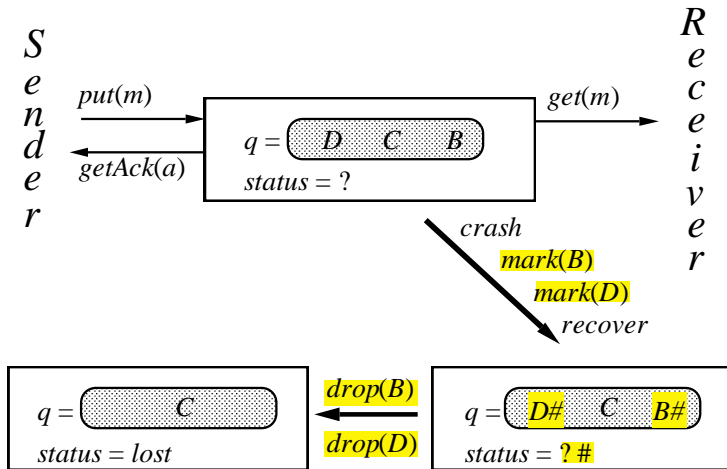
Show that Y *simulates* X:

- 1)  $f$  maps initial states of Y to initial states of X.
- 2) For each Y-action and each state  $y$   
there is a sequence of X-actions that is the same externally,  
such that the diagram commutes.



This always works!

# Delayed-Decision Spec: Example

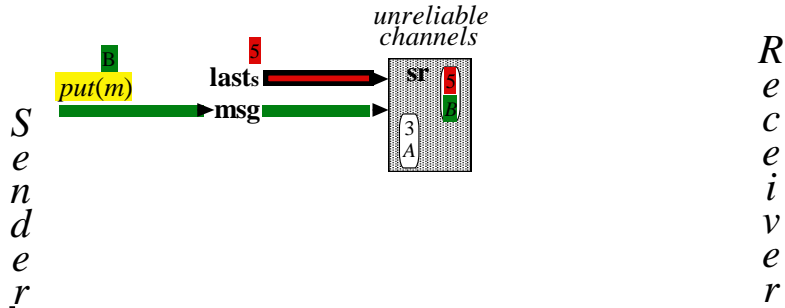


The implementer wants the spec as non-deterministic as possible, to give him more freedom and make it easier to show correctness.

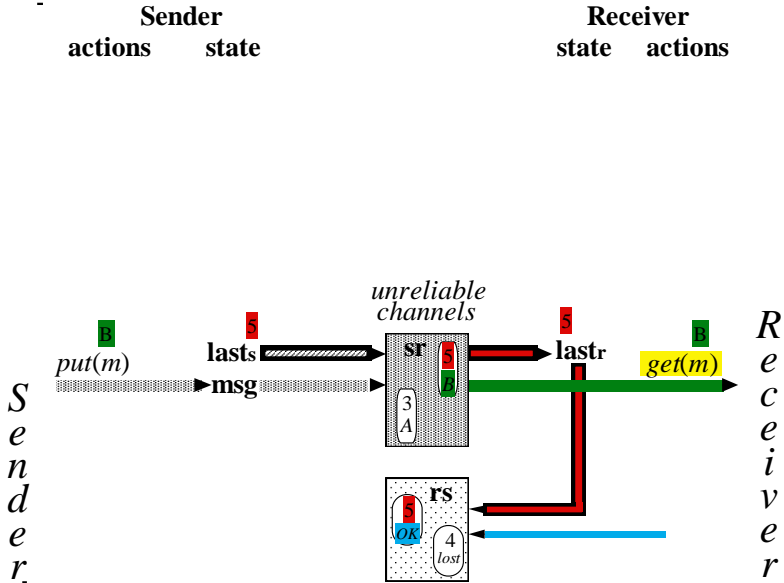
# A Generic Protocol G (1)

