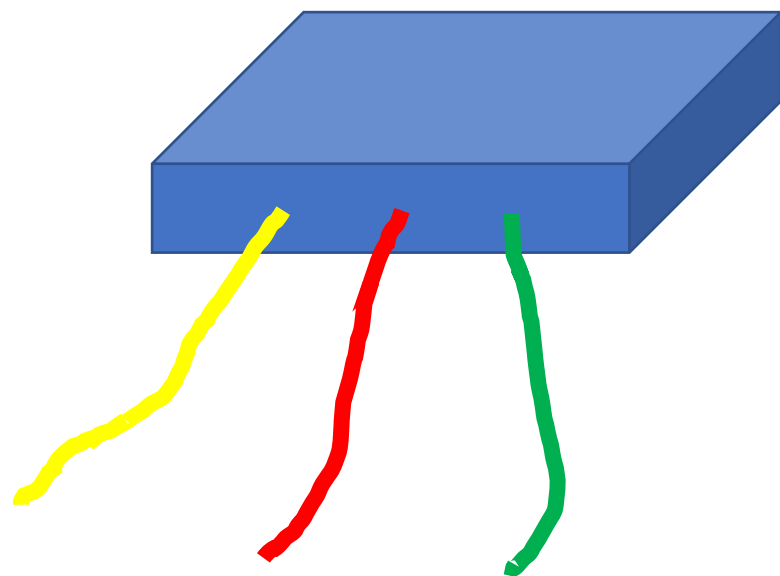


# Lecture 4: Refinement

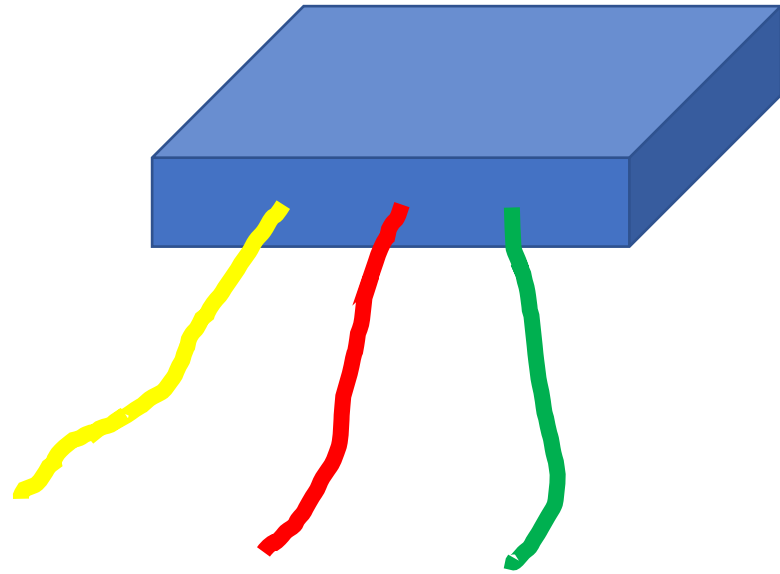
Based on material from Section 10.8, Specifying Systems by Leslie Lamport

You ask for:



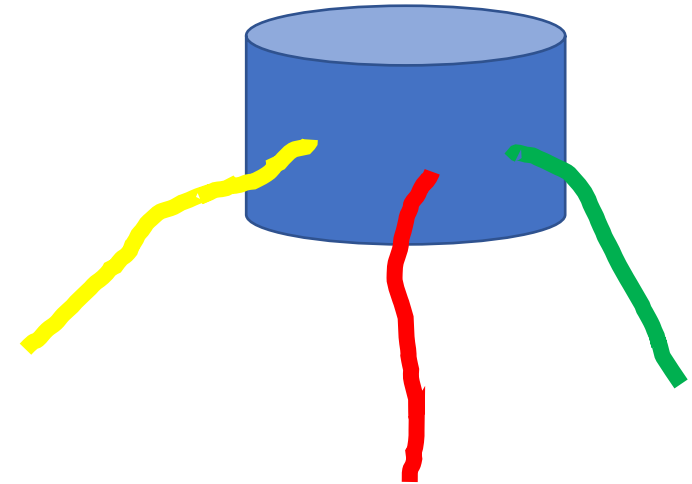
Specification

You ask for:



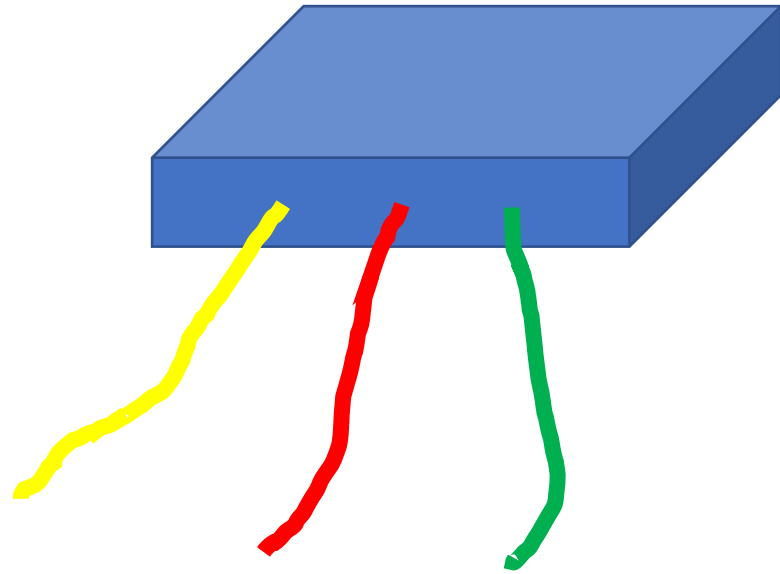
Specification

You get:



