

Validating Formal Descriptions of TCP/IP

Experimental Formal Semantics

Michael Norrish

Michael.Norrish@nicta.com.au

NICTA

September 2007



Background

- This is joint work with Peter Sewell, Keith Wansbrough, Tom Ridge, Steve Bishop, Andrei Serjantov, Michael Fairbairn and Michael Smith (all at University of Cambridge).
- Project began in 2001(?), with work on UDP by Sewell, Serjantov and Wansbrough.
- UDP is a simple protocol, but work on a detailed semantics became overwhelming for pen-and-paper techniques.
- I joined the project to see if mechanical support might be helpful (hence the choice of HOL4...)

Work on UDP

- The mechanised semantics for UDP was not **too** large: I proved a simple safety property in HOL.
- By hand, we also proved
 - a timing property
 - correctness properties for a heart-beat program built using the sockets interface to UDP
- Next step was clear: move to TCP.

Introduction

Beginning a TCP
Specification

The Segment Level Specification

Use of HOL

Specification Validation

As a Theorem
Proving Problem

The High Level Specification

Conclusion

1 Introduction

Beginning a TCP Specification

2 The Segment Level Specification

Use of HOL

3 Specification Validation

As a Theorem Proving Problem

4 The High Level Specification

5 Conclusion

Introduction

Beginning a TCP
SpecificationThe Segment Level
Specification

Use of HOL

Specification
ValidationAs a Theorem
Proving ProblemThe High Level
Specification

Conclusion

Motivation

TCP is critical Internet infrastructure.

Time spent specifying it is time well-spent:

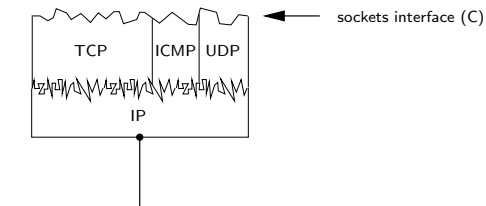
- Users of the API know what to expect
- Future implementors have a better idea of what they have to do
- It can be studied mathematically, not just empirically.

...implementing TCP correctly is very difficult
—Vern Paxson (SIGCOMM'97)

To which we add:

Using TCP and the Sockets API correctly is also difficult

Networking: The TCP/IP Protocols



IP (Internet Protocol): unreliable asynchronous small messages, delivered to IP addresses such as 128.34.1.14.

UDP (User Datagram Protocol): as above, but delivered to IP address/Port pairs.

TCP (Transmission Control Protocol): duplex streams, with retransmission, flow control, congestion control, etc. Messages between IP/Port pairs.

Introduction

Beginning a TCP
Specification

The Segment Level Specification

Use of HOL

Specification Validation

As a Theorem
Proving Problem

The High Level Specification

Conclusion

Networking: The Sockets API

Introduction

Beginning a TCP
Specification

The Segment Level Specification

Use of HOL

Specification Validation

As a Theorem
Proving Problem

The High Level Specification

Conclusion

The Sockets API gives programmers levers with which to control of the various internet protocols:

- Expression in C is ghastly (`ntohs`, `struct inaddr *`...)
- Even stripped of C-isms, there are a plethora of confusing entry-points (`bind`, `listen`, `accept`, `connect`, `close`, `socket`, `dupfd`...)
- Specified in POSIX and the various implementations (including Windows)

Post Hoc Specification

Introduction

Beginning a TCP
Specification

The Segment Level Specification

Use of HOL

Specification Validation

As a Theorem
Proving Problem

The High Level Specification

Conclusion

We cannot tell the world's TCP and OS implementors what to do.

Our specification must reflect not only the existing “specifications”:

- RFCs, POSIX, ...

But also current practice, as enshrined in representative implementations: Windows, BSD and Linux.

Writing a Huge Specification

This project was a test of theoretical machinery:

