# Verifying Reliable Sessions Over an Unreliable Network in Distributed Separation Logic

**Léon Gondelman &**
**Jonas Kastberg Hinrichsen**

**The Second Iris Workshop**

**May 2, 2022**

# I. Reliable Communication in Distributed Systems

# Communicating processes

- Network communication & message-passing concurrency:

  > coordination is done via exchanging messages (not via shared memory)

  > communication protocols and resource transfer play central role

# Fundamental Difference

- Communication over the network is fundamentally unreliable and asynchronous:

  > messages are lost, arrive out of order, got duplicated, or forged by adversary

  > messages arrive from one machine to another with a certain delay

  > network partitions make it impossible to distinguish, in a finite amount of time, between delayed messages and lost messages (e.g. due to remote's crash)

# Fault Tolerance

- Transport layer protocols such as TCP, SCTP and others provide some reliability guarantees (*at-most-once in-order delivery*).

- However, no protocol can guarantee that messages *will arrive in-order & without duplicates exactly once.*

- In the presence of network partitions/broken connections, TCP is no better than UDP: *in fine,* reliability is achieved at the application level.

- Many reasons to build fault-tolerance on top of UDP:
  *> gaming community, Google QUIC (2013), Ensemble (Haiden 98)*

# Verification Perspective

- Two research directions:

  > **Assume** fault-tolerance **to reason** about high-level problems/algorithms:
  *map-reduce, deadlock freedom, op-based CRDTs, …*

  > **Model** network with faults **to build** fault-tolerance:
  *consensus algorithms, reliable causal broadcast, client-server sessions.*

- Longstanding goal: a unified framework where high-level abstractions meet realistic fault-tolerant implementations.

- The story of this work: one step towards this goal.

- Actris Session Type-based Reasoning

  > provides a high-level model of reliable communication (Actris Ghost Theory)

  > has been applied so far only to reason about message-passing concurrency, where the communication layer itself is reliable.

$$\{c \rightarrowtail \,!\,\vec{x}:\vec{\tau}\,\langle v\rangle\{P\}.\,prot * P[\vec{t}/\vec{x}]\}$$
$$\quad \texttt{send } c \ (v[\vec{t}/\vec{x}])$$
$$\{c \rightarrowtail prot[\vec{t}/\vec{x}]\}$$

$$\{c \rightarrowtail \,?\,\vec{x}:\vec{\tau}\,\langle v\rangle\{P\}.\,prot\}$$
$$\quad \texttt{recv } c$$
$$\{w.\,\exists(\vec{y}:\vec{\tau}).\,(w = v[\vec{y}/\vec{x}]) *$$
$$\qquad\qquad P[\vec{y}/\vec{x}] * c \rightarrowtail prot[\vec{y}/\vec{x}]\}$$

- **Aneris Distributed Separation Logic**

  > provides rules to reason about unreliable unconnected communication;

  > had no native/library support for reliable/connected communication (i.e. each time reliability/sessions had to be built in ad-hoc way).

HT-SEND
$$\left\{ \begin{array}{c} sh \xrightarrow{m.\mathrm{src}_{\mathrm{ip}}} (\mathsf{Some}(m.\mathrm{src}), b) * m.\mathrm{dst} \mapsto \Phi * \\ m.\mathrm{src} \rightsquigarrow (R, T) * (m \notin T \Rightarrow \Phi\ m) \end{array} \right\}$$
$$\langle m.\mathrm{src}_{\mathrm{ip}};\ \mathtt{sendto}\ sh\ m.\mathrm{str}\ m.\mathrm{dst} \rangle$$
$$\left\{ \begin{array}{c} w.\ w = |m.\mathrm{src}| * m.\mathrm{src} \rightsquigarrow (R, T \cup \{m\}) * \\ sh \xrightarrow{m.\mathrm{src}_{\mathrm{ip}}} (\mathsf{Some}(m.\mathrm{src}), b) \end{array} \right\}$$

HT-RECV
$$\left\{ sh \xrightarrow{sa_{\mathrm{ip}}} (\mathsf{Some}(sa), b) * sa \rightsquigarrow (R, T) * sa \Mapsto \Phi \right\}$$
$$\langle sa_{\mathrm{ip}};\ \mathtt{receivefrom}\ sh \rangle$$
$$\left\{ \begin{array}{c} w.\ sh \xrightarrow{sa_{\mathrm{ip}}} (\mathsf{Some}(sa), b) * \\ (b = \mathsf{false} * w = \mathsf{None} * sa \rightsquigarrow (R, T)) \vee \\ (\exists m.\ w = \mathsf{Some}\ (m.\mathrm{str}, m.\mathrm{src}) * m.\mathrm{dst} = sa * \\ sa \rightsquigarrow (R \cup \{m\}, T) * (m \notin R \Rightarrow \Phi\ m)) \end{array} \right\}$$

*(a) socket handle resource*   $sh \xrightarrow{sa_{\mathrm{ip}}} (\mathsf{Some}(sa), b)$

- **Aneris Distributed Separation Logic**

  > provides rules to reason about unreliable unconnected communication;

  > had no native/library support for reliable/connected communication
  (i.e. each time reliability/sessions had to be built in ad-hoc way).