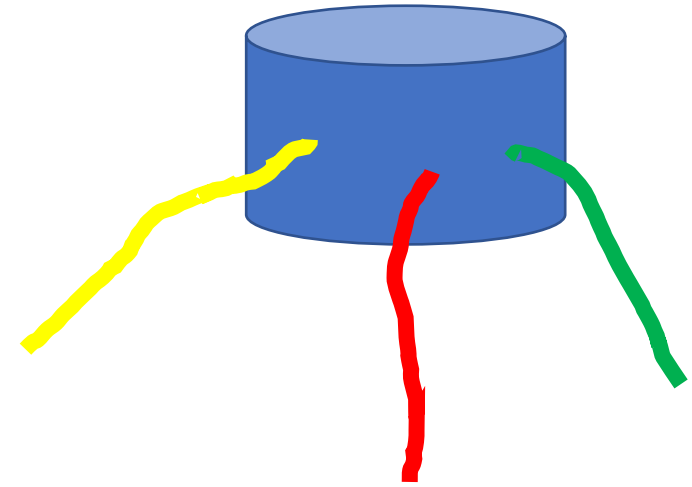
Implementation

You ask for:



Specification

You get:

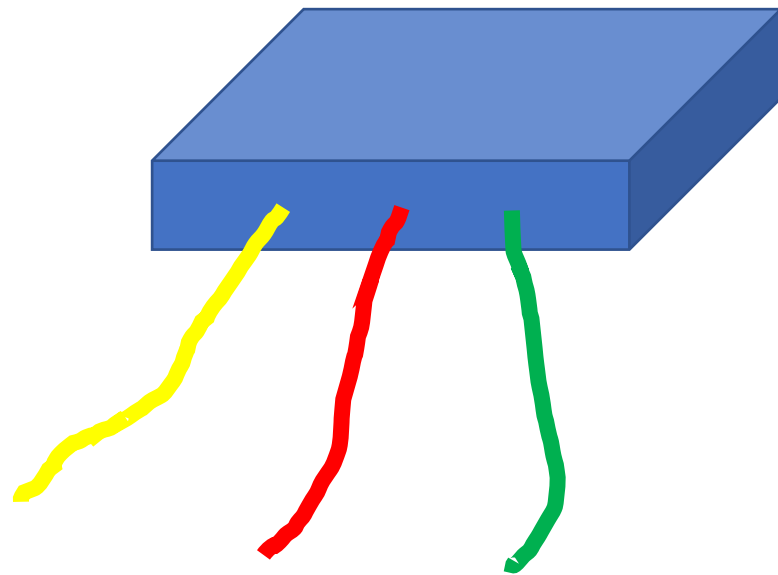


Implementation

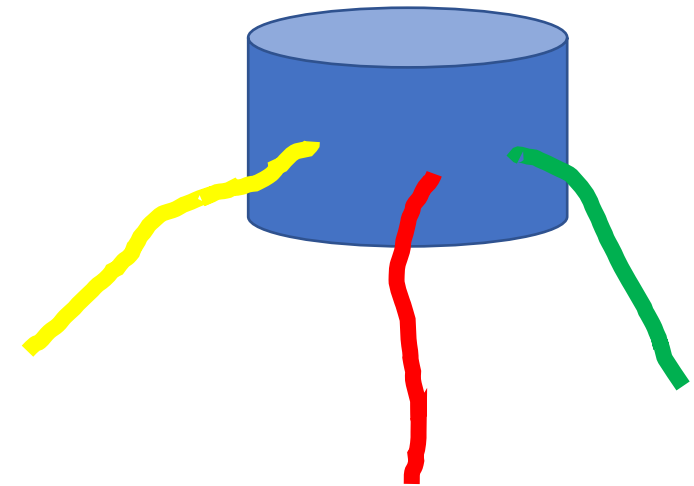
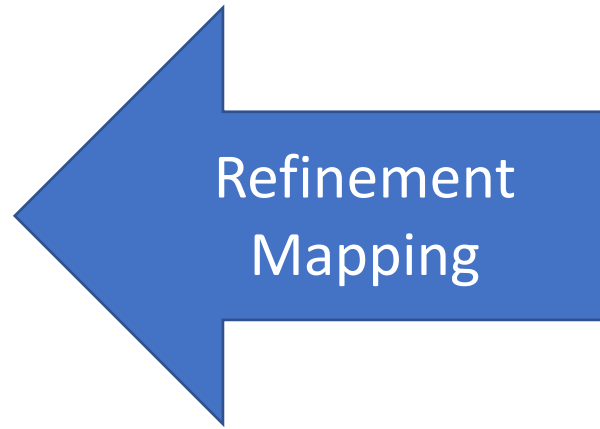
*Is every behavior of the implementation also a behavior of the specification?*

You ask for:

You get:



