



Offload Sockets Framework and Sockets Direct Protocol High Level Design

Draft 2

June 2002

Revision History and Disclaimers

Rev.	Date	Notes
Draft 1	March 2002	First cut.
Draft 2	June 2002	Updated with new OFFLOAD address family architecture and naming. Combined Socket Framework and SDP/IBT HLD's into one document.. All content included and reviewed by internal design team.

THIS SPECIFICATION IS PROVIDED "AS IS" WITH NO WARRANTIES WHATSOEVER, INCLUDING ANY WARRANTY OF MERCHANTABILITY, NONINFRINGEMENT, FITNESS FOR ANY PARTICULAR PURPOSE, OR ANY WARRANTY OTHERWISE ARISING OUT OF ANY PROPOSAL, SPECIFICATION OR SAMPLE. Intel disclaims all liability, including liability for infringement of any proprietary rights, relating to use of information in this specification. No license, express or implied, by estoppel or otherwise, to any intellectual property rights is granted herein.

This Specification as well as the software described in it is furnished under license and may only be used or copied in accordance with the terms of the license. The information in this document is furnished for informational use only, is subject to change without notice, and should not be construed as a commitment by Intel Corporation. Intel Corporation assumes no responsibility or liability for any errors or inaccuracies that may appear in this document or any software that may be provided in association with this document.

Except as permitted by such license, no part of this document may be reproduced, stored in a retrieval system, or transmitted in any form or by any means without the express written consent of Intel Corporation.

Intel is a trademark or registered trademark of Intel Corporation or its subsidiaries in the United States and other countries.

***Other names and brands may be claimed as the property of others.**

Copyright © 2002 Intel Corporation.

Approval

Role	Signature	Date
Responsible Engineer		
Engineering Group Leader		
Software Engineering Manager		

Abstract

The Offload Sockets Framework (OSF) enables network applications to utilize Linux sockets and File I/O APIs to communicate with remote endpoints across a system area network (SAN), while bypassing the kernel resident TCP/IP protocol stack. The offload sockets framework is completely transport and protocol independent and can be used to support multiple offload technologies. For the rest of this document, as an application of the Offload Sockets Framework, Sockets Direct Protocol (SDP) is used as the target protocol and InfiniBand as the target transport. However, other transport technologies (such as TCP Offload Engines - TOE) and protocols (such as iSCSI) can easily make use of the offload sockets framework.

The Sockets Direct Protocol (SDP) is an InfiniBand specific Protocol defined by the Software Working Group (SWG) of the InfiniBand Trade Association (IBTA). It defines a standard wire protocol over IBA fabric to support user-mode stream sockets-compatible (SOCK_STREAM) networking over IBA. SDP utilizes various InfiniBand features (such as remote DMA (RDMA), memory windows, solicited events etc.) for high-performance zero-copy data transfers. SDP is a pure wire-protocol level specification and does not go into any socket API or implementation specifics. Work is under way to separate the SDP wire-protocol, which is transport and O/S agnostic, from the IBA transport dependencies. This separation will enable transports other than IBA fabrics.

SDP facilitates direct mapping of stream connections to InfiniBand reliable connections (or virtual circuits). IP-Over-IB (IPoIB) specifies a mapping of IP (both v4 & v6) protocols over IBA fabric and treats the IBA fabric simply as the link layer. IPoIB specification is currently being worked on by IETF working group and will be published as an IETF RFC. The IPoIB RFC will define the packet level format of IP packets on the IBA fabric, and also describes an InfiniBand address resolution protocol (IBARP). Conceptually, the IPoIB driver in Linux will look like a network driver and will plug-in underneath the IP stack as any standard Linux network device. The IPoIB driver exposes a network device per IBA port on the host system, and these devices could be used to assign (statically or dynamically - using protocols such as DHCP) IP addresses. The SDP stack simply makes use of these IP assignments for endpoint identifications.

SDP only deals with stream sockets (TCP), and if installed in a system, allow bypassing the OS resident TCP stack for all TCP connections between any endpoints on the IBA fabric. All other socket types (such as datagram, raw, packet etc.) are supported by the Linux IP stack and operate over the IPoIB drivers. The IPoIB stack has no dependency on the SDP stack; however, the SDP stack depends on IPoIB drivers for learning about local IP assignments and for address resolution.

The Linux implementation of SDP includes many interdependent pieces that are referred to as Offload Sockets Framework (OSF). OSF enables applications to use these same standard sockets and File I/O APIs to transparently communicate with remote endpoints/nodes across a system area network (SAN), bypassing the kernel resident TCP/IP protocol stack. This framework will also allow applications to bypass the resident TCP/IP protocol stack while using the default address family of AF_INET.

OSF includes an Offload Protocol Switch (OPS) module (Ipv4 internet protocol switch) that allows integration of Offload Protocol (OP) modules along side the existing TCP/IP stack, an OP module that supports SDP wire protocol, an Offload Transport Interface that provides a common interface to Offload Transport (OT) module that supports InfiniBand hardware. This design document covers the specifics of all three components.

Contents

1.	Introduction	1-1
1.1	Purpose and Scope	1-1
1.2	Audience	1-1
1.3	Acronyms and Terms	1-1
1.4	References	1-1
1.5	Conventions	1-2
1.6	Before You Begin	1-2
2.	Features	2-1
3.	Design Overview	3-1
4.	Design Assumptions & Rules	4-1
4.1	Kernel and API	4-1
4.2	Native AF_INET stack support	4-1
4.2.1	OPS and Protocol Module loading	4-1
4.2.2	Protocol ID's	4-1
4.2.3	Protocol Interfaces	4-1
4.2.4	IP port space	4-2
4.2.5	Generic socket options	4-2
4.2.6	Existing IP protocol specific socket options	4-2
4.2.7	New offload protocol specific socket options	4-3
5.	Design Details	5-4
5.1	Offload Protocol Switch Module	5-4
5.1.1	Linux Kernel Modifications	5-4
5.1.2	Initialization	5-4
5.1.3	Shutdown	5-4
5.1.4	Socket Driver Modifications	5-5
5.1.5	Offload Data Structures	5-7
5.1.6	Socket Create	5-1
5.1.7	Socket Bind	5-1
5.1.8	Socket Listen	5-1
5.1.9	Socket Connect	5-1
5.1.10	Socket Accept and Poll on a Listening Socket	5-3
5.1.11	Socket Operations on a Connected Socket	5-4
5.1.12	Socket Options Processing	5-1
5.2	Sockets Direct Protocol Module	5-1
5.2.1	Linux Kernel Modifications	5-1
5.2.2	Initialization	5-1
5.2.3	Shutdown	5-1
5.2.4	Buffer Strategy	5-2
5.2.5	Connection Services	5-3
5.2.6	Completion Model	5-3
5.2.7	Data Transfer Models	5-3
5.2.8	Locking and Threading Model	5-4

5.3	InfiniBand Offload Transport Module.....	5-5
5.3.1	Linux Kernel Modifications.....	5-6
5.3.2	Initialization.....	5-6
5.3.3	Shutdown.....	5-6
5.3.4	Connection Services.....	5-6
5.3.5	Buffer Strategy.....	5-6
5.3.6	Data Transfer Services.....	5-7
5.3.7	Completion Services.....	5-7
5.3.8	Locking and Threading Model.....	5-8
6.	System Resource Usage.....	6-9
6.1	Sockets Direct Protocol Module.....	6-9
6.2	InfiniBand Transport Module.....	6-9
7.	Internal Compatibility.....	7-1
7.1	Interaction with Other Components.....	7-1
7.1.1	Socket Driver and OPS Module Interaction.....	7-1
7.1.2	OPS and SDP Module Interaction.....	7-2
7.1.3	SDP and InfiniBand Transport Interaction.....	7-4
8.	External Compatibility.....	8-5
8.1	Standards.....	8-5
8.1.1	Sockets Direct Protocol Module.....	8-5
8.1.2	InfiniBand Offload Transport Module.....	8-5
9.	Other Dependencies.....	9-1
9.1	Offload Protocol Switch Module.....	9-1
9.2	Sockets Direct Protocol Module.....	9-1
9.3	InfiniBand Transport Module.....	9-1
10.	Installation and Configuration.....	10-1
10.1	Offload Protocol Switch Module.....	10-1
10.1.1	Installation.....	10-1
10.1.2	Configuration.....	10-1
10.2	Sockets Direct Protocol Module.....	10-1
10.2.1	Installing.....	10-1
10.2.2	Configuring.....	10-1
10.3	InfiniBand Transport Module.....	10-1
10.3.1	Installing.....	10-1
10.3.2	Configuring.....	10-1
11.	Unresolved Issues.....	11-1
11.1	Offload Protocol Switch Module.....	11-1
11.2	Sockets Direct Protocol Module.....	11-1
11.3	InfiniBand Transport Module.....	11-1
12.	Data Structures and APIs.....	12-1
12.1	Offload Protocol Switch Definitions.....	12-1
12.2	SDP Definitions.....	12-1
12.3	InfiniBand OT Definitions.....	12-1

Figures

Figure 1: High Performance Sockets in the kernel	3-1
Figure 2: Linux Socket Driver Modifications	5-5
Figure 3: Socket structures for offload support.....	5-7
Figure 4: Socket Create	5-1
Figure 5: Socket bind	5-1
Figure 6: Socket Listen	5-1
Figure 7: Socket Connect	5-1
Figure 8: Socket Accept & Poll	5-3
Figure 9: Socket Options	5-1
Figure 10: SDP Buffer Mode Overview.....	5-2
Figure 11: InfiniBand offload Transport Service Components.....	5-5
Figure 12: Proposed Patch to Linux Socket Driver.....	7-1
Figure 13: OPS and SDP interaction	7-3
Figure 14 SDP and InfiniBand Transport.....	7-4

1. Introduction

1.1 Purpose and Scope

This HLD defines the implementation of all offload components described in the “*Offload Sockets Framework and Sockets Direct Protocol Architecture Specification*”, including inter-component dependencies, and provides sufficient design detail that will satisfy the product requirements as specified.

1.2 Audience

Anyone interested in understanding this implementation of the Architecture Specification should read this document, including:

- Software developers who are integrating the separate modules into their own software projects
- Hardware developers who need an understanding of the software behavior to optimize their designs
- Evaluation engineers who are developing tests for InfiniBand-compliant devices
- Others in similar roles who need more than a basic understanding of the software

1.3 Acronyms and Terms

OSF:	Offload Sockets Framework (Software components that enables protocol offloading)
OPS:	Offload Protocol Switch (A logical software module that performs protocol switching)
OT:	Offload Transport (An entity that exports reliable transport semantics)
OTI:	Offload Transport Interface
OP:	Offload Protocol (Any upper layer protocol run over OTs. E.g. Sockets Direct Protocol)
SDP:	Sockets Direct Protocol (A Socket emulation protocol specified for InfiniBand)
TOE:	TCP Offload Engine (Hardware that supports offloading TCP/IP protocol from host)
IBA:	InfiniBand Architecture
IPoIB:	IP-over-InfiniBand (and IETF defined RFC to send IP packets on InfiniBand fabric)

1.4 References

InfiniBand

InfiniBand Architecture Specification, Version 1.0a, <http://www.infinibandta.org/>

IP over IB IETF draft: <http://www.ietf.org/ids.by.wg/ipoib.html>

InfiniBand Specification Annex A4 - Sockets Direct Protocol (SDP), Release 1.0.a

Sockets

Fast Sockets reference: <http://www.cs.purdue.edu/homes/yau/cs690y/fastsocket.ps>

Stream Socket on Shrimp reference: <http://www.cs.princeton.edu/shrimp/Papers/canpc97SS.ps>

Device Drivers

Rubini, Alessandro and Corbet, Johathan. Linux Device Drivers Book, 2nd Edition: O'reilly, June 2001. ISBN: 0-59600-008-1. <http://www.xml.com/ldd/chapter/book/>

1.5 Conventions

This document uses the following typographical conventions and icons:

<i>Italic</i>	is used for book titles, manual titles, URLs, and new terms.
Bold	is used for user input (in the Installation section).
Fixed width	is used for code definitions, data structures, function definitions, and system console output. Fixed width text is always in Courier font.



NOTE

Is used to alert you to an item of special interest.



DESIGN ISSUE

Is used to alert you to unresolved design issues that may impact the module's design, function, or usage.

1.6 Before You Begin

Please note the following:

This document assumes that you are familiar with the *InfiniBand Architecture Specification*, which is available from the InfiniBand Trade Association at <http://www.infinibandta.org>.

2. Features

This section lists a set of features and goals for Offload Sockets Framework (OSF) support in Linux. The items listed in this section are by no means complete and may need to be further refined before finalizing on the best solution.