$$
\text{HT-SEND}
$$
$$
\left\{
\begin{array}{l}
sh \xrightarrow{m.src_{ip}} (\text{Some}(m.src), b) * m.dst \Mapsto \Phi * \\
m.src \rightsquigarrow (R, T) * (m \notin T \Rightarrow \Phi\, m)
\end{array}
\right\}
$$
$$
\langle m.src_{ip}; \text{ sendto } sh\ m.str\ m.dst \rangle
$$
$$
\left\{
\begin{array}{l}
w.\ w = |m.src| * m.src \rightsquigarrow (R, T \cup \{m\}) * \\
\quad sh \xrightarrow{m.src_{ip}} (\text{Some}(m.src), b)
\end{array}
\right\}
$$

$$
\text{HT-RECV}
$$
$$
\left\{ sh \xrightarrow{sa_{ip}} (\text{Some}(sa), b) * sa \rightsquigarrow (R, T) * sa \Mapsto \Phi \right\}
$$
$$
\langle sa_{ip}; \text{ receivefrom } sh \rangle
$$
$$
\left\{
\begin{array}{l}
w.\ sh \xrightarrow{sa_{ip}} (\text{Some}(sa), b) * \\
\quad (b = \text{false} * w = \text{None} * sa \rightsquigarrow (R, T)) \vee \\
\quad (\exists m.\ w = \text{Some}\ (m.str, m.src) * m.dst = sa * \\
\qquad sa \rightsquigarrow (R \cup \{m\}, T) * (m \notin R \Rightarrow \Phi\, m))
\end{array}
\right\}
$$

*(b) message history resources* $sa \rightsquigarrow (R, T)$

- Aneris Distributed Separation Logic

  > provides rules to reason about unreliable unconnected communication;

  > had no native/library support for reliable/connected communication
  (i.e. each time reliability/sessions had to be built in ad-hoc way).

$\text{HT-SEND}$
$$\left\{\begin{array}{l} sh \xrightarrow{m.\text{src}_{\text{ip}}} (\text{Some}(m.\text{src}), b) * m.\text{dst} \Longmapsto \Phi \ * \\ m.\text{src} \rightsquigarrow (R, T) * (m \notin T \Rightarrow \Phi\ m) \end{array}\right\}$$
$$\langle m.\text{src}_{\text{ip}};\ \texttt{sendto}\ sh\ m.\text{str}\ m.\text{dst}\rangle$$
$$\left\{\begin{array}{l} w.\ w = |m.\text{src}| * m.\text{src} \rightsquigarrow (R, T \cup \{m\}) \ * \\ sh \xrightarrow{m.\text{src}_{\text{ip}}} (\text{Some}(m.\text{src}), b) \end{array}\right\}$$

$\text{HT-RECV}$
$$\left\{sh \xrightarrow{sa_{\text{ip}}} (\text{Some}(sa), b)\ *\ sa \rightsquigarrow (R, T) * sa \Longmapsto \Phi\right\}$$
$$\langle sa_{\text{ip}};\ \texttt{receivefrom}\ sh\rangle$$
$$\left\{\begin{array}{l} w.\ sh \xrightarrow{sa_{\text{ip}}} (\text{Some}(sa), b)\ * \\ \quad (b = \texttt{false}\ *\ w = \text{None} * sa \rightsquigarrow (R, T)) \vee \\ \quad (\exists m.\ w = \text{Some}\ (m.\text{str}, m.\text{src}) * m.\text{dst} = sa\ * \\ \qquad sa \rightsquigarrow (R \cup \{m\}, T) * (m \notin R \Rightarrow \Phi\ m)) \end{array}\right\}$$