Sender  
actions state

Receiver  
state actions

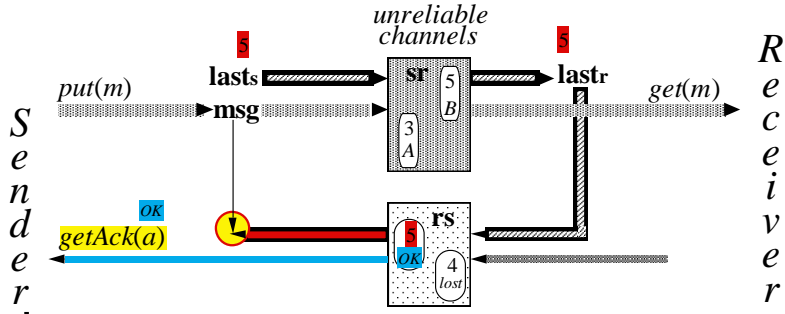


# A Generic Protocol G (2)

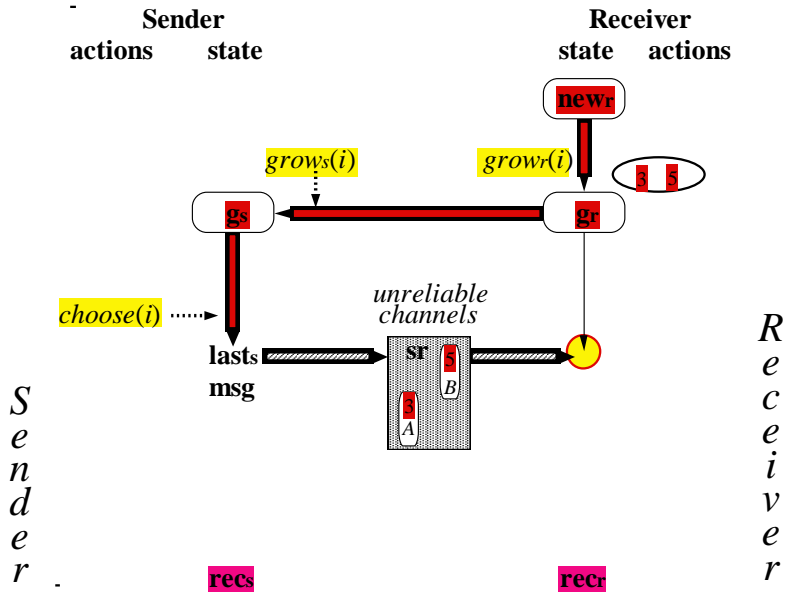


# A Generic Protocol G (3)

	<b>Sender</b>		<b>Receiver</b>
	actions	state	state actions

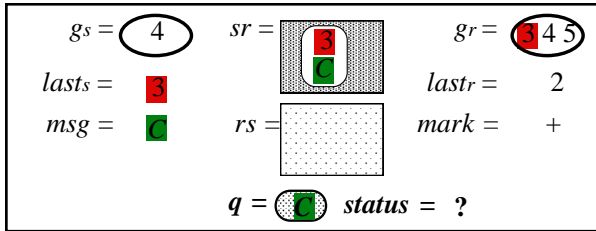


# A Generic Protocol G (4)



# G at Work

*S  
e  
n  
d  
e  
r*

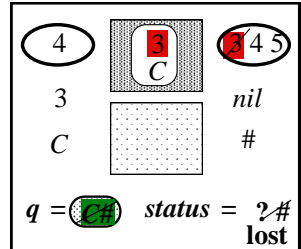
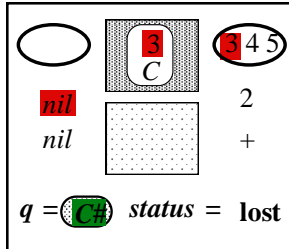
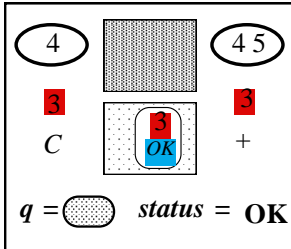


*R  
e  
c  
e  
i  
v  
e  
r*

*get(C)*

*crashes*

*crash<sub>r</sub>; recover  
(before strikeout)  
shrink<sub>r</sub>(3)  
(after strikeout)*





# Abstraction Function for G

$cur-q$  =  $\langle msg \rangle$  if  $msg \neq nil$  and  $(last_s = nil$  or  $last_s \in g_r)$   
           $\langle \rangle$  otherwise