- All offload protocols/transporters need to have a standard Linux network driver. This allows network administrators to use standard tools (like `ipconfig`) to configure and manage the network interfaces and assign IP addresses using static or dynamic methods.
- The offload sockets framework should work with and without kernel patches. To this effect, the offload protocols and transporters will reside under a new offload address family (`AF_INET_OFFLOAD`) module. Applications will be able to create socket instances over this new address family directly. However, for complete application transparency, an optional minimal patch to the Linux kernel (socket driver) can be applied to allow re-direction of `AF_INET` sockets to the new `AF_INET_OFFLOAD` address family. The `AF_INET_OFFLOAD` module will work as a protocol switch and interact with the `AF_INET` address family. The patch also defines a new address family called `AF_INET_DIRECT` for applications that want to be strictly using the OS network stack. This kernel patch can be optional based on distributor and/or customer requirements.
- All standard socket APIs and File I/O APIs that are supported over the OS resident network stack should be supported over offload sockets.
- Support for native [Asynchronous I/O \(AIO\)](#) is being worked in Linux community. The offload framework should utilize this to support newer protocol and transporters that are natively asynchronous. (For example, SDP stack could utilize the AIO support to support PIPELINED mode in SDP)
- Architecture should support a layered design so as to easily support multiple offload technologies, and not just SDP. This insures the offload sockets framework is useful for multiple offload technologies.
- The proposed architecture should support implementations optimized for zero-copy data transfer modes between application buffers across the connection. High performance can be achieved by avoiding data copies and using RDMA support in modern interconnects to do zero copy transfers. This mode is typically useful for large data transfers where the overhead of setting up RDMA is negligible compared to the buffer copying costs.
- The proposed architecture should support implementations optimized for low latency small data transfer operations. Use of send/receive operations incurs lower latency than RDMA operations that needs explicit setup.
- Behavior with signals should be exactly same as with existing standard sockets.
- `listen()` on sockets bound to multiple local interfaces (with `IPADDR_ANY`) on a `AF_INET` socket should listen for connections on all available IP network interfaces in the system (both offloaded, and non-offloaded). This requires the `listen()` call from application with `IPADDR_ANY` to be replicated across all protocol providers including the in-kernel TCP stack.
- Multiplexed I/O operations using API's such as `select()` and `poll()` should work across `AF_INET` socket file descriptors supported by different protocol/transport providers including the

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

in-kernel IP stack. This guarantees complete transparency at the socket layer irrespective of which protocol/transport provider is bound to a socket. .

- Operations over socket connections bound to the in-kernel O/S protocol (TCP/IP) stack should be directed to the TCP/IP stack with minimum overhead. Application bound to kernel network stack should see negligible performance impact because of offload sockets support.
- Ability to fallback to kernel TCP/IP stack dynamically in case of operation/connection failure in direct mapping of stream connections to offloaded protocols/transports. Connection requests for AF_INET sockets that fail over offload stack is automatically retried with the kernel TCP/IP stack. Once a direct mapped connection is established, it cannot be failed back to the TCP stack, and any subsequent failures are reported to application as typical socket operation failures.
- Offload Socket framework enables sockets of type STREAMS only. Other socket types (such as RAW, DATAGRAMS, PACKET etc.) will use only the OS network stack.
- Offload sockets framework will support offloading of stream sessions both within local IP subnet and outside local IP subnet that needs routing. Offload protocols/transports will have the ability to specify if they do self-routing or need routing assistance. Ability to offload stream sessions to remote IP subnet will be useful for TOE vendors in general and for IBA edge router vendors who map SDP sessions on IBA fabric to TCP sessions outside fabric. For protocols/transports that do self-routing, the offload sockets framework simply forwards the requests to them. For protocols/transports that need routing support (such as SDP), the framework utilizes the OS route tables and applies its configurable policies before forwarding requests to offload transports. This enables the use of O/S managed route tables to configure both offload and non-offload stacks.
- Since the socket API extensions defined by the [Interconnect Software Consortium \(ICSC\)](#) in the open group is work-in-progress at this time, the offload sockets framework will not attempt to address them in this phase. This could be attempted at a later phase.
- Offload sockets framework should not affect any existing Linux application designs that uses standard OS abstractions and features (such as `fork()`, `exec()`, `dup()`, `clone()`, etc.). Transparency to applications should be maintained.
- Offload sockets framework should support both user-mode and kernel-mode socket clients and maintain the existing socket semantics for existing user mode or kernel mode clients.
- The offload sockets framework currently deals with only IPv4 address family. Even though the same offload concepts can be equally applied to offload IPv6 family, it is deferred for later stages of the project.

3. Design Overview

The Offload Sockets Framework architecture implements high-performance socket support in the kernel by providing a new Offload Protocol Switch (OPS) via a new address family (AF_INET_OFFLOAD) module. This new address family module will support dynamic binding of offload protocols and transports under the offload address family. The offload protocols will register with the offload family and the transport modules will register with the offload protocol modules. The offload sockets architecture is shown in Figure 1.

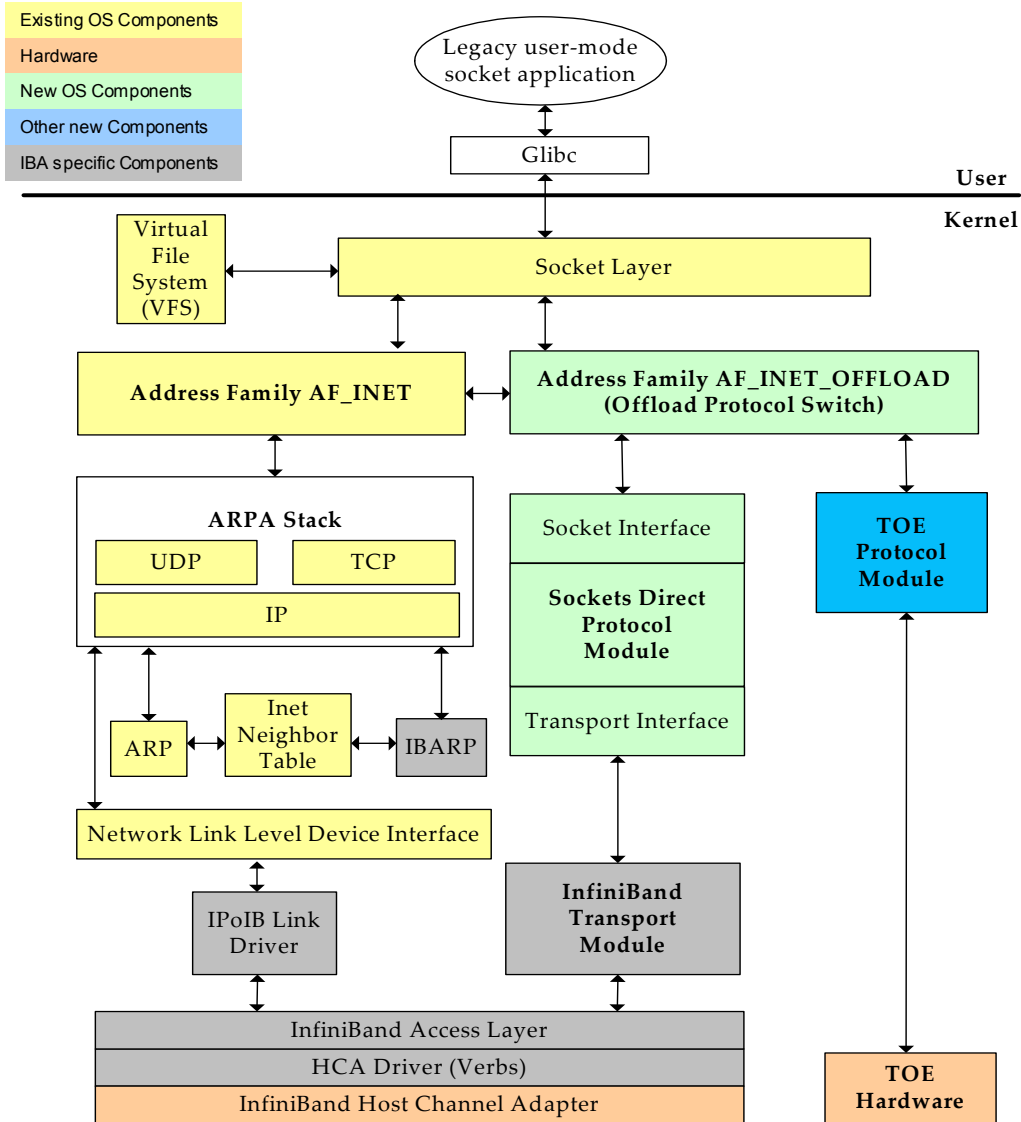


Figure 1: High Performance Sockets in the kernel

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

There are multiple major components to this, such as the Offload Protocol Switch (OPS) module, the SDP offload protocol (OP) module (which include the socket and transport Interface), and the InfiniBand offload transport (OT) module.

4. Design Assumptions & Rules

4.1 Kernel and API

Design is based on Linux 2.4.x kernel feature set. Kernel modifications will be limited as much as possible and will be isolated with compile time switches. Asynchronous I/O kernel support that is currently under development in the Linux community will be monitored and will be supported as part of this design when it becomes available. Both IA32 and IA64 environments will be supported. This project is specific for IPv4 address family but will design with future IPv6 in mind.

4.2 Native AF_INET stack support

OPS design currently assumes that the only native stream protocol type supported by the native stack is IPPROTO_TCP. (Currently, Linux treats a stream socket created with IPPROTO_IP (any) as sockets over IPPROTO_TCP). However, OPS can be extended to span across multiple native stream protocol types (such as IPPROTO_SCTP) if required.

4.2.1 OPS and Protocol Module loading

It is expected that the AF_INET_OFFLOAD switch and all offload protocol modules be loaded immediately following the native TCP stack. Load order restrictions are the result of sharing port space with native AF_INET stack. AF_INET_OFFLOAD protocol registration will fail if an offload protocol, with the same protocol id supported under the native stack (IPPROTO_TCP), attempts to register after the first create socket call. A protocol that does not share a protocol id with the native stack will be allowed to register, even after the first create socket call. Both types of protocol modules can unload at anytime. The load ordering can be controlled via rc3.d init processing.

4.2.2 Protocol ID's

Offload protocols do not need to have unique protocols IDs. There could be multiple protocol modules registered with the switch using the same protocol id. Protocols supported by offload protocol modules could also overlap with protocols (such as IPPROTO_TCP) registered with the native AF_INET stack.

4.2.3 Protocol Interfaces

The local interfaces (IP addresses) exposed by the protocol modules registered with the switch may overlap or not. But each of the local interfaces exposed by an offload protocol modules need to be exposed by a respective link driver to the native AF_INET TCP/IP stack.

If a local interface is shared across protocol modules (offloaded or not), the OPS will decide which one of them will be picked for a specific socket. This could be done by the switch in the order of notification of local interface UP events by the individual protocol modules, with the native protocol under AF_INET always selected at the end of the search.

4.2.4 IP port space

If the same protocol exists under AF_INET and AF_INET_OFFLOAD families, the port space for that protocol is shared. For example, a TOE protocol module could register IPPROTO_TCP with the switch (AF_INET_OFFLOAD) and if we try to do a bind on a socket (created with protocol of IPPROTO_TCP or IPPROTO_IP (any)), the bind is first done over the native AF_INET TCP stack (to get/confirm the port number) and the same port number is locked down across all other offload protocols that match the address specification in the bind.

If the socket was created specifying a non-native protocol, the bind operation on the socket is spanned (for port space unification) only across those offload protocol modules that support this protocol and match the address specification in the bind. For example, if there are two offload modules M1 and M2 registered with OPS that expose the non-native protocol IPPROTO_SDP, a bind operation with wildcard IP (IPADDR_ANY) on a socket created with protocol type IPPROTO_SDP will span both these protocol modules. The bind will be first applied on M1 (assuming M1 registered first with OPS) to get/confirm the port number and the same port number will be requested in the bind forwarded to M2. Any module that fails the bind will be removed from the protocol list for this socket.

If a bind specifies a specific local interface, the bind operation is forwarded only to those offload protocol modules that registered the specified local interface with the OPS. If bind specifies wildcard local interface (IPADDR_ANY), then the bind is forwarded to all offload protocol modules currently registered with OPS).

If a protocol module fails an operation on behalf of a socket, the failed protocol is no longer included in the protocol list for that socket (for bind and create it is not added to list, for others it would be removed from the list).

4.2.5 Generic socket options

At the socket creation time, if a native socket exists, the OPS will get all generic socket option values from the native sock structure and save them in the offload sock structure. Any get/set sockopt calls from the application for generic socket options are directly returned from this offload sock structure. OPS will provide an up-call for offload protocol modules to get the generic socket level option values at protocol initialization time.

4.2.6 Existing IP protocol specific socket options

Any offload protocol that exposes a native protocol (protocols exposed by native AF_INET stack) or that wish to operate under the IPPROTO_IP wildcard, **must** support all of the protocol specific options currently supported under the native AF_INET stack.

Also at the socket creation time, if a native socket exists, the OPS will get the default protocol specific options from the native socket and will attempt to set these default protocol options over all offload protocols. The offload protocol modules are expected to support these protocol specific options and default values. If setting the default option on an offload protocol module fails, it is removed from the list of OPS supported offload protocol modules.