*(c) socket protocol predicate*   $sa \Longmapsto \Phi$

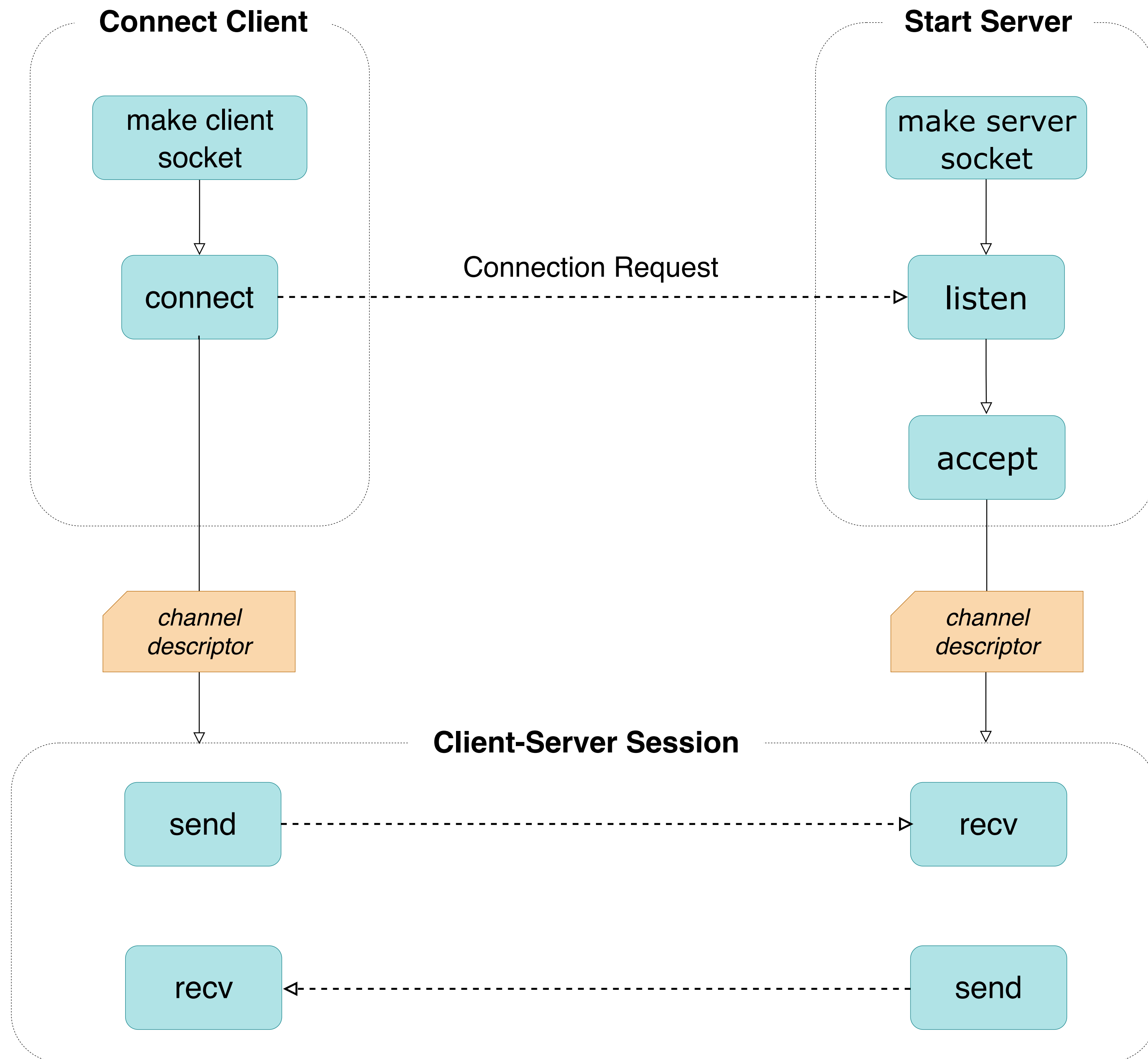Let Aneris and Actris projects meet to enable reasoning
about reliable network communication!


…The rendez-vous point is our verified client-server library.

# II. The API of the library

# Our Library

- BSD sockets-like primitives

- 4-handshake connection

- buffered bidirectional channels

- sequence-ids/acknowledgments/
retransmission mechanisms

- ~ 350 lines of OCaml

- distinction between active/passive
sockets and channels

- data transfer of serialisable values

**Connect Client**

make client socket

connect

**Start Server**

make server socket

listen

accept

Connection Request

*channel descriptor*

*channel descriptor*

**Client-Server Session**

send → recv

recv ← send

# OCaml API

Explicit distinction between *active/passive socket* and *channel descriptor* datatypes

```ocaml
open Ast

  type ('a, 'b) client_skt
  type ('a, 'b) server_skt
  type ('a, 'b) chan_descr
  val make_client_skt : 'a serializer -> 'b serializer -> saddr -> ('a, 'b) client_skt
  val make_server_skt : 'a serializer -> 'b serializer -> saddr -> ('a, 'b) server_skt
  val server_listen : ('a, 'b) server_skt -> unit
  val accept : ('a, 'b) server_skt -> ('a, 'b) chan_descr * saddr
  val connect : ('a, 'b) client_skt -> saddr -> ('a, 'b) chan_descr
  val send : ('a, 'b) chan_descr -> 'a -> unit
  val try_recv : ('a, 'b) chan_descr -> 'b option
  val recv : ('a, 'b) chan_descr -> 'b
```