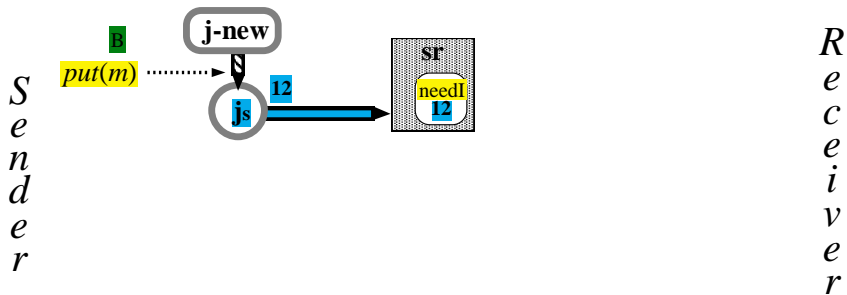
$old-q$  = the messages in  $sr$  with  $i$ 's that are good and not =  $last_s$

$q$  =  $old-q + cur-q$

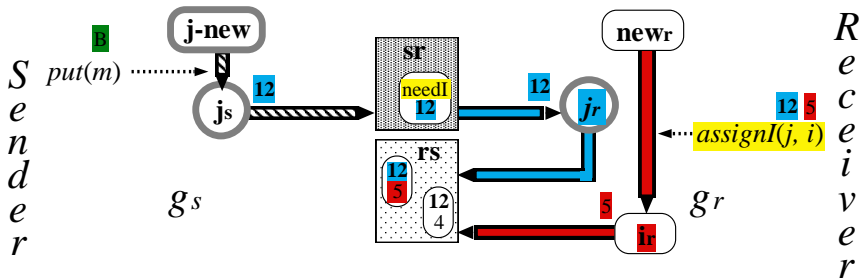
$status$  ? if  $cur-q \neq \langle \rangle$   
          OK if  $last_s = last_r \neq nil$   
          lost if  $last_s \notin (g_r \cup \{last_r\})$  or  $last_s = nil$

$rec_{s/r}$  =  $rec_{s/r}$

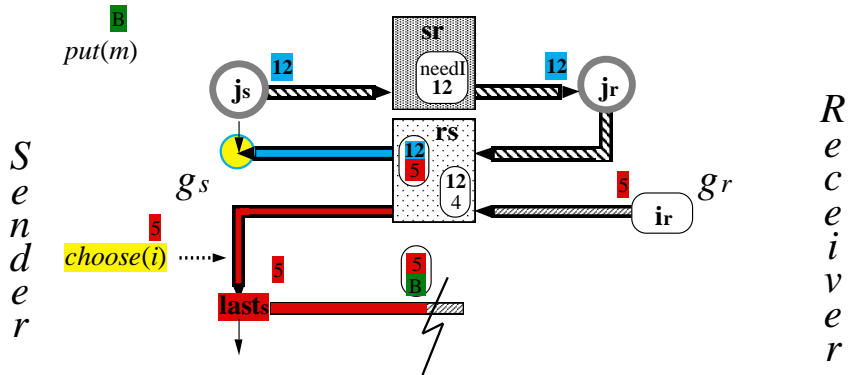
# The Handshake Protocol H (1)



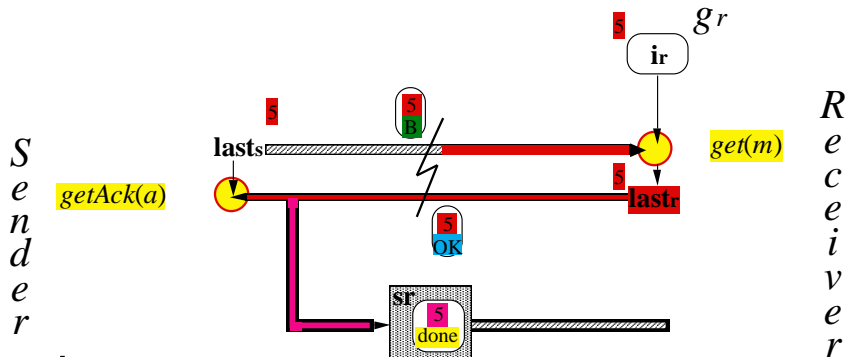
# The Handshake Protocol H (2)