If a native socket exists, any subsequent setsockopt calls for protocol specific options from the application will be first applied over the native socket. Any error on the set operation on the native sock is treated as failure and returned as error to the application. If the native set succeeds, the set

operation is spanned across all the offload protocol modules. If any of the offload protocol modules fail the set operation, OPS will take it off the protocol list for this socket.

If a native socket exists, all getsockopt operations from the application for protocol specific options are simply returned from the native socket values.

4.2.7 **New offload protocol specific socket options**

If multiple offload modules exist that support a common non-native protocol (such as IPROTO_SDP), the non-native protocol specific set option operation results should always match across the modules. Also, for protocol specific getsockopt operations, the returned values must match across modules. Any mismatch in get or set operations will result in OPS returning an error to the get or set sockopt calls from application.

5. Design Details

5.1 Offload Protocol Switch Module

The Offload Protocol Switch (OPS) module exposes a new address family (`AF_INET_OFFLOAD`) and preserves the socket operation semantics with the kernel socket driver at the top. At the bottom it provides a new offload protocol structure and interface for offload protocols. These offload protocols will register with the `AF_INET_OFFLOAD` address family during initialization using an exported registration call.

The OPS module is logically a thin veneer between the socket driver and the INET offload protocols. OPS internal interfaces will be defined to provide routing information for the offload protocol modules and transport interfaces. All external interfaces to `AF_INET_OFFLOAD` are defined by the Linux socket calls in `include/linux/net.h` and `include/linux/sock.h`.

This module is capable of switching socket traffic across multiple offload protocol modules underneath via supplied protocol address information. The module will also be capable of failing back to standard `AF_INET` stack if failures occur during connection initiation over a offload protocol stack. This switching to the standard stack requires no changes to the `AF_INET` stack.

5.1.1 Linux Kernel Modifications

The OPS module will be coded as a new Linux `AF_INET_OFFLOAD` address family. This requires new definitions, `AF_INET_OFFLOAD` and `AF_INET_DIRECT`, in the kernel `linux/include/linux/in.h` include file.

5.1.2 Initialization

The OPS kernel module will be demand loaded during system initialization. The OPS module will have a dependency on the standard Linux socket driver. Actual OPS initialization is performed in accordance with the standard Linux kernel module load procedure; the `init_module()` function is called once the OPS module has been successfully loaded into the kernel. See section 7.1 for OPS initialization process and interaction with the offload protocol modules.

5.1.3 Shutdown

OPS shutdown can be initiated in one of two ways:

- During normal system shutdown procedures.
- The OPS module is forcibly removed using the Linux system administration command `rmmod` (remove module).

With either method, the OPS module unload function `cleanup_module()` is called per the standard Linux kernel module unload procedure. The OPS `cleanup_module()` function will notify the protocol modules of the pending OPS module shutdown. It is expected that all offload protocol modules will shutdown by releasing allocated system resources, halting and unloading. See section 7.1 for OPS shutdown process and interaction with offload protocol modules.

5.1.4 Socket Driver Modifications

This section describes the Linux kernel patch envisioned to support the offload protocol switch (OPS). The OPS is packaged as a separate address family (AF_INET_OFFLOAD) module. This allows OPS to be used without any Linux kernel patches, as long as applications explicitly specified the AF_INET_OFFLOAD address family during socket creation. However, this solution is not suitable for applications that need binary/legacy compatibility (uses only AF_INET family sockets) and flexibility to run over offloaded or non-offloaded stacks without modifications. A minor Linux socket driver patch achieves the binary backward compatibility of applications.

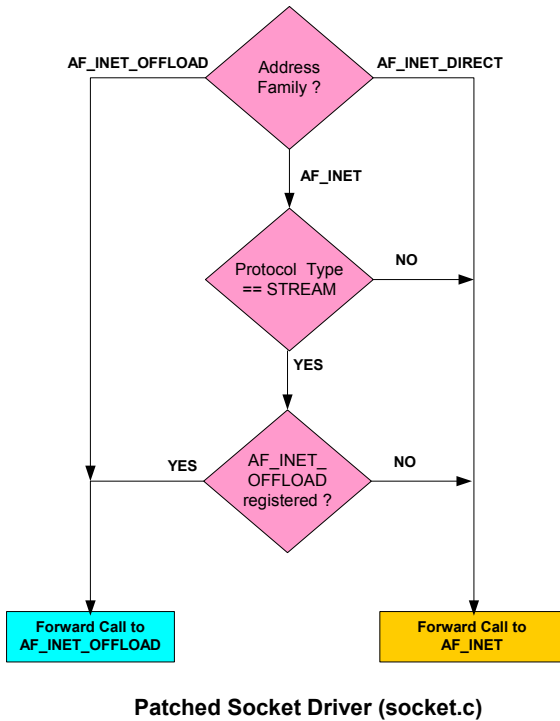
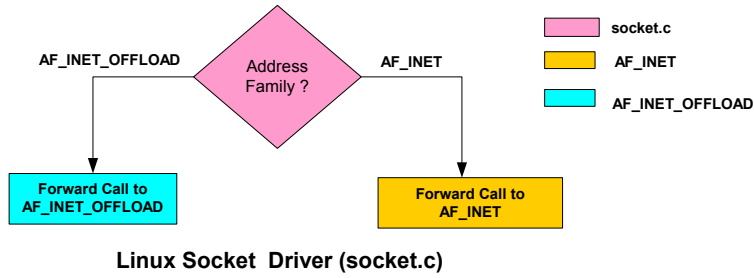


Figure 2: Linux Socket Driver Modifications

The current Linux socket driver address family switching logic is shown above, followed by the patched socket driver logic. With the socket driver patch, whenever a AF_INET socket is created by an application, if it is of type SOCK_STREAM and if the OPS module is loaded and registered with the socket driver, the socket creation calls are forwarded to the OPS module. The patch also exposes as new address family called AF_INET_DIRECT if specific applications wanted to specify during socket

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

creation not to attempt offloading them. Section 5.1.6 explains another reason to expose the AF_INET_DIRECT address family to avoid endless loop in the socket driver.

5.1.5 Offload Data Structures

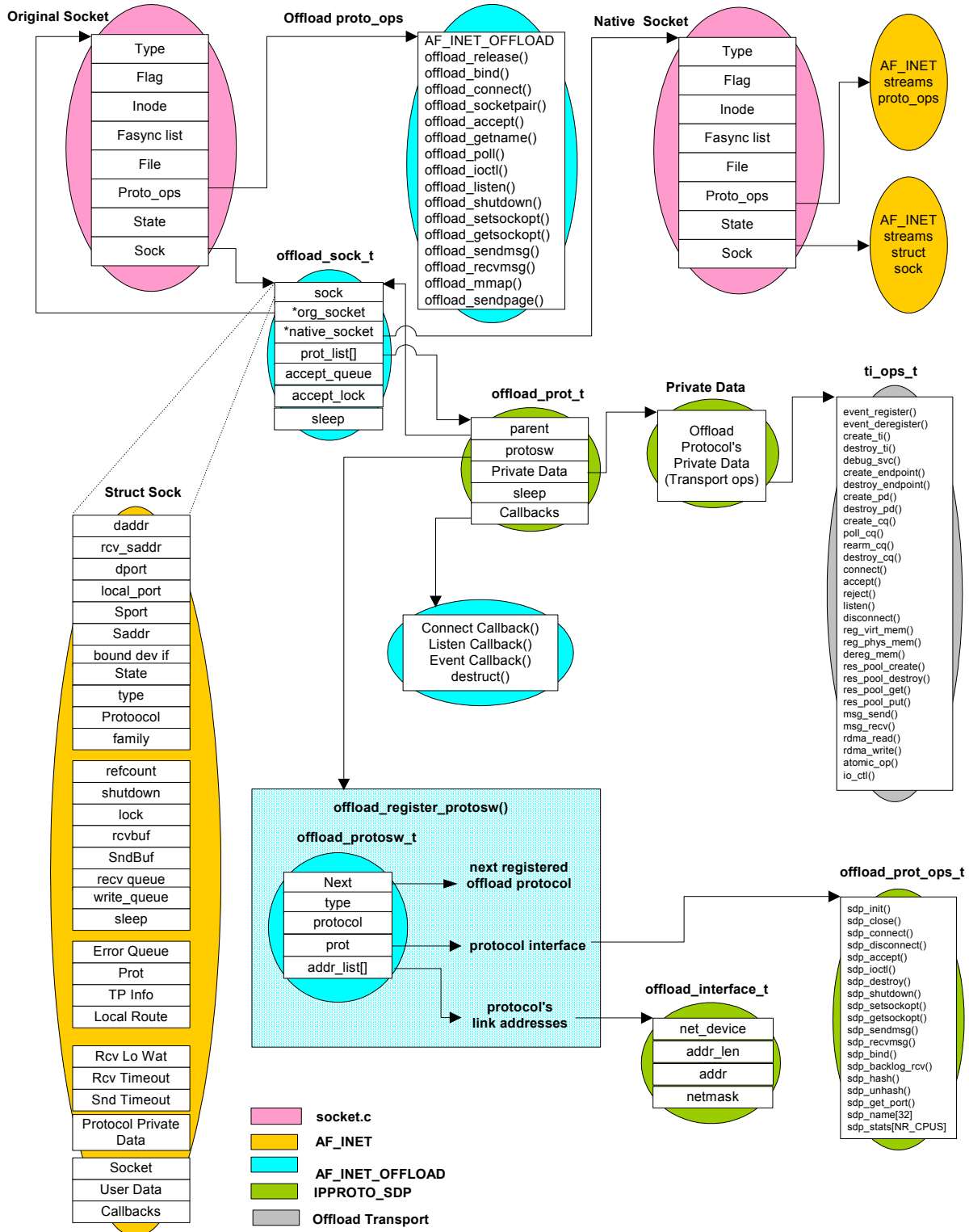


Figure 3: Socket structures for offload support

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

This section attempts to describe the offload data structures and how they relate with the existing Linux `socket` and `sock` data structures. Linux networking stack utilizes the `socket` and `sock` data structures for all socket management functions across the various layers. The `socket` structure is typically created by the socket driver and the `sock` structure is created by the address family driver.

Typically a `socket` structure has a single `sock` pointer in it which points to the address family and protocol specific operations on that socket. The primary reason for this is that during socket creation the address family and protocols are well specified which allows the socket to bind itself during socket creation to a particular protocol. However with offload socket framework, to deal with `IPADDR_ANY` and `IPPROTO_ANY` conditions where a specific protocol or local endpoint is not specified upfront, there is need for multiple offload protocols to be linked to the same socket until enough information comes (during bind or connect) to bind the socket to a specific protocol. To fulfill this requirement the `socket` structure is transparently modified (by adding more members to the end of the buffer pointed by the `sock` member element). Thus a socket created by calling the OPS module has the `socket->sock` field point to an `offload_sock_t` structure. The `offload_sock_t` structure starts with a `sock` element followed by other elements. The `offload_sock_t` contains an `prot_list` which points to a linked list of `offload_prot_t` structures each representing a offload protocol module registered with the OPS. The `offload_prot_t` structures keep all the data and state needed by the offload protocols and transports.

Since the sockets created by the OPS can span over multiple protocols (including the `AF_INET TCP` protocol), a higher level accept queue is maintained at the `offload_sock_t` level. The `offload_sock_t` also points to a 'Native socket' which points to a fully qualified socket that is created using the same credentials as the 'Original socket' but over the standard Linux `AF_INET` address family.

For specific details of the `offload_sock_t` structure, please refer the accompanying data structure definitions in section 12.