How **client** serialises values
to be send to the **server**

How **server** deserialises values
received from the **client**

```ocaml
open Ast

  type ('a, 'b) client_skt
  type ('a, 'b) server_skt
  type ('a, 'b) chan_descr
  val make_client_skt : 'a serializer -> 'b serializer -> saddr -> ('a, 'b) client_skt
  val make_server_skt : 'a serializer -> 'b serializer -> saddr -> ('a, 'b) server_skt
  val server_listen : ('a, 'b) server_skt -> unit
  val accept : ('a, 'b) server_skt -> ('a, 'b) chan_descr * saddr
  val connect : ('a, 'b) client_skt -> saddr -> ('a, 'b) chan_descr
  val send : ('a, 'b) chan_descr -> 'a -> unit
  val try_recv : ('a, 'b) chan_descr -> 'b option
  val recv : ('a, 'b) chan_descr -> 'b
```
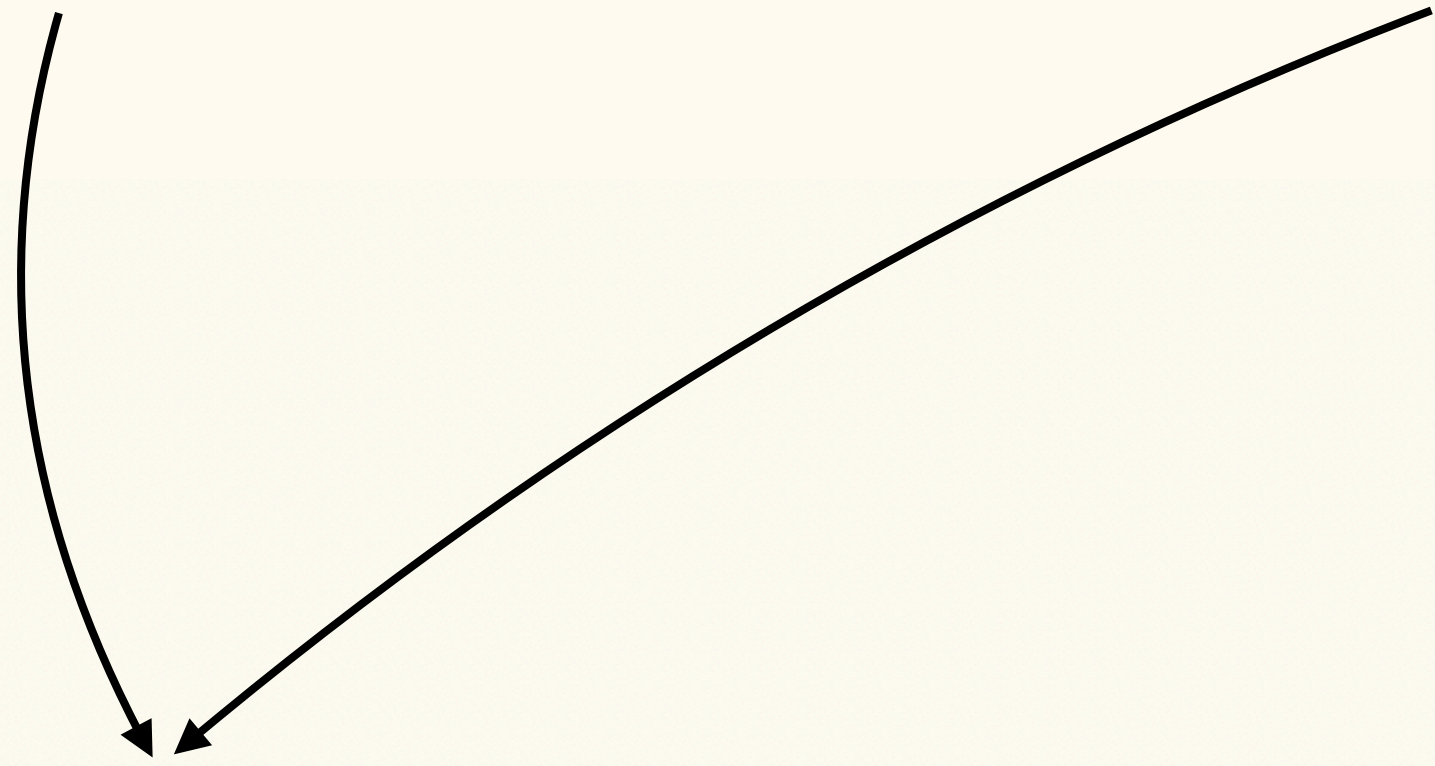
How **server** serialises values
to be send to the **client**

How **client** deserialises values
received from the **server**

```ocaml
open Ast

    type ('a, 'b) client_skt
    type ('a, 'b) server_skt
    type ('a, 'b) chan_descr
    val make_client_skt : 'a serializer -> 'b serializer -> saddr -> ('a, 'b) client_skt
    val make_server_skt : 'a serializer -> 'b serializer -> saddr -> ('a, 'b) server_skt
    val server_listen : ('a, 'b) server_skt -> unit
    val accept : ('a, 'b) server_skt -> ('a, 'b) chan_descr * saddr
    val connect : ('a, 'b) client_skt -> saddr -> ('a, 'b) chan_descr
    val send : ('a, 'b) chan_descr -> 'a -> unit
    val try_recv : ('a, 'b) chan_descr -> 'b option
    val recv : ('a, 'b) chan_descr -> 'b
```

# Example: echo server

```
open Ast
open Serialization_code
open Client_server_code

let int_s = int_serializer
let str_s = string_serializer

let rec echo_loop c =
  let req = recv c in
  send c (strlen req);
  echo_loop c

let accept_loop s =
  let rec loop () =
    let c = fst (accept s) in
    fork echo_loop c; loop ()
  in loop ()

let server srv =
  let s = make_server_skt int_s str_s srv in
  server_listen s;
  fork accept_loop s
```

```
let client clt srv s1 s2 =
  let s = make_client_skt str_s int_s clt in
  let c = connect s srv in
  send c s1; send c s2;
  let m1 = recv c in
  let m2 = recv c in
  assert (m1 = strlen s1 && m2 = strlen s2)

let client_0 clt srv =
  client clt srv "carpe" "diem"
```