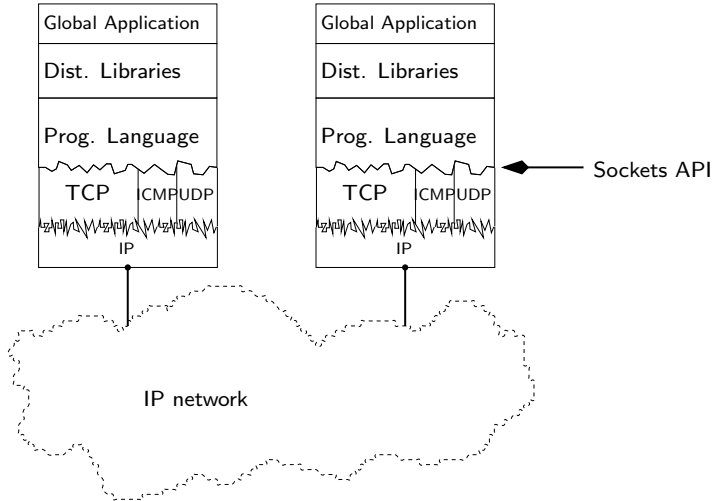
- An operational semantics with almost 200 reduction rules (no recursion though)
- Handling of non-determinism
 - timing: modelling the interleaving of asynchronous and synchronous system calls, as well as packet arrival and dispatch
 - choices of values: under-specified behaviours cause semantic states to take on constrained values for attributes
- Quantified time: TCP is full of timed quantities, and counters representing the passage of time.

Writing a Huge Specification

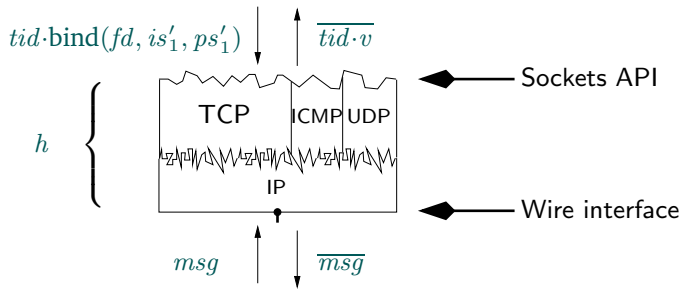
In addition to “inherent complexities”, we contended with 20–30 years of haphazard code evolution, resulting in

- the warped sockets API;
- specifications allowing a great deal of implementation latitude;
- numerous ugly corner cases

Specification: Where to Cut?



Specification: Where to Cut?



The specification describes the evolution of **hosts**, which are involved in six sorts of behaviours:

- system calls **in**; syscall returns **out**;
- messages **in**; messages **out**;
- time elapsing; internal/unobservable state-changes.

Specification Language

The specification is written in higher-order logic.

Expressive:

- Supports natural, mathematical idiom
- Rich types (lists, sets, finite-maps, \mathbb{N} , \mathbb{R})
 - + user defined types (records, algebraic types)
- Captures non-determinism and under-specification easily

Clear:

- Has well-defined semantics
- Easy to write (non-expert CS people picked it up in a week or so)

Introduction

Beginning a TCP
Specification

The Segment Level
Specification

Use of HOL

Specification
Validation

As a Theorem
Proving Problem

The High Level
Specification

Conclusion

The Segment Level Specification

Our segment level specification (the “low level spec”) describes

- the “on the wire” behaviour of hosts: the packets they emit
- what programmers using the sockets API can expect
- the internal (hidden) state of hosts supporting the above

The Segment Level Specification

Our segment level specification (the “low level spec”) describes

- the “on the wire” behaviour of hosts: the packets they emit
- what programmers using the sockets API can expect
- the internal (hidden) state of hosts supporting the above

We have

- validated the specification against real implementations
- found bugs in the implementations

Important Types: Segments

```
tcpSegment = ⟨  
  is1 : ip option;           (* source IP *)  
  is2 : ip option;           (* destination IP *)  
  ps1 : port option;         (* source port *)  
  ps2 : port option;         (* destination port *)  
  seq : tcp_seq_local;       (* sequence number *)  
  URG, ACK : bool;  
  PSH, RST : bool;  
  SYN, FIN : bool;  
  win : word16;              (* window size (unsigned) *)  
  mss : word16 option;       (* maximum segment size *)  
  ...  
⟩
```


Important Types: Hosts

```
host = ⟨  
  arch : arch;                                (* architecture *)  
  privs : bool;  
  socks : sid ↦ socket;  
  ts : tid ↦ hostThreadState timed;          (* threads *)  
  listen : sid list;  
  iq : msg list timed;                        (* messages in *)  
  oq : msg list timed;                        (* messages out *)  
  ticks : ticker  
  ...  
⟩
```