5.1.6 Socket Create

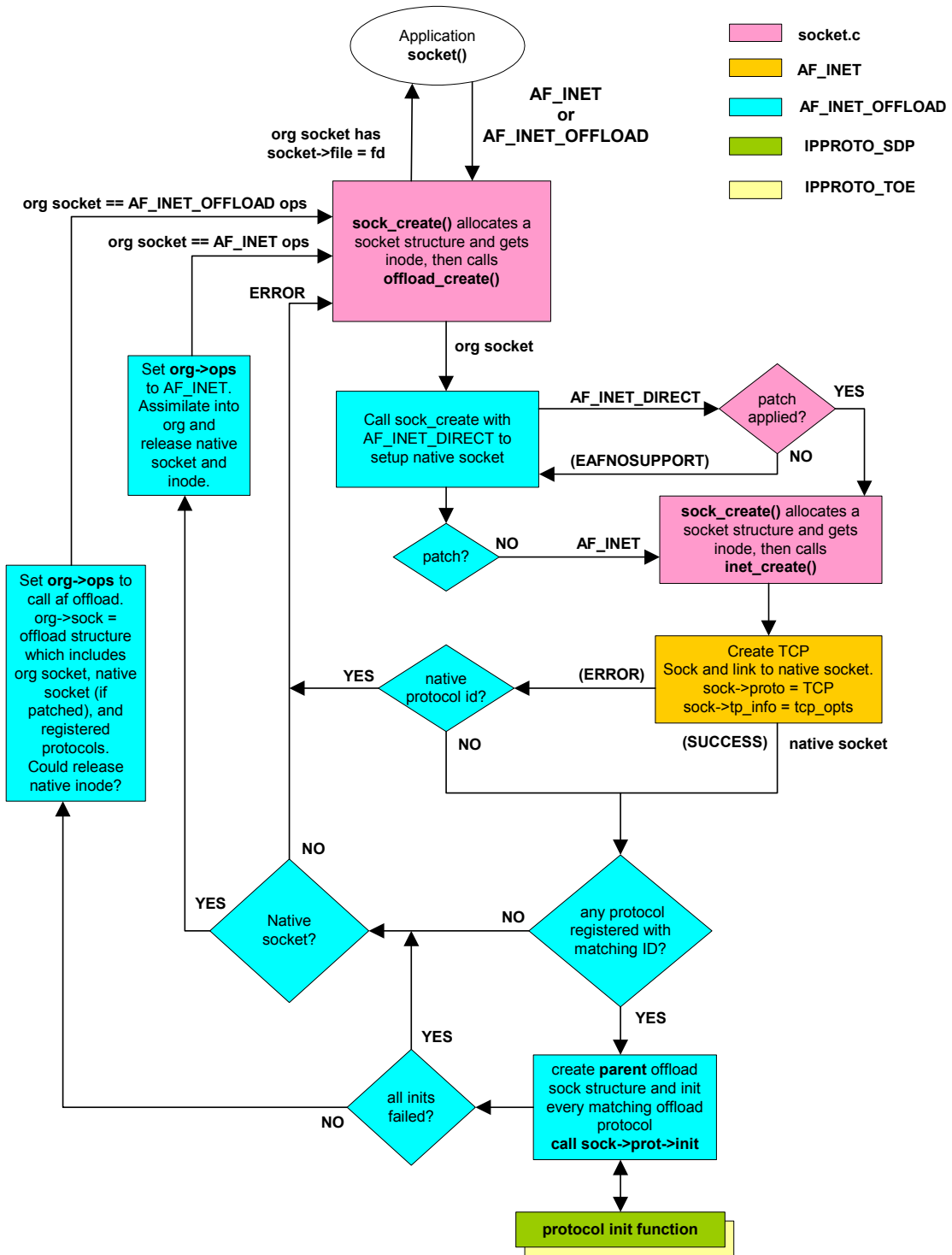


Figure 4: Socket Create

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

This section illustrates the socket creation process using the offload sockets framework. The socket driver `sock_create()` function creates the 'Original socket' structure and forwards the call to the OPS module by calling its `offload_create()` entry point in the `proto_ops` table of the socket. The OPS module calls back the socket driver's `sock_create()` function with the `AF_INET_DIRECT` address family to create a 'Native socket' over the `AF_INET` stack. If the native socket creation was successful, it implies that the socket driver patch is running and the application specified protocol is something supported by the OS `AF_INET` stack. If the native socket creation fails, returns an "address family not found" error, the switch will call `sock_create()` function again with `AF_INET` address family. This is to cover the case where the protocol specified is of native `AF_INET` type (ie `IPPROTO_TCP`) and the port space must be shared. If the native socket creation returns any other error and the protocol specified is of native `AF_INET` type then OPS returns an error to the application.

OPS checks if the specified protocol in the `sock_create()` call is registered with OPS. If a match is not found, the specified protocol is not offload able and the native socket is copied to the original socket to fulfill the socket creation request. If the native socket was not available, then the socket creation call is failed and returned. For each matching protocol, the offload protocol's initialization function is called. `IPPROTO_IP` is the any protocol case and matches all registered protocols regardless of the id. If all protocol modules fail initialization, then the OPS attempts a fail back to the `AF_INET` protocol stack if a native socket was available by assimilating into the original socket, releasing the native socket, and setting the socket `proto_ops` to the native stack. If the protocol id matches any registered offload protocol then the socket `proto_ops` is set to the `proto_ops` defined by the OPS.

Based on the above conditions, the `offload_create()` call from the socket driver could return with either `org_socket->proto_ops` pointing to protocol operation call table of OPS, or protocol operations of `AF_INET`, or a NULL indicating error. In any case, this is completely transparent to the socket driver and the application above it.

5.1.7 Socket Bind

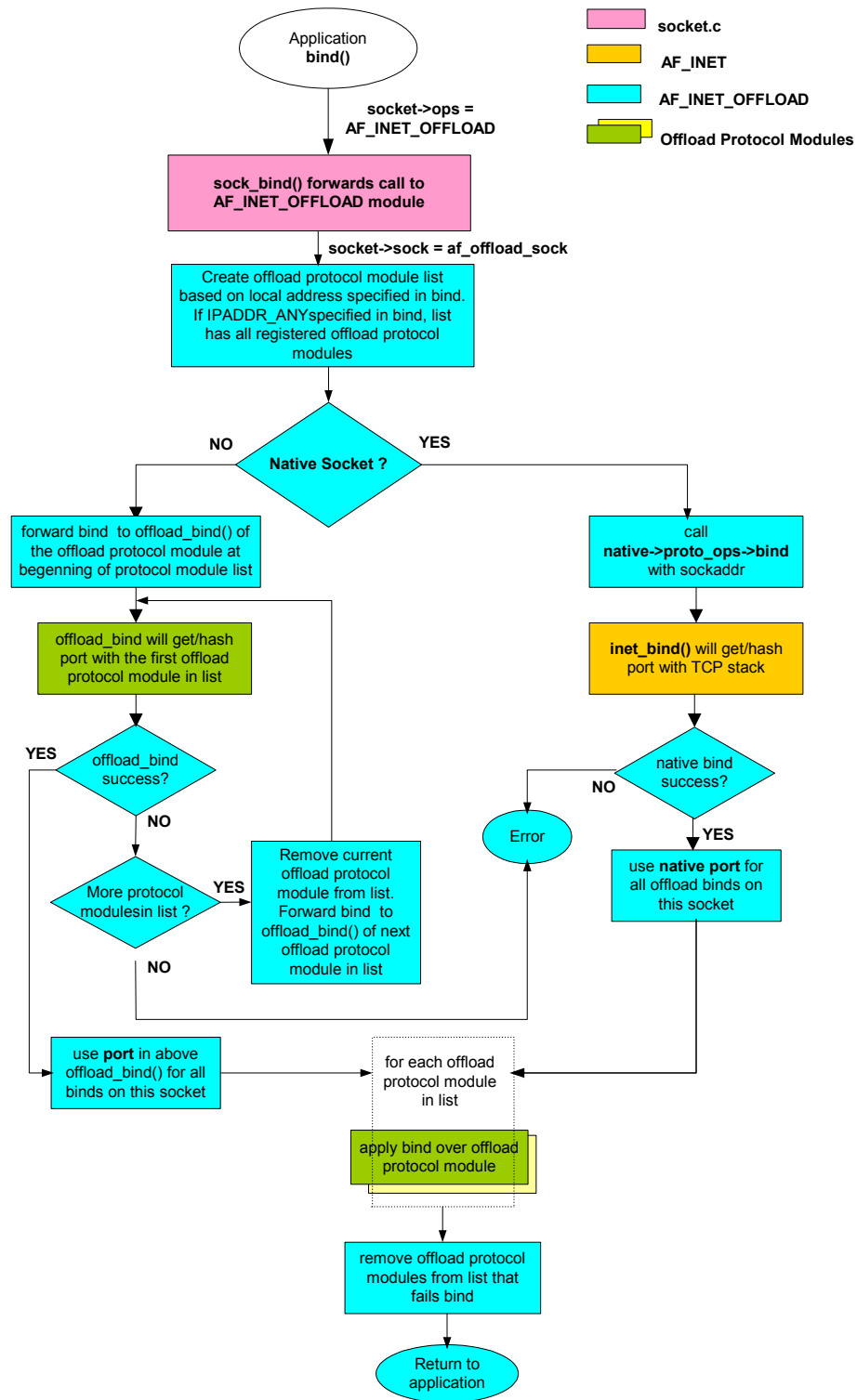


Figure 5: Socket bind

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

This section illustrates the program flow in the OPS for a socket bind operation. Since the socket creation request was handled by OPS (as described in the previous section), the `socket->proto_ops` points to the OPS (AF_INET_OFFLOAD module's) function call. The socket driver simply forwards the bind operation to the AF_INET_OFFLOAD address families bind entry point.

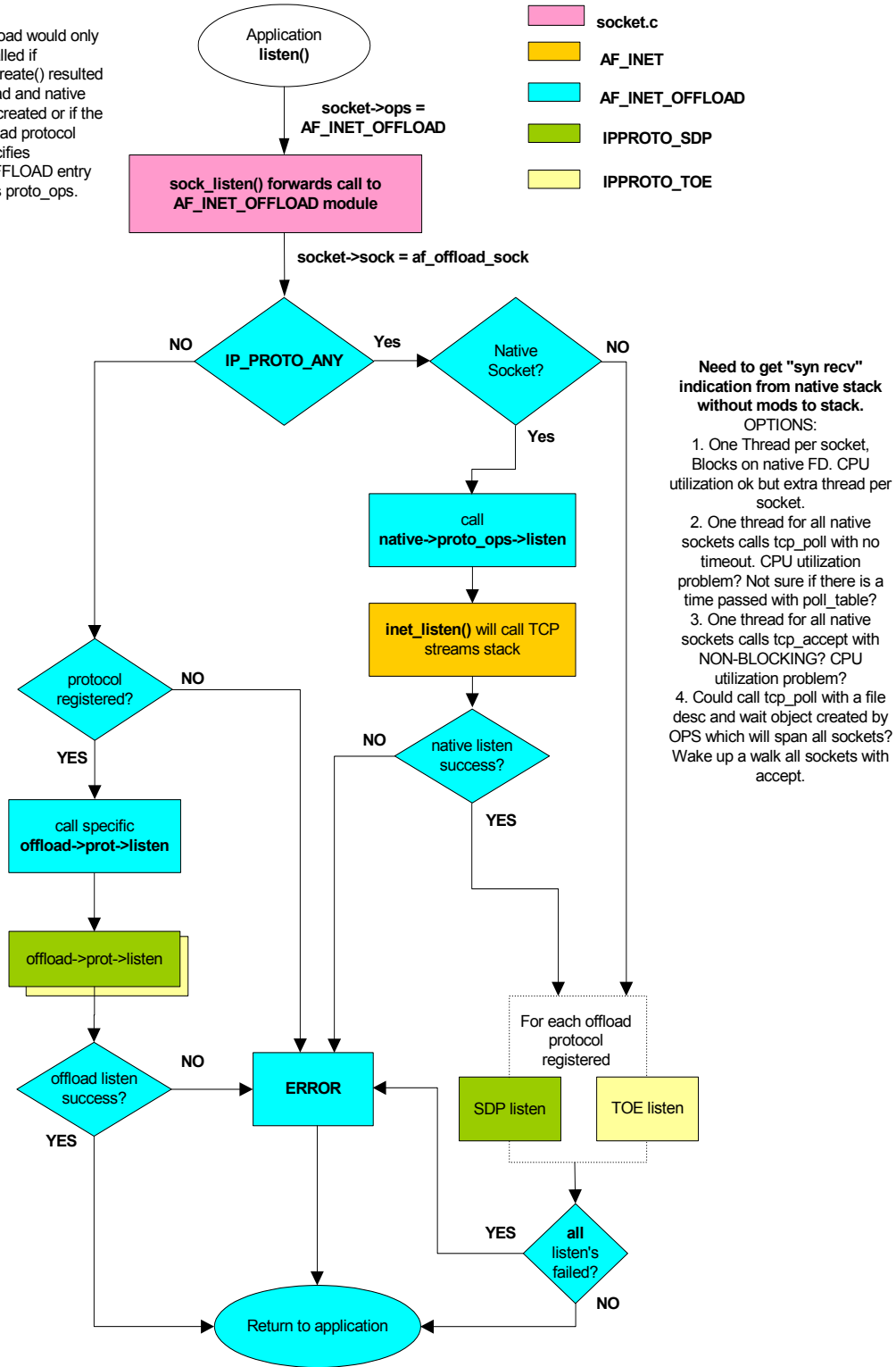
Socket bind operation allows clients to use a wildcard local address (INADDR_ANY) or specify a local IP address. Once the OPS module's bind entry point gets called by the socket driver, the OPS module prunes the list of protocol modules attached to the socket (during socket creation time) based on the local address specified in bind. If a INADDR_ANY is specified the protocol module list for the socket is left untouched. If a specific local address is specified, all protocol modules that do not support the specified address are removed from the protocol module list for this socket.

Once the protocol module list is refined, OPS checks if the application socket has a native socket associated with it. If a native socket exists for this socket, the bind is first invoked over the native socket. The bind operation on a native socket is handled by the AF_INET stack and the TCP/IP protocol suite registered with it. The native socket bind is done first to reserve the application specified port number (or get an unused port number if application uses dynamic ports) from the native AF_INET stack. If the native socket bind fails, an error is returned back to the application. If the native socket bind succeeds, the port number from this native bind is specified in bind operation forwarded to each of the offload protocol modules attached to this socket. This allows the port space to be unified across all offload protocol modules and the native AF_INET stack. Offload protocol modules are supposed to keep track of the endpoints (port numbers and addresses) that are in use and fail the bind requests if they are already in use. Any offload protocol module that fails the bind operation is removed from the protocol module list for this socket.