# III. Specification

**User Parameters:**

$$\text{UserParams} \triangleq$$

$$\left\{\begin{array}{ll} \texttt{srv} : \text{Address}; & \texttt{srv\_ser} : \text{Serialization}; \\ \texttt{prot} : \text{iProto}; & \texttt{clt\_ser} : \text{Serialization}; \end{array}\right\}$$

**Session Resources:**

$$\text{SessionResources}(\text{UP} : \text{UserParams}) \triangleq$$

$$\left\{\begin{array}{ll} \texttt{srv\_si} : \text{Message} \to \text{iProp}; & \texttt{CanConnect} : \text{Val} \to \text{Address}; \to \text{iProp}; \\ \texttt{SrvInit} : \text{iProp}; & c \xrightarrow[ser]{ip} \,!\,\vec{x}{:}\vec{\tau}\,\langle v\rangle\{P\}.\,prot \;\; (\textit{mapsto connective}); \\ \texttt{CanListen} : \text{Val} \to \text{iProp}; & \textit{laws about those resources} \;(\textit{e.g. subprotocols}) \\ \texttt{Listens} : \text{Val} \to \text{iProp}; & \end{array}\right\}$$

*Notations* : $\quad$ S := SessionResources(UP), S.srv := UP.srv

## Client Setup:

Ht-make-client-socket

$$\{\text{FreeAddr}(clt) * clt \rightsquigarrow (\emptyset, \emptyset) * clt \neq S.\text{srv}\}$$
$$\langle clt_{\text{ip}}; \text{mk\_clt\_skt} (S.\text{srv\_ser}) (S.\text{clt\_ser}) clt\rangle$$
$$\{w. \exists skt. w = skt * S.\text{CanConnect } clt\, skt\}$$

Ht-connect

$$\{S.\text{CanConnect } clt\, skt\}$$
$$\langle clt_{\text{ip}}; \text{connect } skt\, S.\text{srv}\rangle$$
$$\{w. \exists c. w = c * c \xrightarrow[S.\text{clt\_ser}]{clt_{\text{ip}}} S.\text{prot}\}$$

## Server Setup:

Ht-make-server-socket

$$\left\{ \begin{array}{l} \text{FreeAddr}(S.\text{srv}) * S.\text{srv} \rightsquigarrow (\emptyset, \emptyset) * \\ S.\text{srv} \mapsto S.\text{srv\_si} * S.\text{SrvInit} \end{array} \right\}$$
$$\langle S.\text{srv}_{\text{ip}}; \text{mk\_srv\_skt } S.\text{srv\_ser } S.\text{clt\_ser } S.\text{srv}\rangle$$
$$\{w. \exists skt. w = skt * S.\text{CanListen } skt\}$$

Ht-listen

$$\{S.\text{CanListen } skt\}$$
$$\langle S.\text{srv}_{\text{ip}}; \text{listen } skt\rangle$$
$$\{S.\text{Listens } skt\}$$

Ht-accept

$$\{S.\text{Listens } skt\} \langle S.\text{srv}_{\text{ip}}; \text{accept } skt\rangle \{w. \exists c. w = (c, sa) * S.\text{Listens } skt * c \xrightarrow[S.\text{srv\_ser}]{S.\text{srv}_{\text{ip}}} \overline{S.\text{prot}}\}$$

## Client Setup:

Ht-make-client-socket

$$\{\mathsf{FreeAddr}(clt) * clt \rightsquigarrow (\emptyset, \emptyset) * clt \neq S.\mathsf{srv} \}$$
$$\langle clt_{\mathrm{ip}}; \mathsf{mk\_clt\_skt}\ (S.\mathsf{srv\_ser})\ (S.\mathsf{clt\_ser})\ clt \rangle$$
$$\{w.\ \exists skt.\ w = skt * S.\mathsf{CanConnect}\ clt\ skt\}$$

Ht-connect

$$\{S.\mathsf{CanConnect}\ clt\ skt\}$$
$$\langle clt_{\mathrm{ip}}; \mathsf{connect}\ skt\ S.\mathsf{srv} \rangle$$
$$\{w.\ \exists c.\ w = c * c >\!\!\xrightarrow[S.\mathsf{clt\_ser}]{clt_{\mathrm{ip}}}\!\!S.\mathsf{prot}\}$$

## Server Setup:

Ht-make-server-socket

$$\left\{ \begin{aligned} &\mathsf{FreeAddr}(S.\mathsf{srv}) * S.\mathsf{srv} \rightsquigarrow (\emptyset, \emptyset) * \\ &S.\mathsf{srv} \mapsto S.\mathsf{srv\_si} * S.\mathsf{SrvInit} \end{aligned} \right\}$$
$$\langle S.\mathsf{srv}_{\mathrm{ip}}; \mathsf{mk\_srv\_skt}\ S.\mathsf{srv\_ser}\ S.\mathsf{clt\_ser}\ S.\mathsf{srv} \rangle$$
$$\{w.\ \exists skt.\ w = skt * S.\mathsf{CanListen}\ skt\}$$

Ht-listen

$$\{S.\mathsf{CanListen}\ skt\}$$
$$\langle S.\mathsf{srv}_{\mathrm{ip}}; \mathsf{listen}\ skt \rangle$$
$$\{S.\mathsf{Listens}\ skt\}$$