Specification

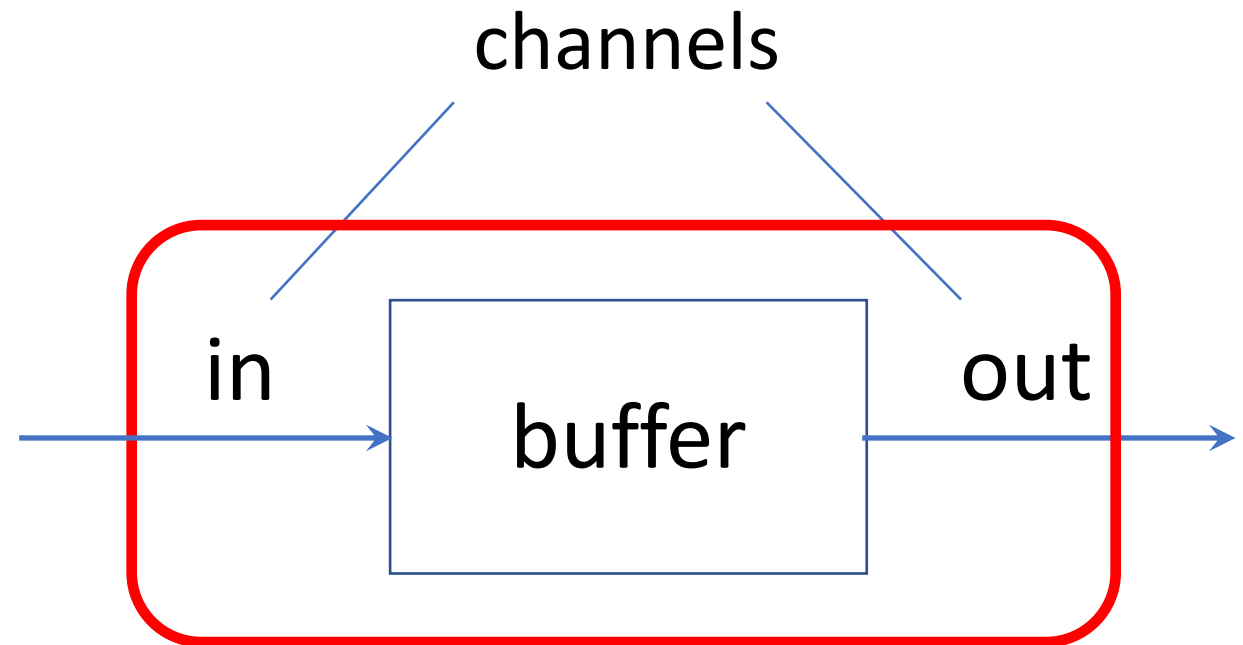


Implementation

*Is every behavior of the implementation also a behavior of the specification?*

# External/internal variables of a state

- A specification has certain *external variables* that can be observed and/or manipulated
- It may also have *internal variables* that are used to describe behaviors but that cannot be observed
- Example: FIFO
  - External variables: in, out
  - Internal variable: buffer



# *Externally visible vs complete* behavior

A system may exhibit externally visible behavior

$$e_1 \rightarrow e_2 \rightarrow e_3 \rightarrow e_4 \rightarrow \dots$$

if there exists a complete behavior

$$(e_1, y_1) \rightarrow (e_2, y_2) \rightarrow (e_3, y_3) \rightarrow (e_4, y_4) \rightarrow$$

that is allowed by the specification

Here  $e_i$  is some externally visible state (for example, in and out channels) and  $y_i$  is internal state (for example, the buffer)

# Stuttering Steps

A specification should allow changes to the internal state that does not change the externally visible state.

For example:

$$(e_1, y_1) \rightarrow (e_2, y_2) \rightarrow (e_2, y_2') \rightarrow (e_3, y_3) \rightarrow (e_4, y_4) \rightarrow$$

leads to external behavior

$$e_1 \rightarrow e_2 \rightarrow e_2 \rightarrow e_3 \rightarrow e_4 \rightarrow \dots$$

which should be identical to

$$e_1 \rightarrow e_2 \rightarrow e_3 \rightarrow e_4 \rightarrow \dots$$



# Proving that an implementation meets the specification

- First note that an implementation is just a specification
- We call the implementation the “lower-level” specification

We need to prove that if an implementation allows the complete behavior

$$(e_1, z_1) \rightarrow (e_2, z_2) \rightarrow (e_3, z_3) \rightarrow (e_4, z_4) \rightarrow$$

then there exists a complete behavior

$$(e_1, y_1) \rightarrow (e_2, y_2) \rightarrow (e_3, y_3) \rightarrow (e_4, y_4) \rightarrow$$

allowed by the specification

A mapping from low-level complete behaviors to high-level complete behaviors is called a “**refinement mapping**”

Note, there may be multiple possible refinement mappings---you only need to show one

# Recall: Module HourClock

MODULE *HourClock*

EXTENDS *Naturals*

VARIABLE *hr*

$HCini \triangleq hr \in (1 .. 12)$

$HCnext \triangleq hr' = \text{IF } hr \neq 12 \text{ THEN } hr + 1 \text{ ELSE } 1$

$HC \triangleq HCini \wedge \square[HCnext]_{hr}$

THEOREM  $HC \Rightarrow \square HCini$

# Implementation

- Suppose we wanted to replace *hr* by a 4-bit binary value
- We need a way to represent *n*-bit binary values
- We also need a function of *n*-bit binary values to numbers

# Functions in TLA+

- A function  $f$  has a domain, written  $\text{DOMAIN } f$
- $f$  assigns to each  $x \in \text{DOMAIN } f$  a value  $f[x]$ 
  - TLA+ uses array notation (square brackets) rather than parentheses
- $f \equiv g$  iff
$$\text{DOMAIN } f = \text{DOMAIN } g \wedge \forall x \in \text{DOMAIN } f: f[x] = g[x]$$
- The range of  $f$  is  $\{ f[x] \mid x \in \text{DOMAIN } f \}$
- $[S \rightarrow T]$  is defined to be the set of functions whose domain is  $S$  and whose range is a subset of  $T$

# Function values

- $[x \in S \mapsto e]$  is defined to be the function  $f$  with domain  $S$  such that
$$\forall x \in S: f[x] = e \quad // \text{ } x \text{ is a free variable in } e$$
- For example
  - $succ \triangleq [n \in 1..12 \mapsto \text{IF } n = 12 \text{ THEN } 1 \text{ ELSE } n + 1]$
  - $prod \triangleq [x \in Real, y \in Real \mapsto x * y]$
  - $double \triangleq [x \in Real \mapsto prod[x][2]]$
- Similar to lambda expressions

# Cool: Records are functions

[ *val*  $\mapsto$  42, *rdy*  $\mapsto$  1, *ack*  $\mapsto$  0 ] is equivalent to