# The Handshake Protocol H (3)

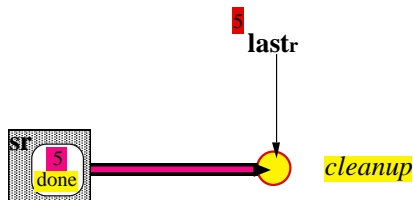


# The Handshake Protocol H (4)



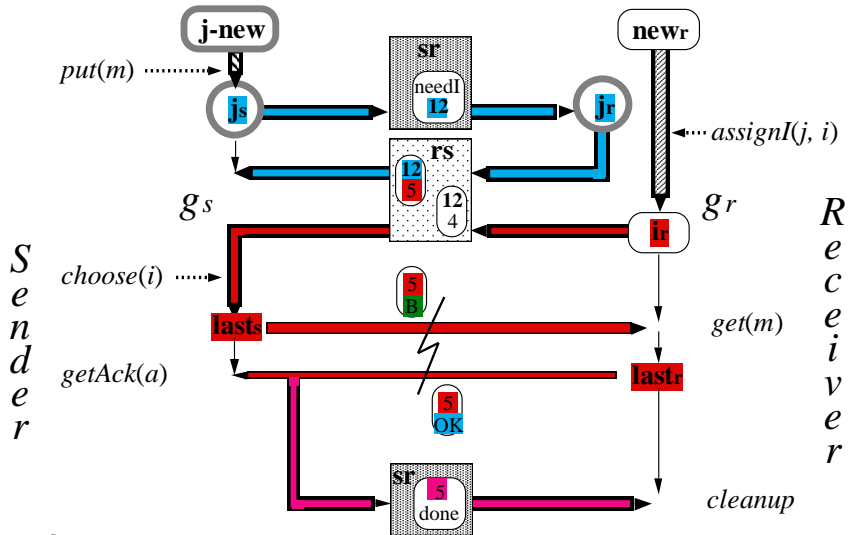
# The Handshake Protocol H (5)

*S  
e  
n  
d  
e  
r*



*R  
e  
c  
e  
i  
v  
e  
r*

# The Handshake Protocol H (6)



# Abstraction Function for H

**G**

**H**

$g_s$  the  $i$ 's with  $(j_s, i)$  in  $rs$

$g_r$   $\{i_r\} - \{nil\}$

$sr$  and  $rs$  the  $(I, M)$  and  $(I, A)$  messages in  $sr$  and  $rs$

$new_{s/r}$ ,  $last_{s/r}$ , and  $msg$  are the same in G and H

$grow_r(i)$  receiver sets  $i_r$  to an identifier from  $new_r$

$grow_s(i)$  receiver sends  $(j_s, i)$

$shrink_s(i)$  channel  $rs$  loses the last copy of  $(j_s, i)$

$shrink_r(i)$  receiver gets  $(i_r, done)$

*An efficient program is an exercise in logical brinksmanship.  
(Dijkstra)*



# Reliable Messages: Summary

## Ideas

Identifiers on messages

Sets of good identifiers, sender's  $\subseteq$  receiver's

Cleanup

**The spec is simple.**

**Implementations are subtle because of crashes.**

The abstraction functions reveal their secrets.

The subtlety can be factored in a precise way.

# Atomic Actions

$S$  : *State*

$X$	$Y$
5	5
	$do(x := x-1)$
4	5
	$do(y := y+1)$
4	6

**Name**

**Guard**

**Effect**

---

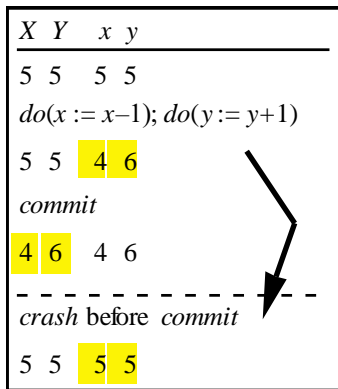
$do(a):Val$

$(S, val) := a(S)$

*A distributed system is a system in which I can't get my work done because a computer has failed that I've never even heard of.*  
*(Lamport)*