Specification: A Sample Rule

bind_5 rp_all: fast fail Fail with EINTR: the socket is already bound to an address and does not support rebinding; or socket has been shutdown for writing on FreeBSD

$$h\langle\langle ts := ts \oplus (tid \mapsto (\text{Run})_d) \rangle\rangle$$

$$\frac{tid.\text{bind}(fd, is1, ps1)}{\rightarrow} h\langle\langle ts := ts \oplus (tid \mapsto (\text{Ret}(\text{FAIL EINTR}))_{\text{sched_timer}}) \rangle\rangle$$

$$\begin{aligned} &fd \in \mathbf{dom}(h.fds) \wedge fid = h.fds[fd] \wedge \\ &h.files[fd] = \text{File}(\text{FT_Socket}(sid), ff) \wedge \\ &h.socks[sid] = sock \wedge \\ &(sock.ps1 \neq * \vee \\ &(\text{bsd_arch } h.arch \wedge sock.pr = \text{TCP_PROTO}(tcp_sock) \wedge \dots)) \end{aligned}$$

Specification Tools

Somewhat more than just “vanilla” HOL4.

I wrote very little of the specification.

Specifiers who were new to HOL wrote the specification.

They:

- were happy with straight HOL syntax (didn't want to develop their own custom language and parse this to HOL);
- developed a type-setting tool that was fed stereotyped HOL script (would have surely been happy with Isabelle's facilities)
- blithely wrote in a style that was elegant, mathematical and hard to execute

Experience was a great opportunity to test and improve the implementation:

- Error-reporting was improved.
E.g., Keith “*welcome-to-the-1980s*” Wansbrough implemented line-number reporting on parse errors.
- Scaling problems were addressed:
 - 150KB string literals
 - records with 40+ fields became possible
(avoid nasty $O(n^2)$ behaviours)
 - de Bruijn indices for term abstractions dropped
(and other kernel experiments)
- Spent a lot of time with the simplifier

How Do We Know It's Right?

It would be surprising if a logical entity of this size did not contain errors.

The specification is not really executable

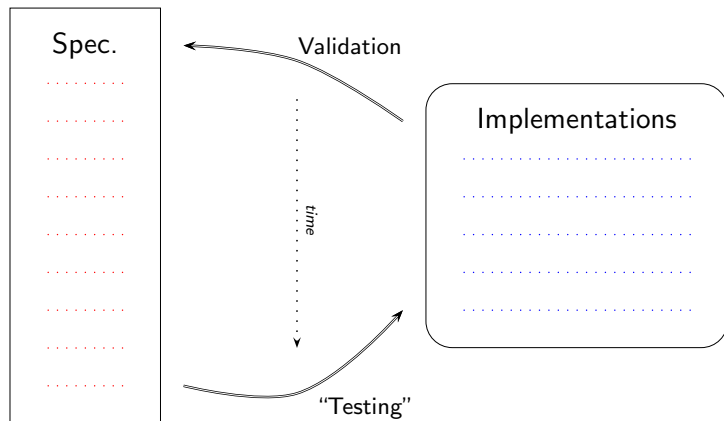
- we did toy with the idea of trying to execute it in concert with specially slowed-down real implementations

But we can **test** implementations against it.

Specification Validation and Testing

Initially, a **post hoc** specification must be held to be in error if it doesn't capture what the implementations do.

Once refined, implementations can be validated against **it**.



Validation/Testing Problem

Observe two machines communicating.

Record trace of events (with time-stamps).

Check that trace of real-world events is compatible with specification.

Validation/Testing Problem

Observe two machines communicating.

Record trace of events (with time-stamps).

Check that trace of real-world events is compatible with specification.

Disagreement indicates

- specification is wrong (initially)
- a possible bug (later)

This is “Just” Testing

Testing programs with simple specifications is easy.

For example, simple observation + calculation can determine if a program for calculating square roots has behaved correctly.

This is “Just” Testing

Testing programs with simple specifications is easy.

For example, simple observation + calculation can determine if a program for calculating square roots has behaved correctly.

With many tools (software model-checking, ESC/Java etc) the programs analysed have implicit specification:

Does not do anything bad.

This is “Just” Testing

Testing programs with simple specifications is easy.

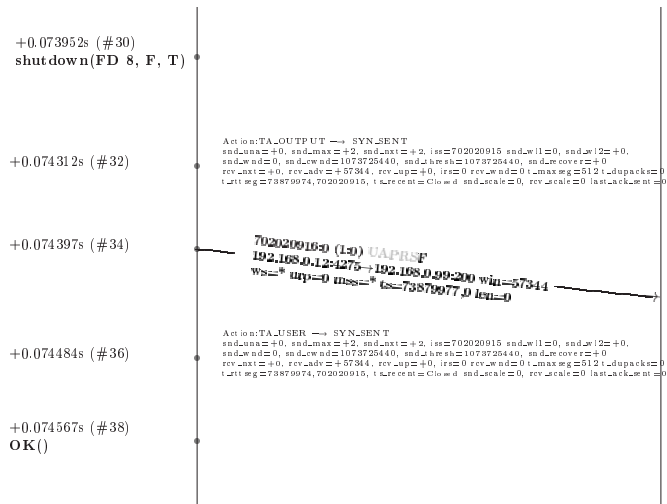
For example, simple observation + calculation can determine if a program for calculating square roots has behaved correctly.