[  $x \in \{ \text{"val"}, \text{"rdy"}, \text{"ack"} \} \mapsto$

IF  $x = \text{"val"}$  THEN 42

ELSE IF  $x = \text{"rdy"}$  THEN 1

ELSE 0 // must be "ack" due to DOMAIN

]

# Choose Operator

CHOOSE  $x$ :  $F$

expression that evaluates to some (possibly unspecific) value  $x$   
that satisfies  $F$  //  $x$  is a free variable in  $F$

CHOOSE  $x \in S$ :  $F \triangleq$  CHOOSE  $x$ :  $x \in S \wedge F$

Undefined if no such  $x$  exists

Example:  $\max(S) \triangleq$  CHOOSE  $x \in S$ :  $\forall y \in S: x \geq y$

the maximum element of  $S$  // undefined if  $S$  is empty

# Choose Operator, cont'd

CHOOSE  $x: F$  always evaluates to the same value. That is,

$$F \equiv G \Rightarrow (\text{CHOOSE } x: F) = (\text{CHOOSE } x: G)$$

Also

$$(v = \text{CHOOSE } x: F) \wedge (w = \text{CHOOSE } x: F) \Rightarrow v = w$$

However, the value of  $v$  is unspecified

Q1: what behaviors are allowed by

$$(x = \text{CHOOSE } n: n \in \text{Nat}) \wedge \Box[x' = \text{CHOOSE } n: n \in \text{Nat}]_x ?$$



# Choose Operator, cont'd

CHOOSE  $x$ :  $F$  always evaluates to the same value. That is,

$$F \equiv G \Rightarrow (\text{CHOOSE } x: F) = (\text{CHOOSE } x: G)$$

Also

$$(v = \text{CHOOSE } x: F) \wedge (w = \text{CHOOSE } x: F) \Rightarrow v = w$$

However, the value of  $v$  is unspecified

Q1: what behaviors are allowed by

$$(x = \text{CHOOSE } n: n \in \text{Nat}) \wedge \Box[x' = \text{CHOOSE } n: n \in \text{Nat}]_x ?$$

**Answer:  $x$  is always the same (but unspecified) natural number**

# Choose Operator, cont'd

CHOOSE  $x: F$  always evaluates to the same value. That is,

$$F \equiv G \Rightarrow (\text{CHOOSE } x: F) = (\text{CHOOSE } x: G)$$

Also

$$(v = \text{CHOOSE } x: F) \wedge (w = \text{CHOOSE } x: F) \Rightarrow v = w$$

However, the value of  $v$  is unspecified

Q1: what behaviors are allowed by

$$(x = \text{CHOOSE } n: n \in \text{Nat}) \wedge \Box[x' = \text{CHOOSE } n: n \in \text{Nat}]_x ?$$

**Answer:  $x$  is always the same (but unspecified) natural number**

Q2: what behaviors are allowed by

$$(x \in \text{Nat}) \wedge \Box[x' \in \text{Nat}]_x ?$$

# Choose Operator, cont'd

CHOOSE  $x: F$  always evaluates to the same value. That is,

$$F \equiv G \Rightarrow (\text{CHOOSE } x: F) = (\text{CHOOSE } x: G)$$

Also

$$(v = \text{CHOOSE } x: F) \wedge (w = \text{CHOOSE } x: F) \Rightarrow v = w$$

However, the value of  $v$  is unspecified

Q1: what behaviors are allowed by

$$(x = \text{CHOOSE } n: n \in \text{Nat}) \wedge \Box[x' = \text{CHOOSE } n: n \in \text{Nat}]_x ?$$

**Answer:  $x$  is always the same (but unspecified) natural number**

Q2: what behaviors are allowed by

$$(x \in \text{Nat}) \wedge \Box[x' \in \text{Nat}]_x ?$$

**Answer:  $x$  can be a different natural number in every state**

# Recursive functions

- $fact \triangleq [n \in Nat \mapsto \text{IF } n = 0 \text{ THEN } 1 \text{ ELSE } n * fact[n - 1]]$

is illegal because  $fact$  is not defined in the expression on the right

Instead:

- $fact \triangleq \text{CHOOSE } f: [n \in Nat \mapsto \text{IF } n = 0 \text{ THEN } 1 \text{ ELSE } n * f[n - 1]]$

Shorthand:

- $fact[n \in Nat] \triangleq \text{IF } x = 0 \text{ THEN } 1 \text{ ELSE } n * fact[n - 1]$

# Aside: Scoping Definitions in TLA+

LET

$d_1 \triangleq e_1$

$d_2 \triangleq e_2$

...

IN

...

- Split complicated formulas into smaller chunks
- Leverage common subexpressions

# Representing an $n$ -bit value

- We can represent an  $n$ -bit value by a function  
$$b \triangleq [x \in 0..(n-1) \mapsto 0..1]$$
- For example, if  $b$  represents 0101 (i.e., 5) then
  - $b[0] = 1$
  - $b[1] = 0$
  - $b[2] = 1$
  - $b[3] = 0$

and  $b$  corresponds to the number

$$b[0] * 2^0 + b[1] * 2^1 + b[2] * 2^2 + b[3] * 2^3$$

Finally: function of  $n$ -bit value  $b$  to number

$BitArrayVal(b) \triangleq$

LET

$n \triangleq$  CHOOSE  $m \in Nat$ : DOMAIN  $b = 0..(m - 1)$

$f[x \in 0..(n - 1)] \triangleq$

IF  $x = 0$  THEN  $b[0]$  ELSE  $b[x] * 2^x + f[x - 1]$

IN

$f[n - 1]$

# BinaryHourClock: (broken) attempt 1

MODULE *BinaryHourClock*

EXTENDS *Naturals*

VARIABLE *bits*

$BitArrayVal(b) \triangleq$

LET

$n \triangleq$  CHOOSE  $m \in Nat : DOMAIN\ b = 0 .. (m - 1)$

$f[x \in 0 .. (n - 1)] \triangleq$  IF  $x = 0$  THEN  $b[0]$  ELSE  $b[x] * 2^x + f[x - 1]$