# Transactions: One Action at a Time

$S$ ,  $s$  : *State*



**Name**

**Guard**

**Effect**

*do*( $a$ ):*Val*

$(s, val) := a(s)$

*commit*

$S := s$

---

*crash*

$s := S$

# Server Failures

$S, s : State$

$\phi : \{nil, run\} := nil$

$X$	$Y$	$x$	$y$	$\phi$
5	5	5	5	nil
<i>do(x := x-1); do(y := y+1)</i>				
5	5	4	6	run
<i>commit</i>				
4	6	4	6	nil
----- <i>crash before commit</i>				
5	5	5	5	nil

Name	Guard	Effect
<i>begin</i>	$\phi = nil$	$\phi := run$

$do(a): \forall a \quad \phi = \text{run} \quad (s, \text{val}) := a(s)$   
 $l$

$commit \quad \phi = \text{run} \quad S := s, \phi := \text{nil}$   

---

 $crash \quad s := S, \phi := \text{nil}$

Note that we clean up the auxiliary state  $\phi$ .

# Incremental State Changes: Logs (1)

$S, s$  : *State*

$L, l$  : *SEQ Action* := <>

$\phi$  : {nil, run} := nil

$S = S + L$

$s, \phi = s, \phi$

$X$	$Y$	$x$	$y$	<i>Logs</i>	$\phi$
5	5	5	5		nil
<i>begin; do(x:=x-1); do(y:=y+1)</i>					
5	5	4	6	$x := 4^*$ $y := 6^*$	run
<i>commit</i>					
5	5	4	6	$x := 4^*$ $y := 6^*$	nil
-----					
<i>crash before commit</i>					
5	5	5	5		nil



<b>Name</b>	<b>Guard</b>	<b>Effect</b>
<i>begin</i>	$\phi = \text{nil}$	$\phi := \text{run}$
<i>do(a):Val</i>	$\phi = \text{run}$	$(s, \text{val}) := a(s), l += a$
<i>commit</i>	$\phi = \text{run}$	$L := l, \phi := \text{nil}$
. . .		
<i>crash</i>		$l := L, s := S+L, \phi := \text{nil}$

# Incremental State Changes: Logs (2)

$S, s$  : *State*  
 $L, l$  : *SEQ Action*  
 $\phi$  : {nil, run}

$$S = S + L$$

$$s, \phi = s, \phi$$

$X$	$Y$	$x$	$y$	<i>Logs</i>	$\phi$
5	5	4	6	$x := 4^*$ $y := 6^*$	nil
<i>apply(x := 4)</i>					
4	5	"	"	$x := 4$ $y := 6^*$	nil
<i>apply(y := 6)</i>					
4	6	"	"	$x := 4$ $y := 6$	nil
<i>cleanLog</i>					
4	6	"	"		nil
-----					
<i>crash after apply(x:=4)</i>					
4	5	"	"	$x := 4^*$ $y := 6^*$	nil

Name	Guard	Effect
<i>begin, do, and commit</i> as before		
<i>apply(a)</i>	$a = \text{head}(l)$	$S := S + a, l := \text{tail}(l)$
<i>cleanLog</i>	$L \text{ in } S$	$L := \langle \rangle$
<i>crash</i>		$l := L, s := S+L, \phi := \text{nil}$

# Incremental Log Changes

$S, s$  : State

$L, l$  : SEQ Action

$\Phi, \phi$  : {nil, run\*, **commit**}

$L = L$  if  $\phi = \text{com}$  else  $\langle \rangle$

$\phi = \phi$  if  $\phi \neq \text{com}$  else nil

$X$	$Y$	$x$	$y$	Logs	$\Phi$	$\phi$
5	5	4	6	$x := 4^*$ $y := 6^*$	nil	run
<i>flush; commit</i>						
5	5	"	"	$x := 4^*$ $y := 6^*$	com	com
<i>apply(x := 4); apply(y := 6)</i>						
4	6	"	"	$x := 4$ $y := 6$	com	com
<i>cleanLog; cleanup</i>						
4	6	"	"		nil	nil
-----						
<i>crash after flush</i> ←						
4	5	"	"	$x := 4^*$ $y := 6^*$	nil	nil