With many tools (software model-checking, ESC/Java etc) the programs analysed have implicit specification:

Does not do anything bad.

With a specification as complicated as TCP's, testing is a hard problem.

Sample Trace Fragment



Outline

- 1 Introduction
Beginning a TCP Specification
- 2 The Segment Level Specification
Use of HOL
- 3 Specification Validation
As a Theorem Proving Problem
- 4 The High Level Specification
- 5 Conclusion

The Trace Checking Problem

Problem: Given an initial host h_0 and a captured trace l_1, \dots, l_n , determine if

$$\exists h_1 \dots h_n. h_0 \xrightarrow{\tau^*} \xrightarrow{l_1} h_1 \xrightarrow{\tau^*} \xrightarrow{l_2} h_2 \dots \xrightarrow{\tau^*} \xrightarrow{l_n} h_n$$

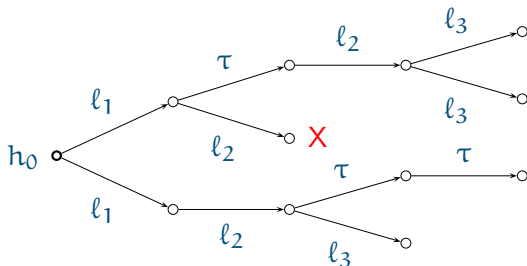
The reduction relation, defining $\xrightarrow{\ell}$ (including $\ell = \tau$), is given by the people writing the specification, and is subject to change as validation proceeds.

This is “existential model checking”: can a path of the specified type be exhibited by the given model?

Why Validation is Hard

Specification is not (can not be) deterministic.

Though we know the initial state, and the observed behaviours, there may be multiple paths possible



The X indicates a state from which no τ or l_3 transition is possible.

Backtracking is required from such a state.

More Non-determinism: Accumulating Constraint Sets

In addition to branching (“competition” between many reduction rules), non-determinism arises within a single rule.

E.g.:

$$\frac{0 \leq i < h.fld_2}{h \xrightarrow{\ell} h\langle fld_1 := i \rangle}$$

Every transition becomes associated with a set of constraints that have to be true for that transition to have occurred.

States come to have symbolic values.

Accumulating Constraint Sets

The result is the proof of a sequence of theorems :

$$\begin{aligned}\Gamma_0 &\vdash h_0 \xrightarrow{\ell_0} h_1 \\ \Gamma_0 \cup \Gamma_1 &\vdash h_1 \xrightarrow{\ell_1} h_2 \\ \Gamma_0 \cup \Gamma_1 \cup \Gamma_2 &\vdash h_2 \xrightarrow{\ell_2} h_3 \\ &\dots\end{aligned}$$

with each Γ_i being the set of constraints associated with the i th transition.

The growing constraint sets have to be checked for satisfiability at each step.

(This is computationally horrific.)

Accumulating Constraint Sets

The result is the proof of a sequence of theorems :

$$\begin{aligned}\Gamma_0 &\vdash h_0 \xrightarrow{\ell_0} h_1 \\ \Gamma_0 \cup \Gamma_1 &\vdash h_1 \xrightarrow{\ell_1} h_2 \\ \Gamma_0 \cup \Gamma_1 \cup \Gamma_2 &\vdash h_2 \xrightarrow{\ell_2} h_3 \\ &\dots\end{aligned}$$

with each Γ_i being the set of constraints associated with the i th transition.

The growing constraint sets have to be checked for satisfiability at each step.

At end of process can “ground” complete trace by finding values satisfying all constraints.