If there was no native socket associated with the application socket, then the port space is unified only across the set of offload protocol modules attached to this socket. In this case, one of the offload protocol modules from the list is selected (possibly in the order of their registration with OPS) and the bind operation executed on it to reserve/get a port number. If the bind operation fails, this offload protocol module is removed from the list and retried on the next offload protocol module in the list. Once a successful bind operation is completed on an offload module, the same port number is used for binds on all other protocol modules in the list for this socket. Similar to the case when native socket is present, any offload protocol module that fails the bind operation is removed from the protocol module list for this socket.

5.1.8 Socket Listen

NOTE: AF offload would only be called if `offload_sock_create()` resulted in both offload and native sockets being created or if the bound offload protocol specifies `AF_INET_OFFLOAD` entry points in its `proto_ops`.



Need to get "syn rcv" indication from native stack without mods to stack.
 OPTIONS:
 1. One Thread per socket, Blocks on native FD. CPU utilization ok but extra thread per socket.
 2. One thread for all native sockets calls `tcp_poll` with no timeout. CPU utilization problem? Not sure if there is a time passed with `poll_table`?
 3. One thread for all native sockets calls `tcp_accept` with NON-BLOCKING? CPU utilization problem?
 4. Could call `tcp_poll` with a file desc and wait object created by OPS which will span all sockets? Wake up a walk all sockets with accept.

Figure 6: Socket Listen

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

OPS handles socket listen very similar to how bind is handled. Listens on sockets that are bound to a specific protocol are forwarded to the protocol module handling the specific protocol. For wildcard protocol bound sockets, the listens are posted on all the available protocol modules bound to the socket, including TCP if a 'Native socket' is attached to the 'Original socket'. Since TCP interface uses blocking semantics for operations such as accept(), poll() etc., the OPS module utilizes a worker thread to process these operations. Also, a global accept queue is constructed at the OPS module to process incoming connection requests through any of the protocol modules, including TCP. The accept/poll section describes details on these operations.

5.1.9 Socket Connect

NOTE: AF offload would only be called if offload_sock_create() resulted in both offload and native sockets being created or if the bound offload protocol specifies AF_INET_OFFLOAD entry points in its proto_ops.

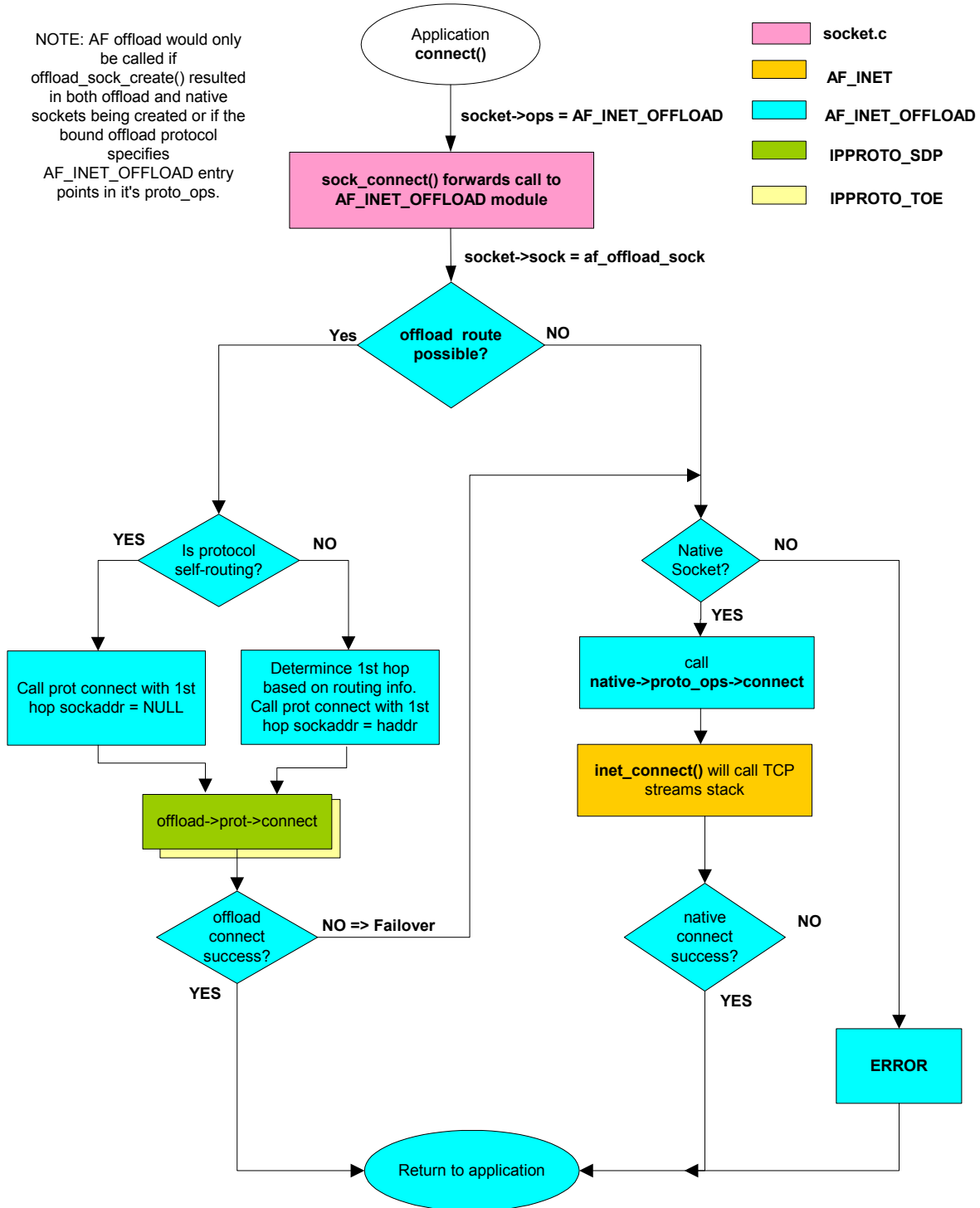


Figure 7: Socket Connect

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

This section illustrates the program flow in the OPS for a socket `connect ()` operation. Since the socket creation request was forwarded to the OPS module by the socket driver simply forwards the `connect ()` operation to the `AF_INET_OFFLOAD` address families `connect ()` entry point as exposed in the `socket->proto_ops` member.

Once the OPS module's `connect` entry point gets called by the socket driver, it checks it's current routing table and based on the `sockaddr` passed with the `connect`, routes to the proper offload protocol. If no routes are found via registered offload protocol providers then the OPS module will forward the `connect` request to the native stack. Otherwise, it will check to see if the protocol provider needs first hop information and will pass the appropriate 1st hop and final destination address information via the protocols `connect` entry point. Any errors returned by the offload protocol stack or by the native protocol stack will be returned to the application.

5.1.10 Socket Accept and Poll on a Listening Socket

NOTE: AF offload would only be called if offload_sock_create() resulted in both offload and native sockets being created

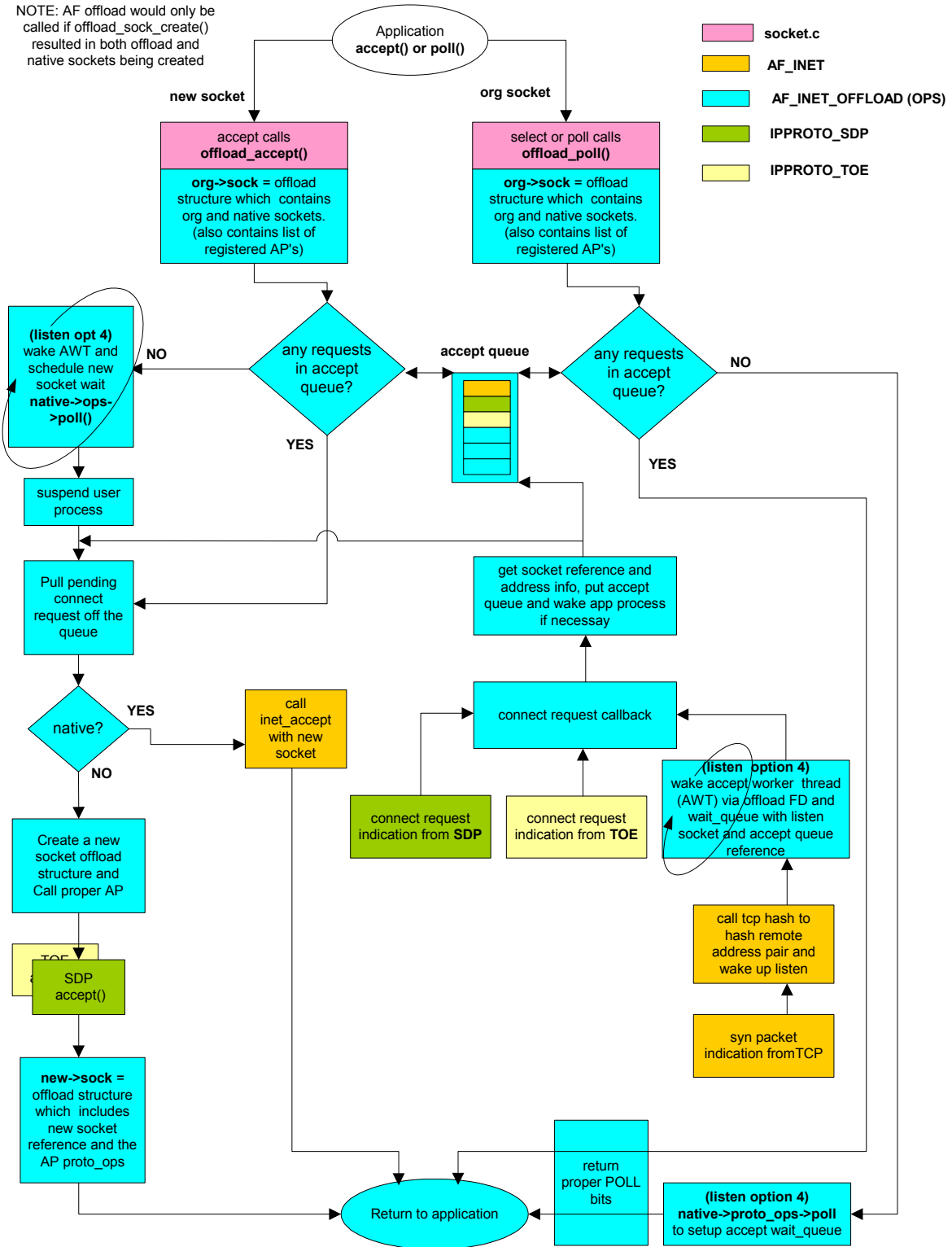


Figure 8: Socket Accept & Poll

This section illustrates the program flow in the OPS for a socket `accept()` and `poll()` operation. Since the socket creation request was forwarded to the OPS module by the socket driver, the `socket->proto_ops` simply forwards the `accept()` or `poll()` operation to the `AF_INET_OFFLOAD` address families `accept()` or `poll()` member respectively.

The OPS poll entry point is called as a result of either a select or poll at the application level. The kernel `do_select` is called to handle any timeout or multiple listen FD's specified by the application. This diagram covers the details of listen FD's only. The OPS will check for listen state and for any pending requests on the accept queue. If there are accept requests pending, it will return the proper POLL bits to the application to indicate the pending connection. Otherwise, it will wake up the Accept Working Thread (AWT) to setup the native accept wait_queue to process native connections (SYN). Any errors returned by the native protocol stack will be returned to the application.

The OPS accept entry point is called as a result of a application accept on a listen FD after a poll or select has indicated a connection request event on a previous listen. The OPS will check for listen state and for any pending requests on the accept queue. If there are accept requests pending, it will pull the accept information off of the accept queue. Otherwise, it will wake up the Accept Working Thread (AWT) to setup the native accept wait_queue and will suspend the user process.

If the accept request is a native request, then the native accept is called with the new socket. Otherwise, if the accept request is an offload protocol request, the OPS will create a new socket offload structure and call the appropriate offload protocol accept. Once the offload protocol `accept()` returns, the OPS will setup the `proto_ops` and the new socket reference and return to the application. Any errors returned by the offload protocol stack or by the native protocol stack will be returned to the application. .

5.1.11 Socket Operations on a Connected Socket

Once a socket transitions into a connected state, it is also bound to a specific protocol that fulfilled the connection setup. Once a socket is connected the operations on this socket are forwarded blindly to the protocol module it is bound to through the `proto_operations` exposed by the protocol modules. The OPS module could also try to remove itself from the function call path of a connected socket by pointing the `proto_ops` member of the socket structure to directly point to the `offload_proto_ops_t` structure of the bound protocol underneath at the time of connection setup completion. Any data transfer operations (`send/rcv/poll/select`) are handed exactly same as how Linux stack handles them currently. If the socket is connected through the Linux `AF_INET` TCP stack, the OPS module makes the original socket resemble a native socket with the `proto_ops` structures pointing to the `AF_INET/TCP` protocol operation call tables. With this model, the socket driver and application above it always sees a transparent interface.