IN  $f[n - 1]$

$B \triangleq$  INSTANCE *HourClock* WITH  $hr \leftarrow BitArrayVal(bits)$

Substitute  
*BitArrayVal(bits)* for *hr*

$Spec \triangleq B!HC$



# What's the (subtle) issue?

- $BitArrayVal(b)$  is undefined unless  $b$  is a function  $b$  with domain  $0..n - 1$  for some  $n$
- $BitArrayVal(\text{"Fred"})$  is undefined
- Perhaps  $BitArrayVal(\text{"Fred"}) = 7$ . If so, "Fred" would be an allowed initial value of  $bits$ . Probably not what we intended to specify

Fix:  $HourVal(b) \triangleq \text{IF } b \in [(0..3) \rightarrow (0..1)] \text{ THEN } BitArrayVal(b) \text{ ELSE } 13$

$B \triangleq \text{INSTANCE } HourClock \text{ WITH } hr \leftarrow HourVal(bits)$

$Spec \triangleq B!HC$

Any value  
but 1..12

Because  $HC$  is never satisfied by a state in which  $hr = 13$ ,  $bits$  has to be in  $[0..3 \rightarrow 0..1]$

# A little more elegant solution

MODULE *BinaryHourClock*

EXTENDS *Naturals*

VARIABLE *bits*

$BitArrayVal(b) \triangleq$

LET

$n \triangleq \text{CHOOSE } m \in Nat : \text{DOMAIN } b = 0 .. (m - 1)$

$f[x \in 0 .. (n - 1)] \triangleq \text{IF } x = 0 \text{ THEN } b[0] \text{ ELSE } b[x] * 2^x + f[x - 1]$

IN  $f[n - 1]$

$ErrorVal \triangleq \text{CHOOSE } v : v \notin 1 .. 12$

$HourVal(b) \triangleq \text{IF } b \in [(0 .. 3) \rightarrow (0 .. 1)] \text{ THEN } BitArrayVal(b) \text{ ELSE } ErrorVal$

$B \triangleq \text{INSTANCE } HourClock \text{ WITH } hr \leftarrow HourVal(bits)$

$Spec \triangleq B!HC$

# A better way of doing it (instead of substitution)

MODULE *BinaryHourClock*

EXTENDS *Naturals*

VARIABLE *bits*

$H(hr) \triangleq$  INSTANCE *HourClock*

$BitArrayVal(b) \triangleq$

LET

$n \triangleq$  CHOOSE  $m \in Nat : DOMAIN\ b = 0 .. (m - 1)$

$f[x \in 0 .. (n - 1)] \triangleq$  IF  $x = 0$  THEN  $b[0]$  ELSE  $b[x] * 2^x + f[x - 1]$

IN  $f[n - 1]$

$ErrorVal \triangleq$  CHOOSE  $v : v \notin 1 .. 12$

$HourVal(b) \triangleq$  IF  $b \in [(0 .. 3) \rightarrow (0 .. 1)]$  THEN  $BitArrayVal(b)$  ELSE  $ErrorVal$

$B \triangleq$  INSTANCE *HourClock* WITH  $hr \leftarrow HourVal(bits)$

$IR(h) \triangleq$   $\square(h = HourVal(bits))$

$Spec \triangleq$   $\exists hr : IR(hr) \wedge H(hr) ! HC$

*bits* and *hr* keep the same time  
(IR stands for Interface Refinement)

Hide *hr*

# Discussion

- Here we *composed* two specifications:
  - An HourClock with values from the set 1..12
  - A BinaryHourClock with values from the set  $[(0..3) \rightarrow (0..1)]$
- *Composition is a conjunction of specifications*
  - A behavior of the composition is a behavior of each of the components
- $IR(hr)$  asserts that *bits* is always the 4-bit value representing *hr*
- $IR(hr) \wedge H(hr)! HC$  asserts that *hr* and *bits* keep the same time

More precisely:

- $IR(hr)$  defines *hr* as a function of *bits*, but does not constrain *bits*
- $IR(hr) \wedge H(hr)! HC$  constrains behaviors involving *bits* by requiring that they have to map to behaviors involving *hr*

# Interface Refinement

- BinaryHourClock is an *implementation* of HourClock
- One has to exhibit a mapping from the "low-level" implementation to the "high-level" specification
  - Map the low-level state to the high-level state:  $hr = HourVal(bits)$
  - Map each low-level step to a high-level step or to a high-level stuttering step
    - In the case of the BinaryHourClock, the low-level and high-level steps are the same
    - If so, that is a special case of interface refinement called "*Data Refinement*"
  - (leaving out liveness for now)
- Each behavior of the implementation is also a behavior of the specification

# We already saw an example of data refinement

MODULE *AsynchInterface*

EXTENDS *Naturals*

CONSTANT *Data*

VARIABLES *val, rdy, ack*

*TypeInvariant*  $\triangleq$   $\wedge val \in Data$   
 $\wedge rdy \in \{0, 1\}$   
 $\wedge ack \in \{0, 1\}$

*Init*  $\triangleq$   $\wedge val \in Data$   
 $\wedge rdy \in \{0, 1\}$   
 $\wedge ack = rdy$

*Send*  $\triangleq$   $\wedge rdy = ack$   
 $\wedge val' \in Data$   
 $\wedge rdy' = 1 - rdy$   
 $\wedge UNCHANGED\ ack$