Ht-accept

$$\{S.\mathsf{Listens}\ skt\}\ \langle S.\mathsf{srv}_{\mathrm{ip}}; \mathsf{accept}\ skt \rangle\ \{w.\ \exists c.\ w = (c, sa) * S.\mathsf{Listens}\ skt\ * c >\!\!\xrightarrow[S.\mathsf{srv\_ser}]{S.\mathsf{srv}_{\mathrm{ip}}}\!\!\overline{S.\mathsf{prot}}\}$$

HT-RELIABLE-RECV

$$\{c \xrightarrow[ser]{ip} ?\vec{x}:\vec{\tau}\langle v\rangle\{P\}.\,prot\}\,\langle ip;\,\mathsf{recv}\ c\rangle\,\{w.\,\exists\vec{y}.\,w = v[\vec{y}/\vec{x}] * c \xrightarrow[ser]{ip} prot[\vec{y}/\vec{x}] * P[\vec{y}/\vec{x}]\}$$

HT-RELIABLE-SEND

$$\{c \xrightarrow[ser]{ip} !\vec{x}:\vec{\tau}\langle v\rangle\{P\}.\,prot\ *\ P[\vec{t}/\vec{x}] * \mathsf{Ser}\ ser\ (v[\vec{t}/\vec{x}])\}\,\langle ip;\,\mathsf{send}\ c\ (v[\vec{t}/\vec{x}])\rangle\,\{c \xrightarrow[ser]{ip} prot[\vec{t}/\vec{x}]\}$$

Init-setup

$\text{True} \Rightarrow \exists\, S : \text{SessionResources}(\text{UP}).$

$\qquad S.\text{SrvInit} *$

$\qquad (\forall sa, \text{Ht-make-client-socket}[S](sa)) *$

$\qquad \text{Ht-make-server-socket}[S] *$

$\qquad (\forall skt\ sa, \text{Ht-connect}[S](skt, sa)) *$

$\qquad (\textit{specs for listen, accept, send, recv, try\_recv})$

```ocaml
let rec echo_loop c =
  let req = recv c in
  send c (strlen req);
  echo_loop c
```

*OCaml function*

```coq
Definition echo_loop : val :=
  rec: "echo_loop" "c" :=
  let: "req" := recv "c" in
  send "c" (strlen "req");;
  "echo_loop" "c".
```

*Generated Coq definition*

```coq
Definition prot_aux (rec : iProto Σ) : iProto Σ :=
  (<! (s : string)> MSG #s ; <? (n : ℕ) > MSG #n {{ ⌜String.length s = n⌝ }}; rec)%proto.
```

*Protocol*

```coq
Lemma wp_echo_loop c :
  {{{ c ↠{S.srv_saddr_ip, S.srv_ser} iProto_dual S.protocol }}}
      echo_loop c @[S.srv_saddr_ip]
  {{{ v, RET v ; ⊥ }}}.
Proof.
  iIntros (Φ) "Hci HΦ". iLöb as "IH". wp_lam.
  wp_recv (s₁) as "_". wp_send with "[//]".
  wp_seq.by iApply ("IH" with "[$Hci]").
Qed.
```
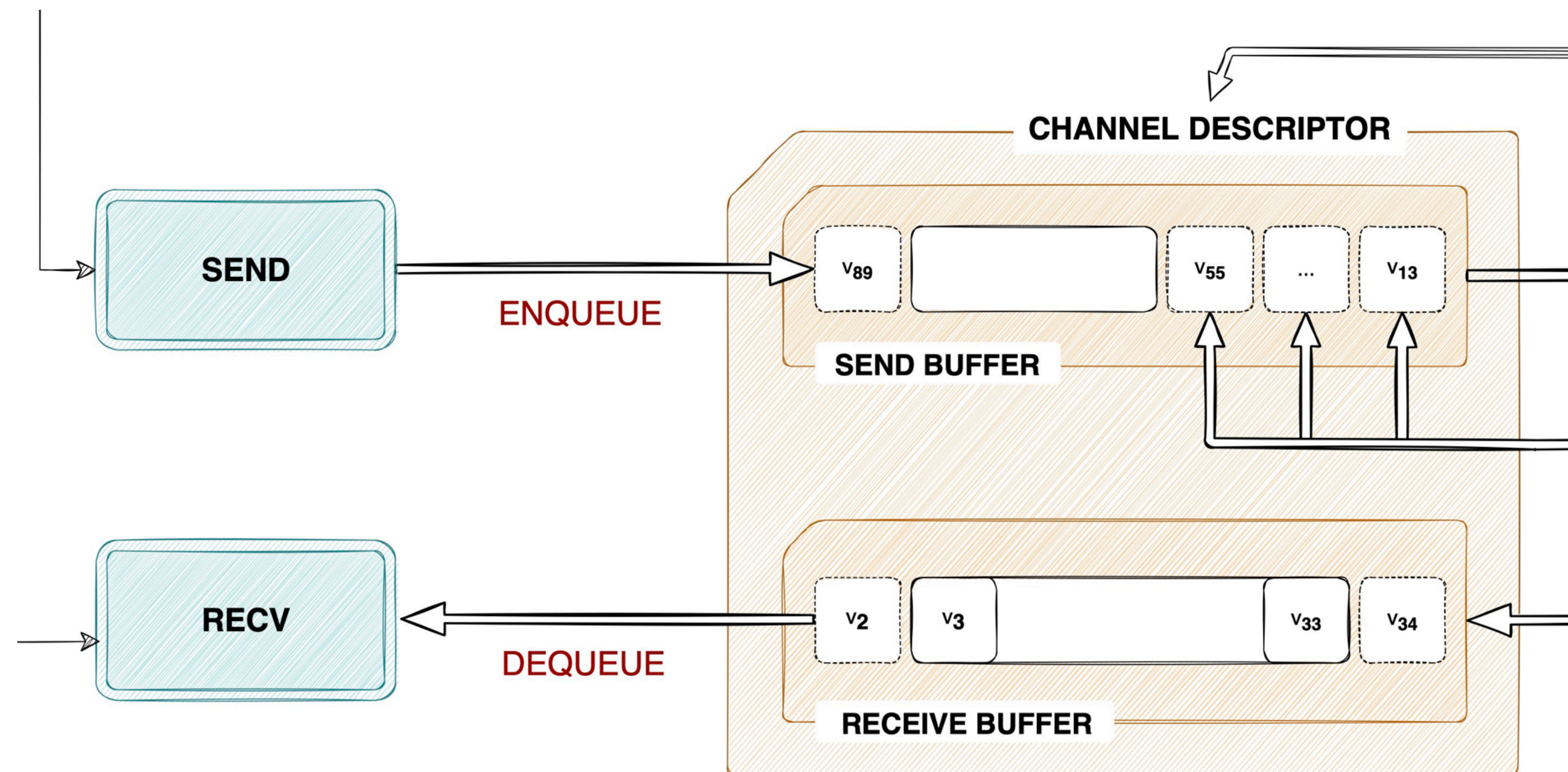
*Proof of echo_loop*