5.1.12 Socket Options Processing

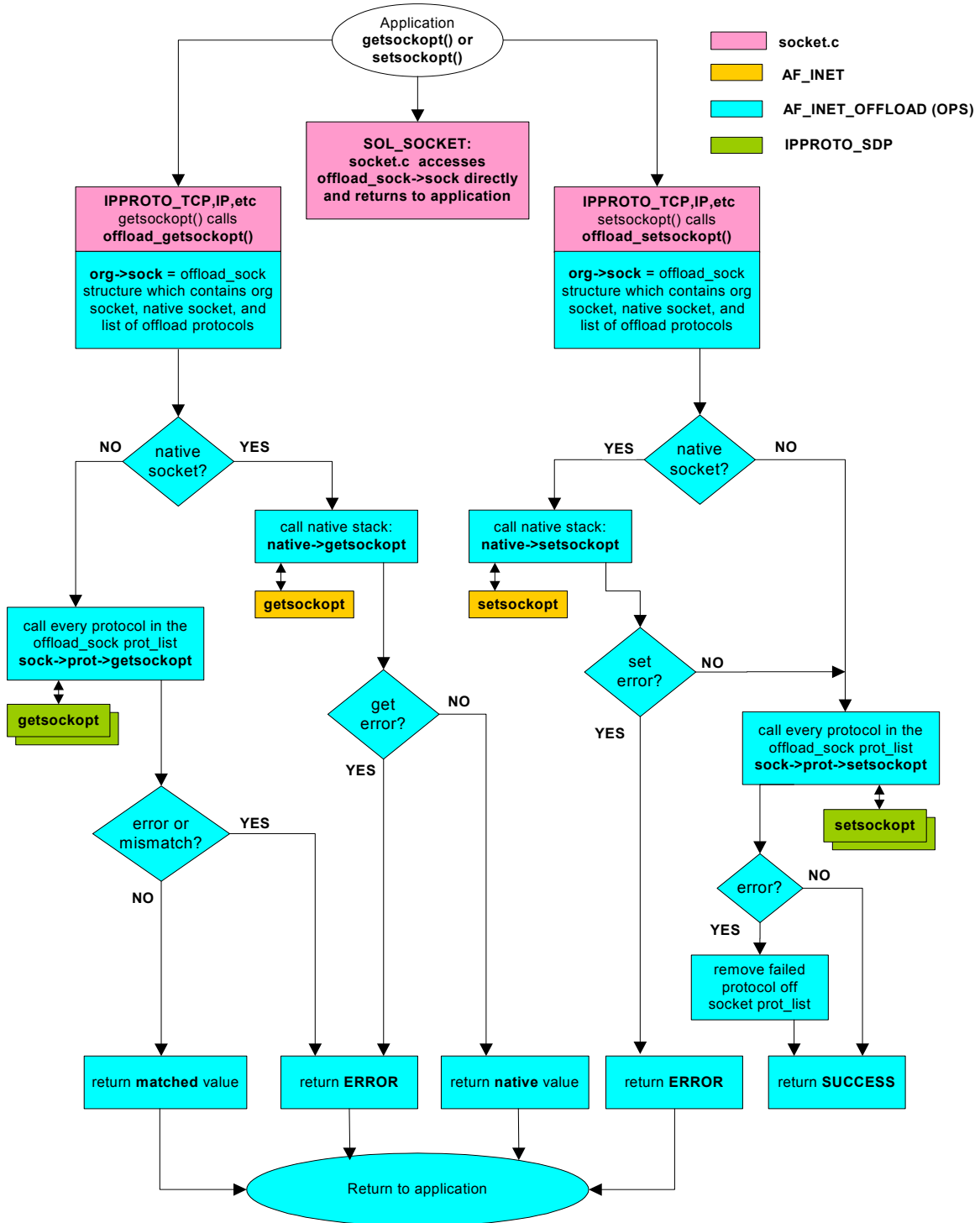


Figure 9: Socket Options

IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design

At the socket creation time, if a native socket exists, the OPS will get all generic socket option values from the native sock structure and save them in the offload sock structure. Any get/set sockopt calls from the application for generic socket options are directly returned from this offload sock structure. OPS will provide an up-call for offload protocol modules to get the generic socket level option values at protocol initialization time.

Any offload protocol that exposes a native protocol (protocols exposed by Linux AF_INET) or that wish to operate under the IPPROTO_IP wildcard, must support all of the protocol specific options currently supported under the native AF_INET stack.

Also at the socket creation time, if a native socket exists, the OPS will get the default protocol specific options from the native socket and will attempt to set these default protocol options over all offload protocols. The offload protocol modules are expected to support these protocol specific options and default values. If setting the default option on an offload protocol module fails, it is removed from the list of OPS supported offload protocol modules.

If a native socket exists, any subsequent setsockopt calls for protocol specific options from the application will be first applied over the native socket. Any error on the set operation on the native sock is treated as failure and returned as error to the application. If the native set succeeds, the set operation is spanned across all the offload protocol modules. If any of the offload protocol modules fail the set operation, OPS will take it off the protocol list for this socket.

If a native socket exists, all getsockopt operations from the application for protocol specific options are simply returned from the native socket values.

If multiple offload modules exist that support a common non-native protocol (such as IPPROTO_SDP modules), protocol specific set option operation results should always match across the modules. For protocol specific get option operations, the returned values must match across modules. Any mismatch in get or set operations will result in OPS returning an error to the get or set sockopt calls from application.

5.2 Sockets Direct Protocol Module

The Linux implementation of SDP includes many interdependent components that are referred to as Offload Sockets Framework (OSF). This framework will allow applications to bypass the resident TCP/IP protocol stack while using the default address family of AF_INET.

Offload Sockets Framework includes an Offload Protocol Switch (OPS) module (Ipv4 internet protocol switch) that allows integration of Offload Protocol (OP) modules along side the existing TCP/IP stack, an OP module that supports SDP on the wire, and an InfiniBand Offload Transport (OT) that provides the transport abstraction to an InfiniBand fabric.

There are 3 major internal components of the SDP module, socket interface layer at the top, the SDP state machine in the middle, and the transport interface at the bottom. The SDP module exposes the standard `proto_ops` call interface at top. In addition to the standard `proto_ops` call interface at the top, SDP also supports the new extended proto operations, `offload_ops`, as defined by the OPS interface. These new interfaces include a new connect call that supports a first hop IP address and a new sock structure that defines specific offload protocol and transport details. This section covers the design details of the Socket Direct Protocol Module.

5.2.1 Linux Kernel Modifications

The SDP module will be coded as a new Linux AF_INET network protocol type. This requires a new definition, `IPPROTO_SDP`, in the kernel `linux/include/linux/in.h` include file. In addition to the new SDP protocol type, the sock structure in `linux/include/net/sock.h` will need to be modified to include a new `sdp_opt` structure qualified with `"#if defined (CONFIG_SDP)"`

5.2.2 Initialization

The SDP kernel module will be demand loaded after the initialization of OPS, normally during system initialization. The SDP module will have a dependency on OPS. Offload transports will have a dependency on SDP.

Actual SDP initialization is performed in accordance with the standard Linux kernel module load procedure; the `'init_module()'` function is called once the SDP module has been successfully loaded into the kernel. See section 7.1 for SDP initialization process and interaction with the OPS and the OT.

5.2.3 Shutdown

SDP shutdown can be initiated in one of two ways:

- During normal system shutdown procedures.
- The SDP module is forcibly removed using the Linux system administration command `'rmmod'` (remove module).

With either method, the SDP module unload function `"cleanup_module()"` is called per the standard Linux kernel module unload procedure.

The SDP `"cleanup_module()"` function will notify the OPS and the offload transports of the pending SDP module shutdown. It is expected that SDP will shutdown by releasing allocated system resources, halting and unloading itself. This will result in a cleanup of all INET addresses in use by the OPSI on

behalf of the SDP module. See section 7.1 for SDP shutdown process and interaction with the OPS and the OTI.

5.2.4 Buffer Strategy

SDP supports several data transfer mechanisms designed to take advantage of reliable transports with RDMA capabilities. Four types of transfer modes are supported with SDP; buffered, combined, write zero copy, and read zero copy. Depending on the application workload and buffer availability, SDP will use the most optimal mode for transferring the data.

Small transfers, determined by a data copy size threshold, will use buffered mode to transfer the data. This mode requires a pre-allocated buffer pool that the application data is copied to/from during each data transfer. This buffer pool is provided by the offload transports but will be sized and managed by SDP on a per connection basis. The transport must provide a mechanism to create/destroy these pools with buffer size, count, and data segments. A free pool queue will be provided by the transport, with a low overhead get/put interface, to retrieve and replace buffer resources. The transport must also provide interfaces, to send buffers, to pre-post receive buffers, and to poll or get callback indications when buffer transfers are complete.

Large transfers, greater than the data copy size threshold, that have available read and write buffers on each end of the connection will use the zero copy read or write mechanism to transfer the data. This zero copy data transfer mode requires a mechanism to initialize RDMA operations with the transport. The transport must provide a mechanism to pre-allocate RDMA transport descriptors that can be used for the data operation to avoid unnecessary overhead in the transfer mode. The transport must also provide interfaces to register and un-register user memory, to initiate RDMA reads/writes, and to poll or get callback indications when transfers are complete.

Figure 10 illustrates the buffered mode and the zero copy paths using a reliable transport:

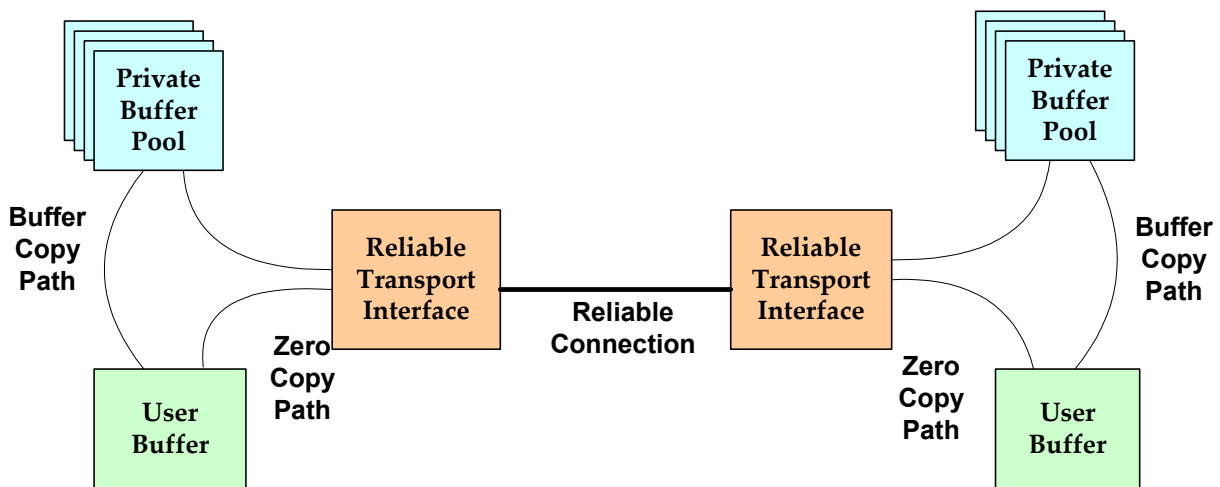


Figure 10: SDP Buffer Mode Overview

Refer to the Data Transfer Mechanism section in the InfiniBand annex A4 Socket Direct Protocol Specification for more details.

5.2.5 Connection Services

Sockets Direct wire protocol defines a standard message protocol that includes a Base Socket Direct Header (BSDH) with all messages. Each message is identified with a Message Identifier (MID) and in addition to BSDH may contain specific message information and actual upper level protocol (ULP) data. For connection establishment the Hello Message (HM) and HelloAck Message are used and to disconnect or abort the DisConn and the AbortConn messages are used. Refer to the SDP Message Formats section in the InfiniBand annex A4 Socket Direct Protocol Specification for more details.

The SDP module is totally transport independent and expects the transport to provide the address resolution and port mapping based on its underlying link. The offload transport interface will provide a mechanism to resolve the endpoints based on a standard IP address/port pairs and return an opaque endpoint data handle back to SDP that can be used for either the listen or connect request.

The transport is also required to provide a mechanism to pass private data (up to 76 bytes) with a connection, listen, accept, and reject calls. These private data buffers are used by SDP to setup connection attributes as defined by the Hello Message (HM) and the HelloAck Message (HAH). The connection reject codes (consumer reject) are exported to the SDP module as specified by the InfiniBand annex A4 Socket Direct Protocol Specification.

InfiniBand's IP over IB mapping and CM REQ does not provide enough granularity to support a listen that multiplexes based on a link address, network address, and a port. It only provides a mechanism to listen on the link address (GID) and the service identifier (mapped IP port). Therefore, the transport interface must supply additional multiplexing fields in the listen call that includes an offset to private data and a size of private data so that SDP can provide the local IP address as the additional multiplexing field. Without this feature, SDP would be required to provide additional multiplexing on top of the listen if there was a GID/SID conflict due to IP aliasing features of an O/S.