*Rcv*  $\triangleq$   $\wedge rdy \neq ack$   
 $\wedge ack' = 1 - ack$   
 $\wedge UNCHANGED\ \langle val, rdy \rangle$

*Next*  $\triangleq Send \vee Rcv$

*Spec*  $\triangleq Init \wedge \square[Next]_{\langle val, rdy, ack \rangle}$

THEOREM  $Spec \Rightarrow \square TypeInvariant$

MODULE *Channel*

EXTENDS *Naturals*

CONSTANT *Data*

VARIABLE *chan*

*TypeInvariant*  $\triangleq chan \in [val : Data, rdy : \{0, 1\}, ack : \{0, 1\}]$

*Init*  $\triangleq$   $\wedge TypeInvariant$   
 $\wedge chan.ack = chan.rdy$

*Send*(*d*)  $\triangleq$   $\wedge chan.rdy = chan.ack$   
 $\wedge chan' = [chan\ EXCEPT\ !.val = d, !.rdy = 1 - @]$

*Rcv*  $\triangleq$   $\wedge chan.rdy \neq chan.ack$   
 $\wedge chan' = [chan\ EXCEPT\ !.ack = 1 - @]$

*Next*  $\triangleq (\exists d \in Data : Send(d)) \vee Rcv$

*Spec*  $\triangleq Init \wedge \square[Next]_{chan}$

THEOREM  $Spec \Rightarrow \square TypeInvariant$







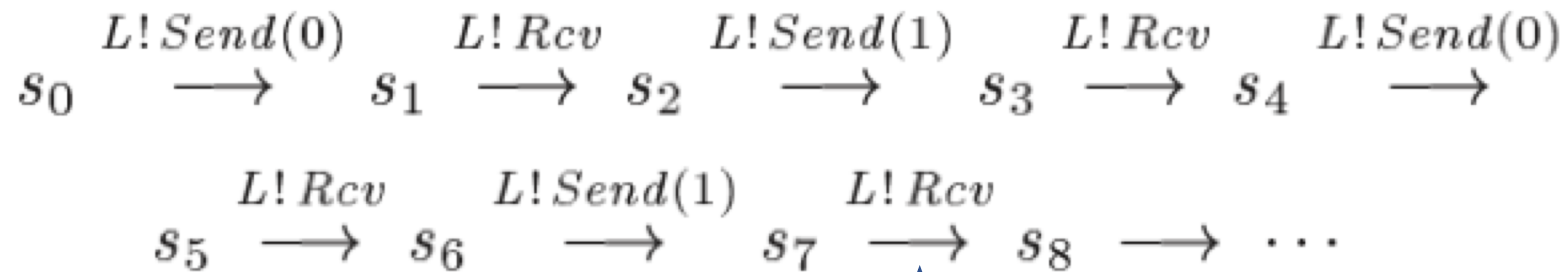


# Channel Refinement

$H \triangleq$  INSTANCE *Channel* WITH  $chan \leftarrow h, Data \leftarrow 1 .. 12$

$L \triangleq$  INSTANCE *Channel* WITH  $chan \leftarrow l, Data \leftarrow \{0, 1\}$

Sending 5 (= 0101):



Corresponds to  
 $H!Send(5)$

Corresponds to  
 $H!Rcv$

# Interface Refinement

Recall definition of  $IR$  for BinaryHourClock a few slides ago:

- $IR$  will specify  $h$  as a function of  $l$ , but does not constrain  $l$

Then, if  $HSpec$  is a high-level spec of the system, we can write the low-level spec as

$$\exists h : IR \wedge HSpec$$

EXTENDS *Naturals, Sequences*

VARIABLES  $h, l$

$ErrorVal \triangleq$  CHOOSE  $v : v \notin [val : 1 .. 12, rdy : \{0, 1\}, ack : \{0, 1\}]$

$BitSeqToNat[s \in Seq(\{0, 1\})] \triangleq$   $BitSeqToNat[\langle b_0, b_1, b_2, b_3 \rangle] = b_0 + 2 * (b_1 + 2 * (b_2 + 2 * b_3))$   
 IF  $s = \langle \rangle$  THEN 0 ELSE  $Head(s) + 2 * BitSeqToNat[Tail(s)]$

$H \triangleq$  INSTANCE *Channel* WITH  $chan \leftarrow h, Data \leftarrow 1 .. 12$

$L \triangleq$  INSTANCE *Channel* WITH  $chan \leftarrow l, Data \leftarrow \{0, 1\}$

$H$  is a channel for sending numbers in  $1 .. 12$ ;  $L$  is a channel for sending bits.



MODULE *ChannelRefinement*

MODULE *Inner*

VARIABLE *bitsSent* The sequence of the bits sent so far for the current number.

$InnerIR \triangleq Init \wedge \square[Next]_{\langle l, h, bitsSent \rangle}$

on previous page

$I(bitsSent) \triangleq \text{INSTANCE } Inner$

$IR \triangleq \exists bitsSent : I(bitsSent)!InnerIR$

$h$  is a function of  $l$

MODULE *LowerSpec*

VARIABLE *lchan*

CONSTANT *Data*

$HS(hchan) \triangleq \text{INSTANCE } Channel \text{ WITH } chan \leftarrow hchan, Data \leftarrow 1 .. 12$

$CR(h) \triangleq \text{INSTANCE } ChannelRefinement \text{ WITH } l \leftarrow lchan$

$LSpec \triangleq \exists h : CR(h)!IR \wedge HS(h)!HSspec$

constrains *lchan*

# Other examples of refinement

# Class Queue

```
class {:autocontracts} Queue {  
  ghost var Contents: seq<int>;  
  var a: array<int>;  
  var hd: int, tl: int;  
  
  predicate Valid() {           // class invariant  
    a.Length > 0 && 0 <= tl <= hd <= a.Length && Contents == a[tl..hd]  
  }  
  
  constructor () ensures Contents == []  
  {  
    a, tl, hd, Contents := new int[10], 0, 0, [];  
  }  
}
```