Typical Constraints

```
pending (cb'_0_t_rttvar =
        MAX (1 / 1600)
            ((1 + (abs delta - 1 / 40) * 10) / 40))
pending (cb'_0_t_srtt = MAX (1 / 3200)
        ((2 + delta * 5) / 40))
pending (delta = r / 100 - 1 / 100 - 1 / 20)
pending (r = real_of_int (ticks_of ticks'22 -
                        SEQ32 Tstamp 152410675w))
MIN 64 (MAX 1 (cb'_0_t_srtt + 4 * cb'_0_t_rttvar)) = 1
r <= 6
3 <= 160 * cb'_0_t_rttvar
19 <= 400 * cb'_0_t_srtt
100 <= 25 * r
800 * cb'_0_t_rttvar <= 19
20 * cb'_0_t_srtt <= 1
```

Typical Constraints

```
pending (cb'_0_t_rttvar =
        MAX (1 / 1600)
            ((1 + (abs delta - 1 / 40) * 10) / 40))
pending (cb'_0_t_srtt = MAX (1 / 3200)
        ((2 + delta * 5) / 40))
pending (delta = r / 100 - 1 / 100 - 1 / 20)
pending (r = real_of_int (ticks_of ticks'22 -
                        SEQ32 Tstamp 152410675w))
MIN 64 (MAX 1 (cb'_0_t_srtt + 4 * cb'_0_t_rttvar)) = 1
r <= 6
3 <= 160 * cb'_0_t_rttvar
19 <= 400 * cb'_0_t_srtt
100 <= 25 * r
800 * cb'_0_t_rttvar <= 19
20 * cb'_0_t_srtt <= 1
```

What is `pending` all about?

Lazy Handling of Equality Constraints

Many constraints in a Γ_i are equalities on freshly introduced variables. (Often in turn derived from `let` expressions in the specification.)

`pending` is used to “tag” equalities that shouldn’t be eliminated, like

$$v = \text{if } P \text{ then } e1 \text{ else } e2$$

Only substitute out “values”. E.g.,

- $v = 4$
- $v = \text{SOME } 4$
- $v = h :: t$

In $h :: t$, sub-expressions h and t may still be complicated.

- Only substitute out the “cons”, inventing fresh variable bindings for h and t

Handling High-Level Specification Idioms

Validating Formal
Descriptions of
TCP/IP

Michael Norrish

Introduction

Beginning a TCP
Specification

The Segment Level
Specification

Use of HOL

Specification
Validation

As a Theorem
Proving Problem

The High Level
Specification

Conclusion

Specification often seemed to be written to cause me maximum pain:

- Use of complicated set comprehensions
- Use of non-injective patterns
- General use of a declarative style
- ...

In all cases, it's important to **keep** the original specification.

- If I force the specifiers to write Prolog programs, we lose HOL's advantages.
- The response is to **prove** logical equivalences, "translating" the high-level expressions into more executable forms.

Non-injective Pattern-Matching

Specifiers like to write things like:

$$\frac{\text{side conditions}}{h\langle fld_1 := v_1 \rangle \xrightarrow{\ell} h\langle fld_1 := v_2; fld_2 := 10 \rangle}$$

When we try to discover if our current host can make this transition, there is no unique h to choose.

Non-injective Pattern-Matching

Specifiers like to write things like:

$$\frac{\text{side conditions}}{h\langle fld_1 := v_1 \rangle \xrightarrow{\ell} h\langle fld_1 := v_2; fld_2 := 10 \rangle}$$

When we try to discover if our current host can make this transition, there is no unique h to choose.

For records (like hosts), programmatically rewrite above to

$$\frac{\text{side conditions}}{\langle fld_1 := v_1; fld_2 := v'_2; \dots fld_n := v_n \rangle \xrightarrow{\ell} \langle fld_1 := v_2; fld_2 := 10; \dots fld_n := v_n \rangle}$$

Non-injective Pattern-Matching & (Finite) Maps

Validating Formal
Descriptions of
TCP/IP

Michael Norrish

If a specifier writes

$$\frac{\text{side conditions}}{\mathbf{h}\langle ts := tids \oplus (t \mapsto v) \rangle \xrightarrow{\ell} \mathbf{h}\langle ts := tids \oplus (t \mapsto v') \rangle}$$

we can't expand the *ts* map into a function with a fixed domain as we did with records.

Introduction

Beginning a TCP
Specification

The Segment Level
Specification

Use of HOL

Specification
Validation

As a Theorem
Proving Problem

The High Level
Specification

Conclusion

Non-injective Pattern-Matching & (Finite) Maps

Validating Formal
Descriptions of
TCP/IP

Michael Norrish

If a specifier writes

$$\frac{\text{side conditions}}{\mathbf{h}\langle ts := tids \oplus (t \mapsto v) \rangle \xrightarrow{\ell} \mathbf{h}\langle ts := tids \oplus (t \mapsto v') \rangle}$$

we can't expand the ts map into a function with a fixed domain as we did with records.