5.2.6 Completion Model

The transport provides a completion queue model that allows the SDP module to either poll for completions or to get an indication of a completion via a callback. It is assumed that the underlying transport and link interface provides the best possible de-serialization and scale-up with the completion callbacks based on the hardware characteristics and workloads running on the system. The SDP module will be capable of running at any priority including interrupt level.

SDP will initialize the CQ and size it according to the maximum work requests expected for sends, receives, and RDMA transfers based on the initial socket creation and socket option settings. The SDP module will setup one CQ for send messages and RDMA requests, and another one for receive messages. This architecture allows for development of independent SDP send and receive completion engines that can work independently.

5.2.7 Data Transfer Models

SDP defines four data transfer mechanisms:

- Bcopy – transfer of user data from send buffers into receive private buffers.

- Read Zcopy – transfer of user data through RDMA reads, preferably directly to and from user buffers.
- Write Zcopy – transfer of user data through RDMA writes, preferably directly to and from user buffers.
- Transaction – an optimized user data transfer for transactions that piggy-backs user data transfers using private buffers on the top of the Write Zcopy mechanism used to transfer the data on the opposite half-connection.

The SDP module will be compliant with all data transfer modes as specified by the InfiniBand annex A4 Socket Direct Protocol Specification.



NOTE

All models will be supported at the final release of this product, however each mode will be implemented in phases in the following order: Bcopy, Read Zcopy and Write Zcopy, and finally transaction. Since Bcopy is the only required model to be compliant with the specification all phases will be compliant with SDP protocol.

5.2.8 Locking and Threading Model

The SDP thread safe design is based on the premise that the SDP module will be driven by external events and thus will not require any dedicated threads of control, which are exclusive to the SDP module. External agents, such as an user application thread, an OPSI initialization thread, or system event notification will make calls into the SDP module thus providing threads of control. Access to internal SDP data structures will be serialized by the use of simple locks.

The offload transport may have dedicated threads of control that can invoke up-calls into SDP module. One such case would be the invocation of an IO completion callback routine. The SDP module is designed to limit the amount of processing in the IO completion callback thread if the transport is calling in a high-priority thread context such as interrupt level. All up-calls will be checked for priority levels and processed accordingly.

5.3 InfiniBand Offload Transport Module

The InfiniBand Offload Transport is designed to abstract reliable hardware transport interface specifics in such a fashion that an offload protocol (OP) module, Sockets Direct Protocol (SDP) in our case, can interface to a consistent transport API over a potentially wide range of differing reliable hardware transports. The InfiniBand Transport module will support any hardware interface that utilizes the InfiniBand Access Layer (IAL).

The InfiniBand OT has at least one Linux network link device associated with it, represented by the Linux `struct netdevice`. Once the InfiniBand OT registers with the SDP module, it will provide link device configuration events to the SDP module. In this manner, Internet link device configuration changes (IP address assignment, address change or device shutdown) will be propagated to the offload SDP agent.

InfiniBand OT service's defines the Internet Protocol Address format as its base level addressing mechanism. Offload protocol modules establish connections or send/receive datagrams based on IPv4 or IPv6 addressing formats. Initially, only the IPv4 address format will be supported. Offload transports are expected to convert an IP address into the transport/link specific address format required for transmission or reception. The InfiniBand OT module is designed to abstract the InfiniBand transport/link into common OT definitions.

The InfiniBand OT module provides four groups of basic services. Each service includes a collection of OT components that provide an implementation instance of that service. That is to say, each component has a well-defined offload transport service interface with a behavior implemented by an offload transport.

Figure 11 shows the OTI components that this specification defines.

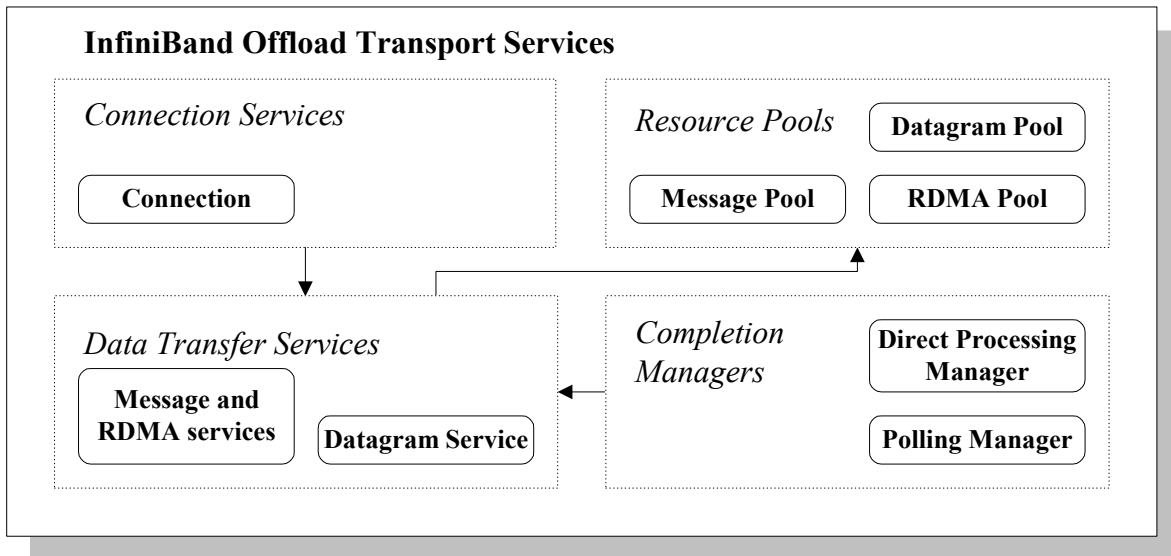


Figure 11: InfiniBand offload Transport Service Components

The transport architecture is flexible enough to allow the creation of additional components for each type of service within the transport. However, the interfaces between internal services and components not shown in Figure 11 are outside the scope of this document.

5.3.1 Linux Kernel Modifications

The InfiniBand OT is designed as a loadable kernel module. No Modifications to the Linux kernel are required.

5.3.2 Initialization

The InfiniBand OT kernel module will be demand loaded after the discovery of offload capable hardware, normally during system initialization. Actual InfiniBand OT initialization is performed in accordance with the standard Linux kernel module load procedure; the ‘init_module()’ function is called once the InfiniBand OT module has been successfully loaded into the kernel.

5.3.3 Shutdown

The InfiniBand OT shutdown can be initiated in one of two ways:

- A network link device, which the InfiniBand OT is monitoring, is shutdown via normal system shutdown procedures. System shutdown results in a link device event notification (system shutdown) being delivered to the InfiniBand OT.
- The InfiniBand OT module is forcibly removed using the Linux system administration command ‘rmmod’ (remove module).

With either method, the InfiniBand OT module unload function ‘cleanup_module()’ is called per the standard Linux kernel module unload procedure.

The InfiniBand OT ‘cleanup_module()’ function will notify the SDP module. The SDP module will notify the OPSI (Offload Protocol Switching Interface) of the interface shutdown, which results in the shutdown of INET addresses in use by the OPSI.

5.3.4 Connection Services

Connection services abstract and simplify the details of using a specific connection protocol. A connection service is responsible for creating data transfer service’s, authenticating their data transfers with remote endpoints, and presenting configured data transfer services to the user. The library provides a single channel connection service that is used to establish point-to-point communication over the transport link medium to a remote transport agent or endpoint.

Communication endpoints are identified by Internet Protocol addresses. Initially only IPv4 (Internet Protocol Version 4) format addresses are supported, with provisions to support IPv6 into the near future.

5.3.5 Buffer Strategy

Resource pools manage the allocation and distribution of data buffers, work requests, and transfer elements. Transfer elements are used to specify data transfer operations to the library. Work requests are required internally by the library to perform data transfers over the fabric interconnection hardware. There are three types of resource pools: datagram pools, message pools, and Remote Direct Memory Access (RDMA) pools.

5.3.5.1 Datagram Buffer Pool

The datagram pool administers datagram elements; queues pending requests for asynchronous operation; and, at the user's option, allocates and manages the distribution of data buffers used for connectionless message-passing operations. Datagram elements are used to request data transfers using connectionless communication and support Segmentation and Reassembly (SAR) of connectionless messages.

5.3.5.2 Message Buffer Pool

The message pool is similar to the datagram pool, but message elements are used to perform data transfers over connected channels. The message pool administers message elements; queues pending requests for asynchronous operation; and, at the user's option, allocates and manages the distribution of data buffers used for message-passing operations.

5.3.5.3 RDMA Resource Pool

The RDMA pool manages data transfer request elements and, optionally, RDMA data buffers. The data transfer request elements are used to signal the library to perform RDMA operations over connected channels. The RDMA pool provides asynchronous operation for requesting resources.

5.3.6 Data Transfer Services

Data transfer services are responsible for posting the following types of requests: sends, receives, RDMA reads, and RDMA writes. These services provide work queue management, perform message-passing flow control and manage scatter-gather data transfer requests. Connection services are used to create and configure data transfer services.

5.3.6.1 Message and RDMA Service

A channel provides a connected communication path between two transport endpoints. Channels support both RDMA and message-passing data transfers and may optionally provide message-level flow control.

5.3.6.2 Datagram Service

A datagram service provides connectionless communication to remote endpoints. Multiple endpoints may be reached through a single datagram service; however, datagram services support only connectionless message data transfers.

5.3.7 Completion Services

Completion services check for and process completed data transfer operations. Completion services are responsible for invoking any necessary post-processing, including user callbacks, and for resuming stalled data transfers.

5.3.7.1 Direct Processing Manager

The direct processing manager processes completions of data transfer requests within the context of the completion callback that delivered the completion notification to the library. This provides the lowest possible latency for processing a single completion by avoiding a thread context switch.

5.3.7.2 Polling Manager

A polling manager is a completion service without active threads. A polling manager manages completions across multiple data transfer services, but it relies on the user to poll for completed requests.

5.3.8 Locking and Threading Model

The InfiniBand OT thread safe design is based on the premise that it will be driven by external events and thus will not require any dedicated threads of control. External agents, such as the SDP module, the InfiniBand Access Layer, or system event notification will make calls into the InfiniBand OT thus providing threads of control through the InfiniBand OT. Access to internal InfiniBand OT data structures will be serialized by the use of simple locks.

6. System Resource Usage

6.1 Sockets Direct Protocol Module

Memory requirements are dependent on the socket allocation size. Private buffers will be allocated based on the socket buffer size as specified by either the default O/S setting or by the `setsockopt` call from the application. Additional overhead on a per socket basis includes the `struct offload_ops_t` (~128bytes).

6.2 InfiniBand Transport Module

The IB Offload Transport will consume additional memory on a per socket basis based on the connection attributes (maximum send and receive depths). Refer to the IB HCA hardware for CQ and QP overhead.

7. Internal Compatibility

7.1 Interaction with Other Components

7.1.1 Socket Driver and OPS Module Interaction

The Offload Protocol Switch (OPS) module provides a dynamic binding facility for offload protocols modules. It is a loadable module that registers dynamically with the Linux socket driver (socket.c). It exposes a new address family (AF_INET_OFFLOAD) but at the same time interacts with the AF_INET address family so that IPPROTO_ANY traffic can be directed to both the offload protocols under AF_INET_OFFLOAD and to the standard AF_INET protocols. Applications can also directly create sockets on AF_INET_OFFLOAD address family and bind to any offload protocols registered with this offload address family, without requiring any kernel patches.

In order to seamlessly support sockets created by applications with AF_INET address family a small patch must be applied to the socket driver (socket.c) sock_create() code to direct AF_INET address traffic to the offload address family module (AF_INET_OFFLOAD) exposed by the OPS module. The OPS module will switch sockets appropriately based on family, protocol type, and protocol. No modifications are needed in the AF_INET stack since the OPS module will interact with the standard AF_INET stack via the address family socket interface. Figure 3 shows the original Linux socket driver address family switching logic and the changes in the proposed kernel patch.