# Class Queue: continued

```
method Enqueue(d: int) ensures Contents == old(Contents) + [d] {  
  if hd == a.Length {  
    var b := a;  
    if tl == 0 { b := new int[2 * a.Length]; }           // a is full  
    forall (i | 0 <= i < hd - tl) { b[i] := a[tl + i]; } // shift  
    a, tl, hd := b, 0, hd - tl;  
  }  
  a[hd], hd, Contents := d, hd + 1, Contents + [d];  
}
```

```
method Dequeue() returns (d: int)  
  requires Contents != []  
  ensures d == old(Contents)[0] && Contents == old(Contents)[1..];  
{  
  d, tl, Contents := a[tl], tl + 1, Contents[1..];  
}
```



## Chain Replication for Supporting High Throughput and Availability

Robbert van Renesse  
rvr@cs.cornell.edu

Fred B. Schneider  
fbs@cs.cornell.edu

*FAST Search & Transfer ASA  
Tromsø, Norway  
and  
Department of Computer Science  
Cornell University  
Ithaca, New York 14853*

### Abstract

Chain replication is a new approach to coordinating clusters of fail-stop storage servers. The approach is intended for supporting large-scale storage services that exhibit high throughput and availability without sacrificing strong consistency guarantees. Besides outlining the chain replication protocols themselves, simulation experiments explore the performance characteristics of a prototype implementation. Throughput, availability, and several object-placement strategies (including schemes based on distributed hash table routing) are discussed.

### 1 Introduction

A *storage system* typically implements operations so that clients can store, retrieve, and/or change data. File systems and database systems are perhaps the best known examples. With a file system, operations (read and write) access a single file and are idempotent; with a database system, operations (transactions) may each access multiple objects and are serializable.

This paper is concerned with storage systems that sit somewhere between file systems and database systems. In particular, we are concerned with storage systems, henceforth called *storage services*, that

- store *objects* (of an unspecified nature),
- support *query* operations to return a value derived from a single object, and
- support *update* operations to atomically change the state of a single object according to some

pre-programmed, possibly non-deterministic, computation involving the prior state of that object.

A file system write is thus a special case of our storage service update which, in turn, is a special case of a database transaction.

Increasingly, we see on-line vendors (like Amazon.com), search engines (like Google's and FAST's), and a host of other information-intensive services provide value by connecting large-scale storage systems to networks. A storage service is the appropriate compromise for such applications, when a database system would be too expensive and a file system lacks rich enough semantics.

One challenge when building a large-scale storage service is maintaining high availability and high throughput despite failures and concomitant changes to the storage service's configuration, as faulty components are detected and replaced.

Consistency guarantees also can be crucial. But even when they are not, the construction of an application that fronts a storage service is often simplified given *strong consistency guarantees*, which assert that (i) operations to query and update individual objects are executed in some sequential order and (ii) the effects of update operations are necessarily reflected in results returned by subsequent query operations.

Strong consistency guarantees are often thought to be in tension with achieving high throughput and high availability. So system designers, reluctant to sacrifice system throughput or availability, regularly decline to support strong consistency guarantees. The Google File System (GFS) illustrates this thinking [11]. In fact, strong consistency guarantees

### State is:

$Hist_{objID}$  : update request sequence

$Pending_{objID}$  : request set

### Transitions are:

T1: Client request  $r$  arrives:

$$Pending_{objID} := Pending_{objID} \cup \{r\}$$

T2: Client request  $r \in Pending_{objID}$  ignored:

$$Pending_{objID} := Pending_{objID} - \{r\}$$

T3: Client request  $r \in Pending_{objID}$  processed:

$$Pending_{objID} := Pending_{objID} - \{r\}$$

if  $r = \text{query}(objId, opts)$  then

reply according options  $opts$  based  
on  $Hist_{objID}$

else if  $r = \text{update}(objId, newVal, opts)$  then

$$Hist_{objID} := Hist_{objID} \cdot r$$

reply according options  $opts$  based  
on  $Hist_{objID}$

Figure 1: Client's View of an Object.

### 3.1 Protocol Details

Clients do not directly read or write variables  $Hist_{objID}$  and  $Pending_{objID}$  of Figure 1, so we are free to implement them in any way that is convenient. When chain replication is used to implement the specification of Figure 1:

- $Hist_{objID}$  is defined to be  $Hist_{objID}^T$ , the value of  $Hist_{objID}$  stored by tail  $T$  of the chain, and
- $Pending_{objID}$  is defined to be the set of client requests received by any server in the chain and not yet processed by the tail.

The chain replication protocols for query processing and update processing are then shown to satisfy the specification of Figure 1 by demonstrating how each state transition made by any server in the chain is equivalent either to a no-op or to allowed transitions T1, T2, or T3.