# IV. Verification

# Anatomy of send & recv

- The implementation of *send* and *recv* is the **same for client and server**.
  In fact, their implementation is also **agnostic of network**.

- This is possible because channels are using **in-** and **out- buffers** as indirection
  (calling *send **enqueues*** to the out-buffer, calling *recv **dequeues** from the in-buffer*)

Crucially, this is also where the *connection* between **Actris Ghost Theory** and the **implementation** takes place. However, this connection is not immediate :

- **the two Actris logical buffers**

  > describe symmetrically for each direction the messages in transit

  > are governed (inside an Iris invariant) by the shared resource prot_ctx $\chi$ $\vec{v}_1$ $\vec{v}_2$

- **the four physical buffers**
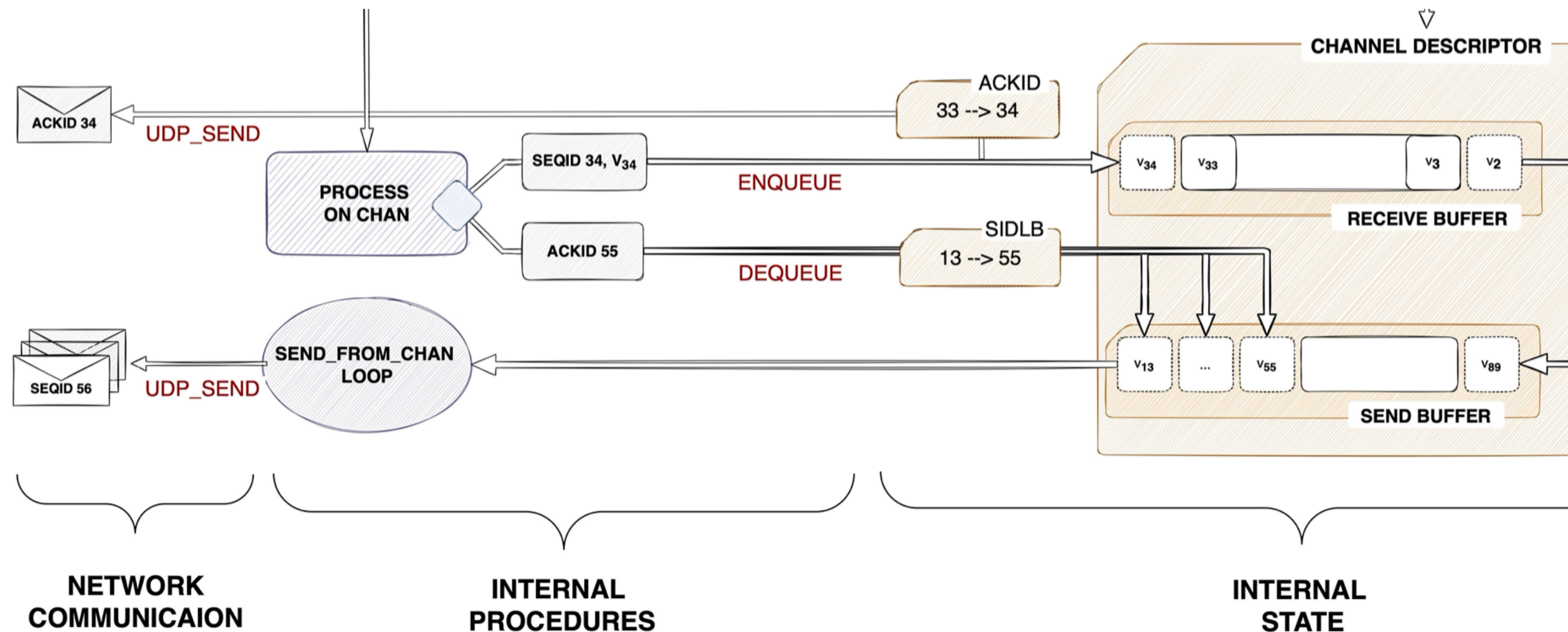
  > play different role (out-buffer simply (re)transmits, in-buffer keeps data for delivery)