Name	Guard	Effect
<i>begin</i> and <i>do</i> as before		
<i>flush</i>	$\phi = \text{run}$	copy <b>some</b> of $l$ to $L$
<i>commit</i>	$\phi = \text{run}, L = l$	$\Phi := \phi := \text{commit}$
<i>apply(a)</i>	$\phi = \text{commit}, "$	$"$
<i>cleanLog</i>	$\text{head}(L) \text{ in } S$ or $\phi = \text{nil}$	$L := \text{tail}(L)$
<i>cleanup</i>	$L = \langle \rangle$	$\Phi := \phi := \text{nil}$
<i>crash</i>		$l := \langle \rangle \text{ if } \Phi = \text{nil} \text{ else } L;$ $s := S + l, \phi := \Phi$

# Distributed State and Log

$S_i, s_i$  : State

$L_i, l_i$  : SEQ Action

$\Phi_i, \phi_i$  : {nil, run\*, commit}

$S, L, \Phi$  are the **products** of the  $S_i, L_i, \Phi_i$

$\phi$  = run if **all**  $\phi_i = \text{run}$   
 com if **any**  $\phi_i = \text{com}$   
 and **any**  $L_i \neq \langle \rangle$   
 nil otherwise

Name	Guard	Effect
<i>begin</i> and <i>do</i> as before		
<i>flush<sub>i</sub></i>	$\phi_i = \text{run}$	copy some of $l_i$ to $L_i$
<i>prepare<sub>i</sub></i>	$\phi_i = \text{run}, L_i \neq l_i$	$\Phi_i := \text{run}$
<i>commit</i>	$\phi = \text{run}, L = l$	<b>some</b> $\Phi_i := \phi_i := \text{commit}$
<i>cleanLog</i> and <i>cleanup</i> as before		
<i>crash<sub>i</sub></i>		$l_i := \langle \rangle$ if $\Phi_i = \text{nil}$ else $L_i$ ; $s_i := S_i + l_i, \phi_i := \Phi_i$

# High Availability

**The  $\Phi = \text{commit}$  is a possible single point of failure.**

With the usual two-phase commit (2PC) this is indeed a limitation on availability.

If data is replicated, an unreplicated commit is a weakness.

**Deal with this by using a highly available *consensus* algorithm for  $\Phi$ .**

Lamport's Paxos algorithm is the best currently known.

# Transactions: Summary

## Ideas

Logs

Commit records

Stable writes at critical points: prepare and commit

Lazy cleanup

**The spec is simple.**

**Implementations are subtle because of crashes.**

The abstraction functions reveal their secrets.

The subtlety can be added one step at a time.



# How to Write a Spec

## Figure out what the state is

Choose it to make the spec clear, not to match the code.

## Describe the actions

What they do to the state

What they return

## Helpful hints

Notation is important; it helps you to think about what's going on.

Invent a suitable vocabulary.

Fewer actions are better.

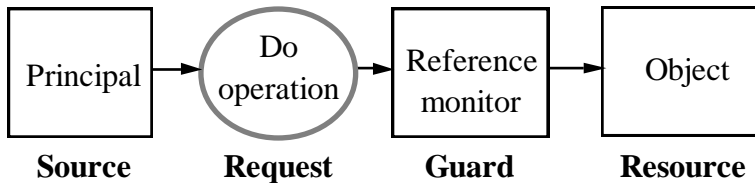
*Less is more.*

More non-determinism is better; it allows more implementations.

*I'm sorry I wrote you such a long letter; I didn't have time to write a short one.*  
*(Pascal)*

# Security: The Access Control Model

Guards control access to valued resources.



Rules control the operations allowed for each principal and object.

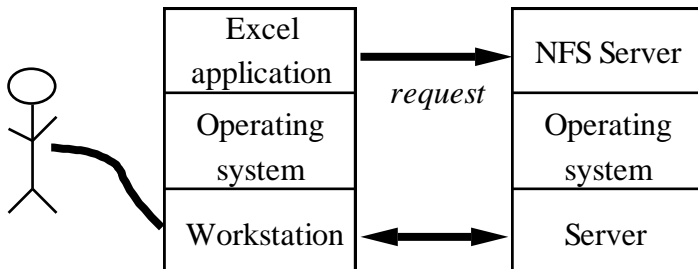
<i>Principal may do</i>	<i>Operation on</i>	<i>Object</i>
Taylor	Read	File "Raises"
Jones	Pay invoice 4325	Account Q34