**State is:**

$Hist_{objID}$  : update request sequence

$Pending_{objID}$  : request set

**Transitions are:**

T1: Client request  $r$  arrives:

$Pending_{objID} := Pending_{objID} \cup \{r\}$

T2: Client request  $r \in Pending_{objID}$  ignored:

$Pending_{objID} := Pending_{objID} - \{r\}$

T3: Client request  $r \in Pending_{objID}$  processed:

$Pending_{objID} := Pending_{objID} - \{r\}$

**if**  $r = \text{query}(objId, opts)$  **then**

**reply** according options  $opts$  based  
on  $Hist_{objID}$

**else if**  $r = \text{update}(objId, newVal, opts)$  **then**

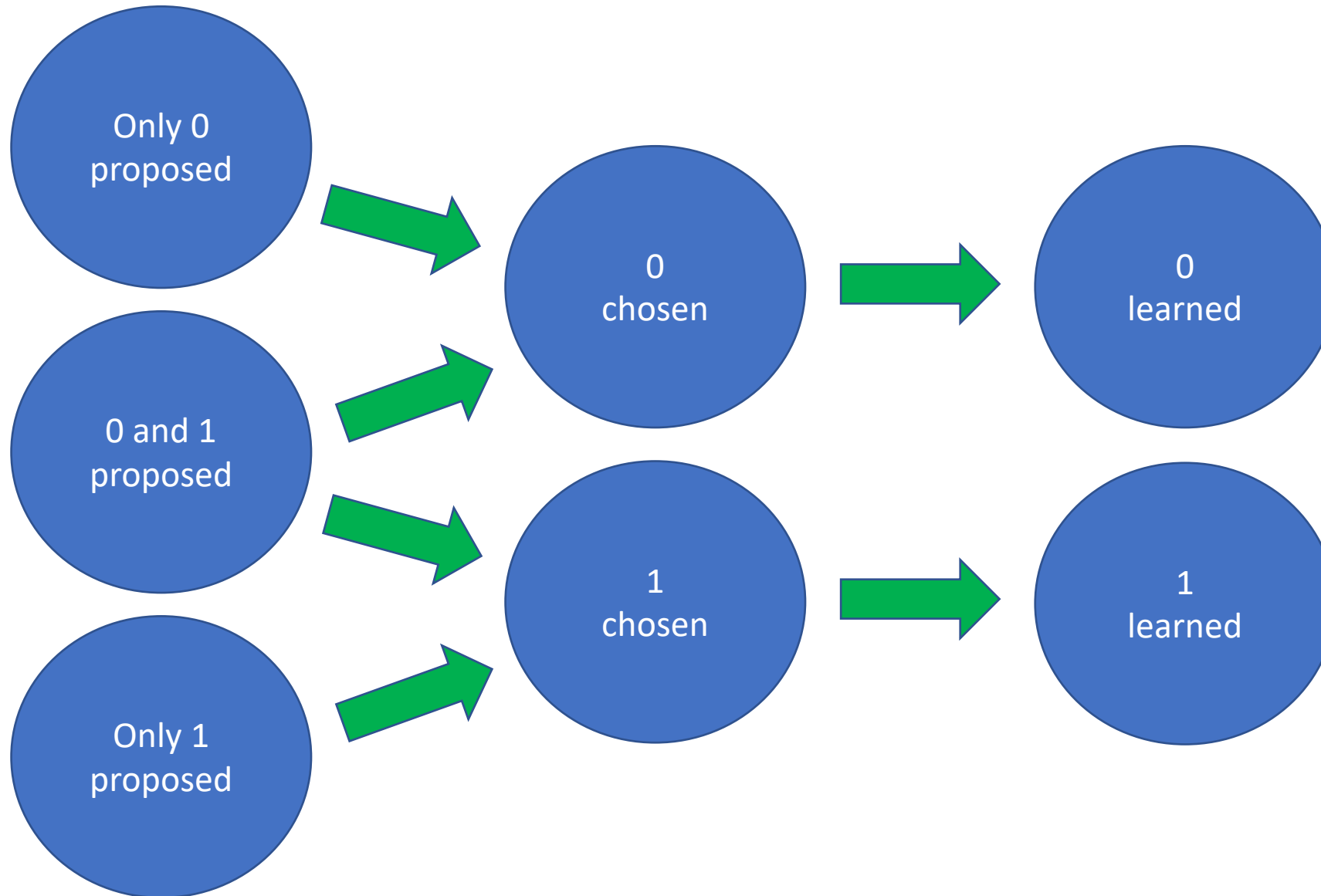
$Hist_{objID} := Hist_{objID} \cdot r$

**reply** according options  $opts$  based  
on  $Hist_{objID}$

Figure 1: Client's View of an Object.

It's not always possible to get a refinement 😞

# Binary Consensus, Specification



# Paxos

- Value is chosen if a quorum of proposers have all accepted the value on the same ballot
- This suggest an easy mapping of the Paxos state to the consensus state

# Problem 1: lack of history

- Unfortunately, Paxos acceptors only remember the latest value they accepted
- So while there may exist a majority that have all accepted the value at time  $t$ , that majority may no longer exist at time  $t+1$ 
  - Even though it is guaranteed that no other value will ever be chosen

# Fix 1: add history variables

- We can add a “ghost variable” to each acceptor that remembers all (value, ballot) pairs it has ever accepted
  - “ghost” means that it does not actually have to be realized
- With this “history variable”, we can exhibit a state mapping

# Problem 2: outrunning the specification

- A refinement mapping maps each step of the low-level specification to either one step of the high-level specification or a stuttering step of the high-level specification
- In Paxos, when  $f=1$  and  $n=3$ , the following scenario is possible:
  - Leader proposes a (value, ballot)
  - Some acceptor accepts (value, ballot)
  - In that one step:
    - The value is chosen
    - The acceptor learns that the value is chosen (decided)
- However, our high-level consensus spec requires two steps:
  - From undecided to chosen and from chosen to learned



# Fix 2: two possibilities

- Change the high-level spec to include a “choose + learn” step
  - i.e., speed up the high-level spec
  - complicates the high-level specification
  - changing the specification may not be allowed
- Add a ghost “prophecy variable” to the low-level specification
  - slow down the low-level spec
  - artificially insert a step between accepting and learning by changing the prophecy variable
  - does not change either the implementation or the high-level spec

# Completeness

- If  $S1$  implements  $S2$  then, possibly by adding history and prophecy variables, there exists a refinement mapping from  $S2$  to  $S1$  (under certain reasonable assumptions)

See Martin Abadi and Leslie Lamport, “The Existence of Refinement Mappings”