  > are local data of each node and are updated asynchronously

# More buffers, seriously ?

- Our solution is to introduce **additional logical buffers** Tl, Rl, Tr, Rr *as a glue.*
*(Tl, Tr) describe the **history of sent** messages;*
*(Rl, Rr) describe the **history of received** messages (by the application)*.

- *Various **relations** must hold between Actris, glue, and physical buffers:*

  - Rr is prefix of Tl and Rl is prefix of Tr $\qquad\qquad$ *(Internal-Coh)*

  - v1 = Tl – Rr and v2 = Tr – Rl $\qquad\qquad$ *(Actris-Coh)*

  - sbufl is suffix of Tl and sbufr is suffix of Tr $\qquad\qquad$ *(SBuf-Coh)*

  - rbufl is prefix of (Tr – Rl) and rbufr is prefix of (Tl – Rr) $\qquad$ *(Rbuf-Coh)*

- The **verification** is then primarily an effort in **preserving these relations**, in the presence of the concurrent accesses of the communication layer.

- The internal procedures that enforce the fault-tolerance are **also (mostly) the same** for clients and servers, and **so are our proofs.**

- The 4-handshake **is different for each side** and requires some effort in verification as it encodes an STS with several edge and absurd cases.

- The implementation/verification of server side is more difficult, because the server must maintain **a table of known clients with their connection state** and a **channel description queue** for the established connections.

# V. Conclusion & Future Directions

# Contributions

# Possible Future Directions

- **Graceful/Abrupt session ending** : *detectable connection failures, reconnection*

- **Cryptography/Security**: *4-way handshake procedure / authentification / QUIC*

- **Network Partitions** : group membership/consensus built on top of our library

- **Group Communication** : *client-service communication*

- **Transparency** : *verified libs for distributed/multithreaded programs (e.g. Functory)*

- (and maybe your insights/ideas !)

# Thank you !

# Backup slides

# Client Implementation



**CONNECTION OPENING**

**DATA TRANSFER**

MAKE CLIENT SOCKET

CONNECT

4-WAY HANDSHAKE

MSG — UDP_RECV

INIT-ACK

INIT — UDP_SEND

COOKIE — UDP_SEND

COOKIE-ACK

CHANNEL DESCRIPTOR

SEND — ENQUEUE

$v_{89}$  $v_{55}$  ...  $v_{13}$

SEND BUFFER

SEND FROM_CHAN LOOP — UDP_SEND — SEQID 56

SIDLB
13 --> 55

ACKID 55 — DEQUEUE

PROCESS ON CHAN LOOP — MSG ID — UDP_RECV

SEQID 34, $v_{34}$ — ENQUEUE

RECV — DEQUEUE

$v_2$  $v_3$  $v_{33}$  $v_{34}$

RECEIVE BUFFER

ACKID
33 --> 34

ACKID 34 — UDP_SEND

USER CALLED METHODS

INTERNAL STATE

INTERNAL PROCEDURES

NETWORK COMMUNICAION

# Server Implementation



**KNOWN CLIENTS TABLE**

saddr ⟼ HALF-OPENED (cookie) ‖
ESTABLISHED (chan_descr, ackid, sidlb)

*shared memory*

**CHANNEL DESCRIPTOR QUEUE**

CHANNEL DESCRIPTOR ... CHANNEL DESCRIPTOR

*shared memory*

**CHANNEL DESCRIPTOR**

$v_{34}$ $v_{33}$ $v_3$ $v_2$

**RECEIVE BUFFER**

$v_{13}$ ... $v_{55}$ $v_{89}$

**SEND BUFFER**

MAKE SERVER SOCKET

LISTEN

ACCEPT

RECV

SEND

**CONNECTION OPENING**

**DATA TRANSFER**

MSG
UDP_RECV

LISTEN_LOOP

INIT — HALF-OPENED

COOKIE — ESTABLISHED

ENQUEUE

DEQUEUE

INIT-ACK
UDP_SEND

COOKIE-ACK
UDP_SEND

ACKID 34
UDP_SEND

ACKID
33 --> 34

SEQID 34, $V_{34}$ — ENQUEUE

PROCESS ON CHAN

ACKID 55 — DEQUEUE

SIDLB
13 --> 55

SEND_FROM_CHAN LOOP

SEQID 56
UDP_SEND

**NETWORK COMMUNICAION**

**INTERNAL PROCEDURES**

**INTERNAL STATE**

**USER CALLED METHODS**