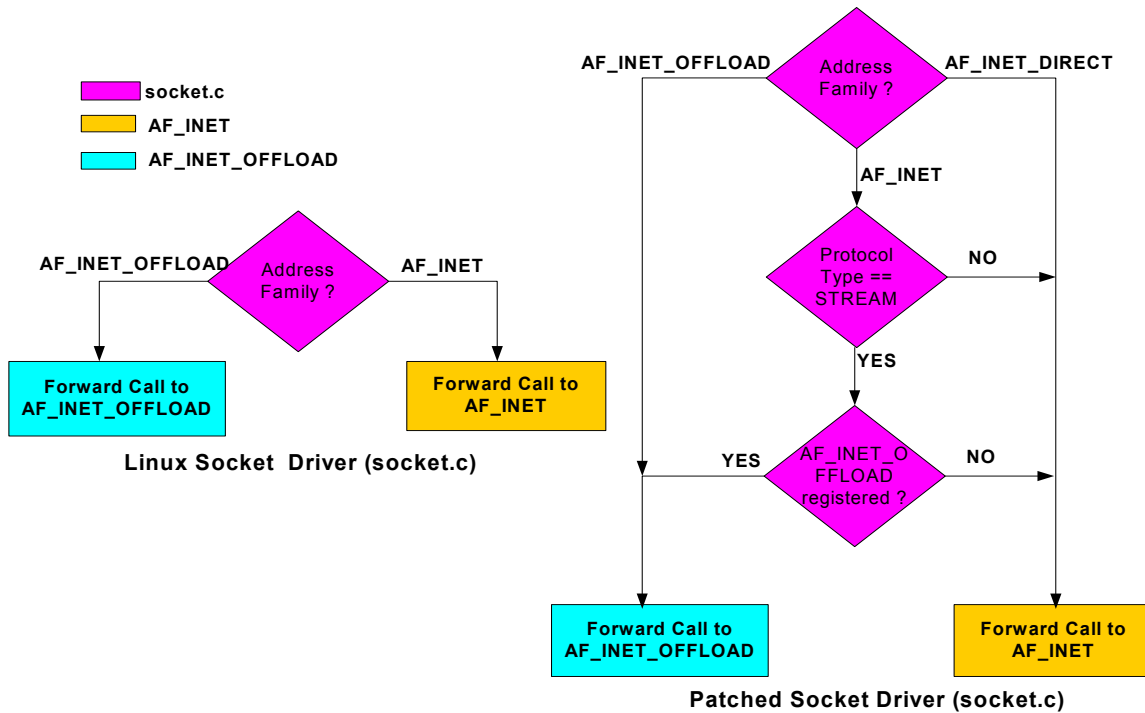


Figure 12: Proposed Patch to Linux Socket Driver

7.1.2 OPS and SDP Module Interaction

Figure 13 shows the various steps involved in the OPS and offload protocol modules initialization and operation

1. The socket.c sock_create() code is modified to forward all AF_INET sock_create calls to AF_INET_OFFLOAD module. The AF_INET_OFFLOAD module will direct the create call based on the offload protocols that have registered. If the create is forwarded back to the AF_INET, the OPS will use AF_INET_DIRECT so that the sock_create() code will know to switch the create directly to the AF_INET path thus sending all subsequent socket calls directly to the correct family.
2. The AF_INET_OFFLOAD module, as part of its initialization, registers its address family AF_INET_OFFLOAD to the socket layer by calling sock_register () call back to the socket layer. All future socket calls, with address family AF_INET_OFFLOAD will be directed by the socket layer to the OPS layer.
3. When Offload Protocol modules get loaded, the offload_register_protosw() is called to register their entry points and capabilities like self-routing, rdma etc to the AF_INET_OFFLOAD module.
4. When a network interface configuration changes (address assignment, address change, etc.), the Offload protocol module which is bound to this network interface notifies the OPS module by calling offload_notification().
5. Depending on the incoming request, the OPS module will then switch to proper protocol module. The switching policy depends on the protocol capabilities, binding priority set by the user. The AF_INET stack is the default stack, and if none of the offload protocol modules are loaded or if none of the module capabilities matches the incoming request, the OPS will forward the request to the AF_INET stack. For protocol modules that do not support self-routing, the AF_INET_OFFLOAD driver will handle the routing issues.

Once an appropriate protocol module is chosen, it is up to the protocol stack to handle the request. For example if SDP protocol module is chosen to service an request, then it is up to the SDP module to establish a session with its peer, pre register buffer element and decide on mode of data transfer (Send vs. RDMA Write for example).

**IBA Software Architecture
Offload Sockets Framework and Sockets Direct Protocol
High Level Design**

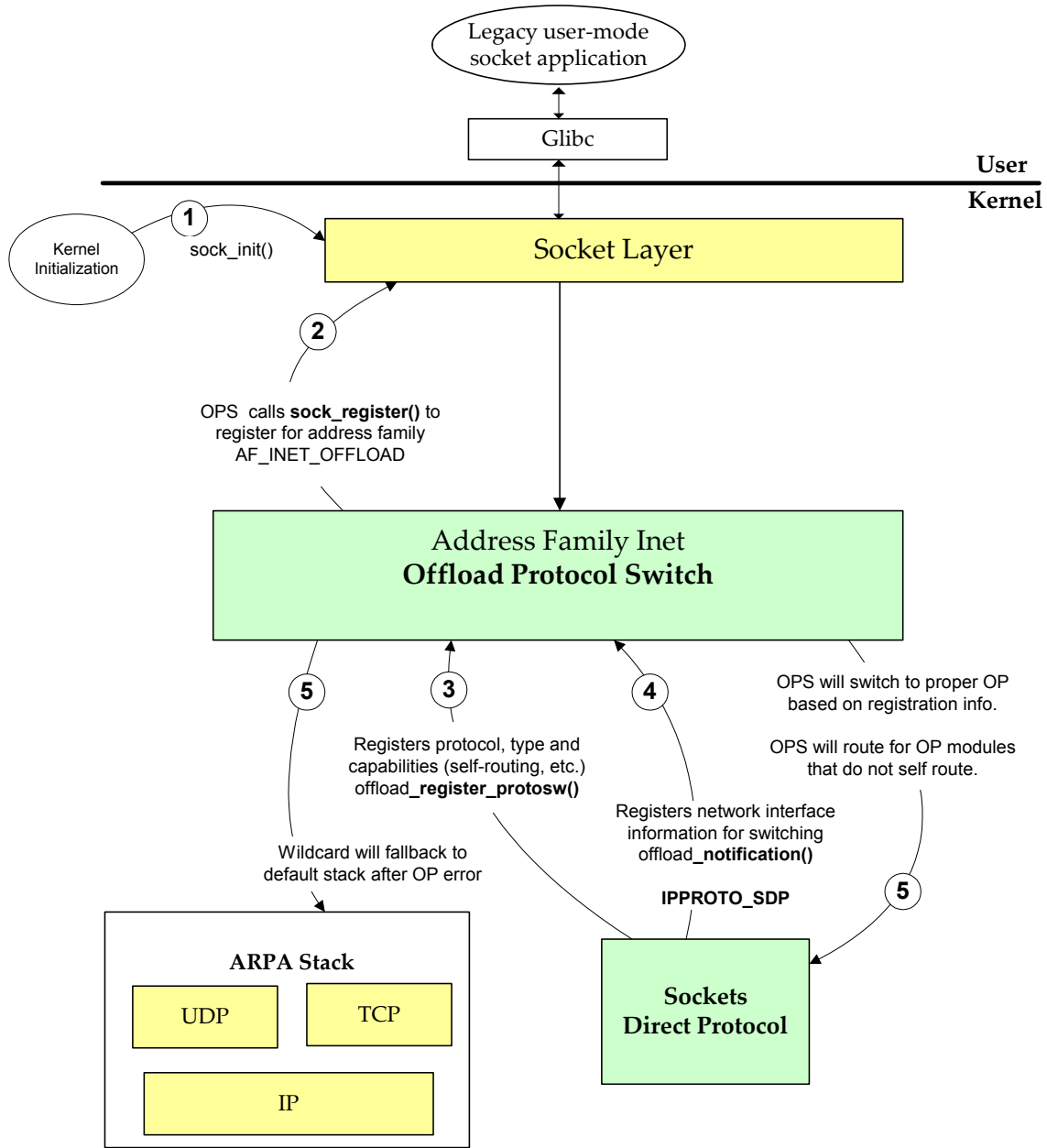


Figure 13: OPS and SDP interaction

7.1.3 SDP and InfiniBand Transport Interaction

Offload Transports (OTs) implement the offload transport interface abstractions. The InfiniBand transport is designed to be a kernel loadable module, which provides an InfiniBand hardware transport specific implementation of the offload transport abstraction defined by the offload transport interface.

Figure 14 shows the various steps involved in the OT module initialization and operation.

1. The InfiniBand (IB) Transport module will register with the SDP module and provide the transport operations interface and capabilities.
2. The InfiniBand Transport module obtains kernel notification services for network devices, IP configuration, and IP routing information using the device name provided by the link driver. When an IP configuration notification occurs for this net device the transport module will forward the via the event notification upcall to SDP, if SDP has registered for events.
3. The SDP module, using the transport register event call, will register for address and net device status changes.
4. The SDP module maintains a list of transports, with address information, and will switch to appropriate transport based on address information during a socket bind.

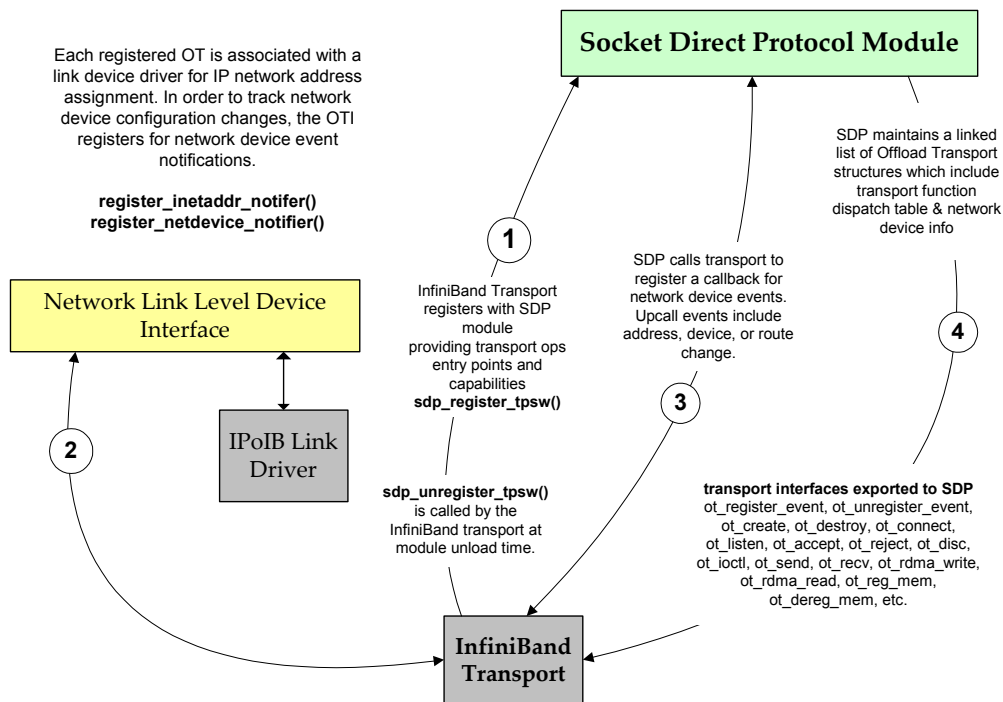


Figure 14 SDP and InfiniBand Transport

8. External Compatibility

8.1 Standards

8.1.1 Sockets Direct Protocol Module

The SDP module and state machine is designed to support all compliance statements highlighted in the InfiniBand Annex A4 Socket Direct Protocol Specification.

8.1.2 InfiniBand Offload Transport Module

The InfiniBand Transport module is designed to support the new Offload Transport Interface (OTI) at the top and interface with the standard InfiniBand Access Layer at the bottom.

9. Other Dependencies

9.1 Offload Protocol Switch Module

9.2 Sockets Direct Protocol Module

SDP has the following dependencies on an IB specific Offload Transport:

- Offload Transport must provide a listen interface that supports all 76 bytes of IB CM REQ private data and will multiplex based on the IP address/port pair provided in the endpoint handle and the specified private data. The transport shall not byte swap any of the private data. This private data will include, in big-endian format, the SDP Hello and the HelloAck messages.
- Offload Transport must support a "consumer reject" error code during connection callbacks and reject calls. IB transports are expected to map this consumer reject code to IB code 28.
- Connect shall be non-blocking.
- Calls that may block include `accept`, `create_endpoint`, `res_pool_get`,
- Accept and Reject calls MUST not be restricted to the connection callback thread.
- IB transports are expected to map the IP port to a SID as detailed in the IBTA "Annex A0 Applications Specific Identifiers" specifications.
- The transport interface cannot abstract the hardware's R_KEY when registering memory for RDMA transfers. The actual 32-bit R_KEY from the IB interface must be provided to the SDP module. This key will be provided to the remote SDP entity via the SDP SinkAvail and the SrcAvail messages.

9.3 InfiniBand Transport Module

No other dependencies at this time.

10. Installation and Configuration

10.1 Offload Protocol Switch Module

TBD: update after coding.

10.1.1 Installation

+

10.1.2 Configuration

10.2 Sockets Direct Protocol Module

TBD: update after coding.

10.2.1 Installing

+

10.2.2 Configuring

10.3 InfiniBand Transport Module

TBD: update after coding.

10.3.1 Installing

+

10.3.2 Configuring

11. Unresolved Issues

11.1 Offload Protocol Switch Module

No unresolved issues at this time.

11.2 Sockets Direct Protocol Module

No unresolved issues at this time.

11.3 InfiniBand Transport Module

No unresolved issues at this time.

12. Data Structures and APIs

12.1 Offload Protocol Switch Definitions

To view the data structures and APIs associated with this component, click on [here](#).

12.2 SDP Definitions

To view the data structures and APIs associated with this component, click on [here](#)

12.3 InfiniBand OT Definitions

To view the data structures and APIs associated with this component, click on [here](#)