Rewrite (again, programmatically!) to

$$\frac{\text{side conditions} \quad t \in \text{dom}(f) \quad f(t) = v}{\mathbf{h}\langle ts := f \rangle \xrightarrow{\ell} \mathbf{h}\langle ts := f \oplus (t \mapsto v') \rangle}$$

(Some finite map idioms are not this simple.)

Introduction

Beginning a TCP
Specification

The Segment Level
Specification

Use of HOL

Specification
Validation

As a Theorem
Proving Problem

The High Level
Specification

Conclusion

The Streams Level Specification

- The specification at the segment level is
 - **good** for implementors
 - **bad** (difficult) for users
- TCP implements an abstraction: two reliable streams between the local and remote end-points
- Can we provide a validated streams level specification of TCP/IP?

First, Write Your Streams Level Specification

The Cambridge specifiers wrote a specification that:

- replaces segment movements with reads and writes to pairs of streams
- drastically reduces the complexity of hosts (TCP control blocks disappear)
- retains accurate descriptions of error cases in the Sockets API

Second, Abstract From Low to High

The next step is to write two abstraction functions:

- One, *low-state* \rightarrow *high-state* removes extraneous detail from hosts, and calculates the state of the streams
- Another, *low-label* \rightarrow *high-label* translates low-level observations into high-level observations.
 - Socket calls remain the same;
 - Packet movements on the network become unobservable τ transitions

Third, Validate High-Level Specification

*Every trace that is to be validated at the high level
has already been validated at the low level.*

Consequently:

- have a sequence of low-level hosts witnessing the trace
- know which reduction rules were chosen at the low level
(high-level rules were written to correspond)

So, high-level validation is much simpler. . .

High-Level Validation Problem

Given concrete h_0 , h , ℓ and the particular rule to check (**SOMERULE**), determine if

$$\frac{\quad}{h_0 \xrightarrow{\ell} h} \text{SOMERULE}$$

- There are no symbolic values in h .
- Can find witnesses for any existential variables in **SOMERULE**'s side conditions immediately

Results

We ran approx. 2000 segment level tests on each of Windows, BSD and Linux.

Over 90% of traces now succeed. (Focus on BSD's TCP.)

So, we're confident our specification is

- Precise
- Accurate

We also found what must be bugs. . .

Sample TCP Bugs

- incorrect RTT estimate after repeated retransmission timeouts
- TCPHAVERCVDFIN wrong — so can SIGURG a closed connection
- initial retransmit timer miscalculation
- simultaneous open responds with ACK instead of SYN,ACK
- receive window updated even for bad segment
- urgent pointer not updated in fastpath (so after 2GB, won't work for 2GB)
- shutdown state changes in pre-established states
- (Linux) sending options in a SYN,ACK that were not received in SYN

Using Our Technology in Other Settings

Validating Formal
Descriptions of
TCP/IP

Michael Norrish

Work on TCP is well-motivated.

But, the approach is well-suited to the design and testing of new protocols too.

Introduction

Beginning a TCP
Specification

The Segment Level
Specification

Use of HOL

Specification
Validation

As a Theorem
Proving Problem

The High Level
Specification

Conclusion

Using Our Technology in Other Settings

Validating Formal
Descriptions of
TCP/IP

Michael Norrish

Work on TCP is well-motivated.

But, the approach is well-suited to the design and testing of new protocols too.

In fact, starting by writing a formal specification may make things easier:

- Can write spec. to be executable
- Avoids dangerous hand-waving
- Gives implementors a clear target

Introduction

Beginning a TCP
Specification

The Segment Level
Specification

Use of HOL

Specification
Validation

As a Theorem
Proving Problem

The High Level
Specification

Conclusion

Functional Specifications

Writing a full functional specification can be a **massive** undertaking.

Circumstances when it can be worthwhile:

- Critical applications
- When multiple (interoperating?) implementations are expected

An ideal specification

- Is mechanically manipulable
- Is written in expressive language
- Does not specify too much
- Has a clear meaning

Mechanising a Specification

Specifications should be

- written in advance; and
- with mechanisation in mind

This is not difficult, and focussing the mind early is of great benefit to a design.

A mechanised specification then allows efficient testing of a new implementation.

This was done successfully by Ridge, Sewell, Dales and Jansen in the design of a new optical switch (with Intel Research).

Introduction

Beginning a TCP
Specification

The Segment Level
Specification

Use of HOL

Specification
Validation

As a Theorem
Proving Problem

The High Level
Specification

Conclusion