Schwarzkopf

Fire three rounds

Bow gun

# A Distributed System



# Principals

**Authentication:**      **Who sent a message?**

**Authorization:**      **Who is trusted?**

**Principal — abstraction of "who":**

People	Lampson, Taylor
Machines	VaxSN12648, Jumbo
Services	SRC-NFS, X-server
Groups	SRC, DEC-Employees
Channels	Key #7438

# Theory of Principals

**Principal says statement**

$P \text{ says } s$

Lampson says “read /SRC/Lampson/foo”

SRC-CA says “Lampson’s key is #7438”

**Principal  $A$  speaks for  $B$**

$A \Rightarrow B$

If  $A$  says something,  $B$  says it too. So  $A$  is stronger than  $B$ .

**A secure channel:**

says things directly

$C \text{ says } s$

If  $P$  is the only sender on  $C$

$C \Rightarrow P$

**Examples**

Lampson  $\Rightarrow$  SRC

Key #7438  $\Rightarrow$  Lampson

# Handing Off Authority

**Handoff rule:** If  $A$  says  $B \Rightarrow A$  then  $B \Rightarrow A$

Reasonable if  $A$  is competent and accessible.

## Examples:

SRC **says** Lampson  $\Rightarrow$  SRC

Node key **says** Channel key  $\Rightarrow$  Node key

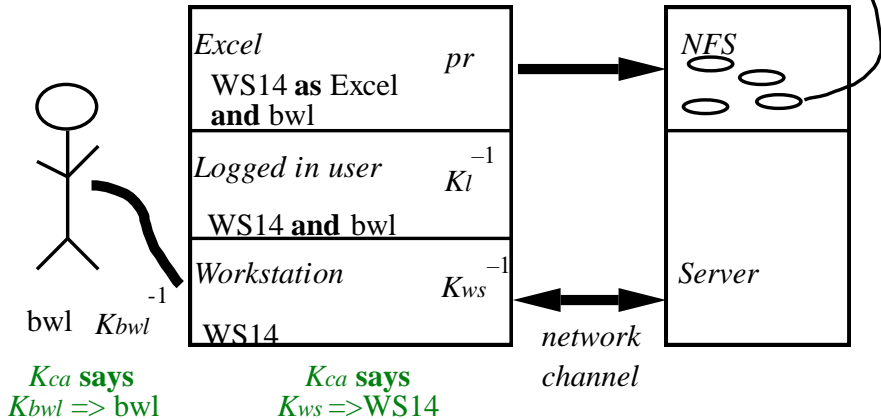
*Any problem in computer science can be solved  
with another level of indirection. (Wheeler).*

# Authenticating to the Server

(SRC-node as Excel) and bwl may read

SRC says WS14  $\Rightarrow$  SRC-node

file foo





# Access Control

## Checking access:

Given a request  $Q$  says read  $O$   
an ACL  $P$  may read  $O$

Check that  $Q$  speaks for  $P$   $Q \Rightarrow P$

## Auditing

Each step is justified by  
a signed statement, or  
a rule

# Authenticating a Channel

**Authentication** — who can send on a channel.

$C \Rightarrow P$ ;  $C$  is the channel,  $P$  the sender.

**To get new  $C \Rightarrow P$  facts**, must trust some principal, a *certification authority*, to tell them to you.

Simplest: trust  $K_{ca}$  to authenticate any name:

$K_{ca} \Rightarrow \text{Anybody}$
-------------------------------------

**Then  $CA$  can authenticate channels:**

$K_{ca}$  **says**  $K_{ws}$   $\Rightarrow$   $WS$

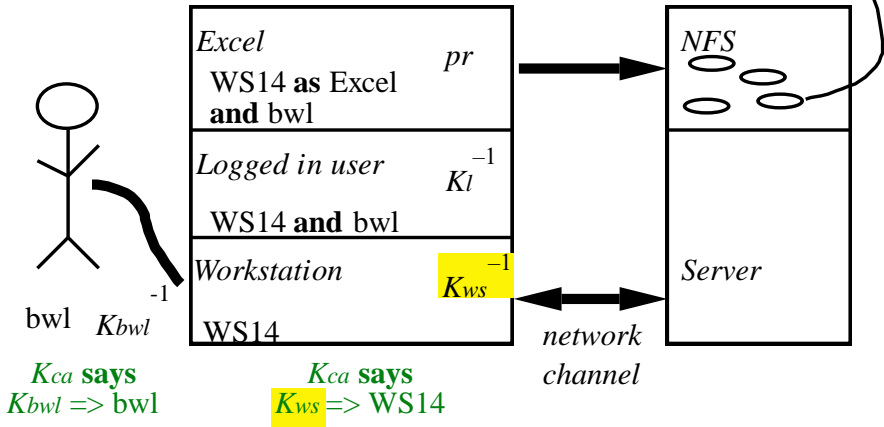
$K_{ca}$  **says**  $K_{bwl}$   $\Rightarrow$   $bwl$

# Authenticated Channels: Example

(SRC-node as Excel) and bwl  
may read

SRC says WS14 => SRC-node

file foo



# Groups and Group Credentials

**Defining groups:** A group is a principal; its members speak for it.

Lampson  $\Rightarrow$  SRC

Taylor  $\Rightarrow$  SRC

...

**Proving group membership:** Use certificates.

$K_{src}$  says Lampson  $\Rightarrow$  SRC

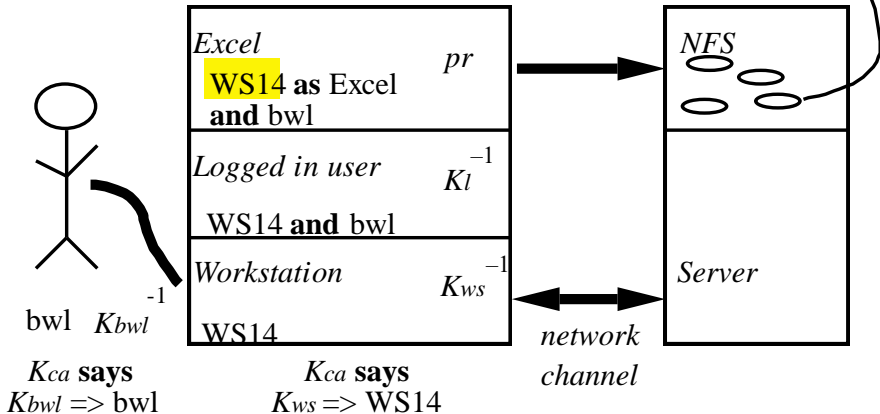
$K_{ca}$  says  $K_{src}$   $\Rightarrow$  SRC

# Authenticating a Group

(SRC-node as Excel) and bwl  
may read

SRC says WS14 => SRC-node

file foo



# Security: Summary

## Ideas

Principals

Channels as principals

“Speaks for” relation

Handoff of authority

**Give precise rules.**

**Apply them to cover many cases.**

# References

- Hints*                      Lampson, Hints for Computer System Design. *IEEE Software*, Jan. 1984.
- Specifications*        Lamport, A simple approach to specifying concurrent systems. *Communications of the ACM*, Jan. 1989.
- Reliable messages*    in Mullender, ed., *Distributed Systems*, Addison-Wesley, 1993 (summer)
- Transactions*         Gray and Reuter, *Transaction Processing: Concepts and Techniques*. Morgan Kaufman, 1993.
- Security*                Lampson, Abadi, Burrows, and Wobber, Authentication in distributed systems: Theory and

practice. *ACM Transactions on Computer Systems*,  
Nov. 1992.



# Collaborators

Charles Simonyi

Bravo: WYSIWYG editor

Bob Sproull

Alto operating system

Dover: laser printer

Interpress: page description language

Mel Pirtle

940 project, Berkeley Computer Corp.

Peter Deutsch

940 operating system

QSPL: system programming language

Chuck Geschke

Mesa: system programming language

Jim Mitchell

Ed Satterthwaite

Jim Horning

Euclid: verifiable programming language

Ron Rider

Ears: laser printer

Gary Starkweather

Severo Ornstein

Dover: laser printer

# Collaborators

Roy Levin

Wildflower: Star workstation prototype  
Vesta: software configuration

Andrew Birrell, Roger Needham, Mike Schroeder

Global name service and authentication

Eric Schmidt

System models: software configuration

Rod Burstall

Pebble: polymorphic typed language