



## **TCPDirect User Guide**

SF-116303-CD , Issue 4

2017/06/28 12:04:32

Solarflare Communications Inc



## TCPDirect User Guide

Copyright © 2017 SOLARFLARE Communications, Inc. All rights reserved.

The software and hardware as applicable (the "Product") described in this document, and this document, are protected by copyright laws, patents and other intellectual property laws and international treaties. The Product described in this document is provided pursuant to a license agreement, evaluation agreement and/or non-disclosure agreement. The Product may be used only in accordance with the terms of such agreement. The software as applicable may be copied only in accordance with the terms of such agreement.

Onload is licensed under the GNU General Public License (Version 2, June 1991). See the LICENSE file in the distribution for details. The Onload Extensions Stub Library is Copyright licensed under the BSD 2-Clause License.

Onload contains algorithms and uses hardware interface techniques which are subject to Solarflare Communications Inc patent applications. Parties interested in licensing Solarflare's IP are encouraged to contact Solarflare's Intellectual Property Licensing Group at:

Director of Intellectual Property Licensing  
Intellectual Property Licensing Group  
Solarflare Communications Inc, 7505 Irvine Center Drive  
Suite 100  
Irvine, California 92618

You will not disclose to a third party the results of any performance tests carried out using Onload or EnterpriseOnload without the prior written consent of Solarflare.

The furnishing of this document to you does not give you any rights or licenses, express or implied, by estoppel or otherwise, with respect to any such Product, or any copyrights, patents or other intellectual property rights covering such Product, and this document does not contain or represent any commitment of any kind on the part of SOLARFLARE Communications, Inc. or its affiliates.

The only warranties granted by SOLARFLARE Communications, Inc. or its affiliates in connection with the Product described in this document are those expressly set forth in the license agreement, evaluation agreement and/or non-disclosure agreement pursuant to which the Product is provided. EXCEPT AS EXPRESSLY SET FORTH IN SUCH AGREEMENT, NEITHER SOLARFLARE COMMUNICATIONS, INC. NOR ITS AFFILIATES MAKE ANY REPRESENTATIONS OR WARRANTIES OF ANY KIND (EXPRESS OR IMPLIED) REGARDING THE PRODUCT OR THIS DOCUMENTATION AND HEREBY DISCLAIM ALL IMPLIED WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NON-INFRINGEMENT, AND ANY WARRANTIES THAT MAY ARISE FROM COURSE OF DEALING, COURSE OF PERFORMANCE OR USAGE OF TRADE. Unless otherwise expressly set forth in such agreement, to the extent allowed by applicable law (a) in no event shall SOLARFLARE Communications, Inc. or its affiliates have any liability under any legal theory for any loss of revenues or profits, loss of use or data, or business interruptions, or for any indirect, special, incidental or consequential damages, even if advised of the possibility of such damages; and (b) the total liability of SOLARFLARE Communications, Inc. or its affiliates arising from or relating to such agreement or the use of this document shall not exceed the amount received by SOLARFLARE Communications, Inc. or its affiliates for that copy of the Product or this document which is the subject of such liability.

The Product is not intended for use in medical, life saving, life sustaining, critical control or safety systems, or in nuclear facility applications.

A list of patents associated with this product is at <http://www.solarflare.com/patent>

SF-116303-CD

Last Revised: 62017

Issue 4



# Contents

<b>1</b>	<b>TCPDirect</b>	<b>1</b>
1.1	Introduction . . . . .	1
<b>2</b>	<b>Overview</b>	<b>3</b>
2.1	Platforms . . . . .	3
2.2	Components . . . . .	4
2.3	Capabilities and Restrictions . . . . .	4
2.3.1	Protocols . . . . .	4
2.3.2	OS . . . . .	5
2.4	How TCPDirect Increases Performance . . . . .	5
2.4.1	Overhead . . . . .	5
2.4.2	Latency . . . . .	6
2.4.3	Bandwidth . . . . .	6
2.4.4	Scalability . . . . .	6
2.5	Requirements . . . . .	6
2.5.1	Adapter . . . . .	6
2.5.2	License . . . . .	7
2.5.3	Onload . . . . .	7
2.5.4	Huge Pages . . . . .	7
2.5.5	PIO . . . . .	7

<b>3</b>	<b>Concepts</b>	<b>9</b>
3.1	Stacks . . . . .	9
3.2	Zockets . . . . .	9
3.2.1	TCP zockets . . . . .	9
3.2.2	UDP zockets . . . . .	10
3.2.3	Waitables . . . . .	10
3.3	Multiplexers . . . . .	10
3.4	TX alternatives . . . . .	10
<b>4</b>	<b>Example Applications</b>	<b>11</b>
4.1	zfudppingpong . . . . .	11
4.1.1	Usage . . . . .	11
4.2	zftcpingpong . . . . .	12
4.2.1	Usage . . . . .	12
4.3	zfaltpingpong . . . . .	12
4.4	zfsink . . . . .	12
4.4.1	Usage . . . . .	12
4.5	zftcpmtpong . . . . .	13
4.5.1	Usage . . . . .	13
4.6	Building the Example Applications . . . . .	13

<b>5</b>	<b>Using TCPDirect</b>	<b>15</b>
5.1	Components . . . . .	15
5.2	Compiling and Linking . . . . .	15
5.2.1	Header files . . . . .	15
5.2.2	Linking . . . . .	15
5.2.3	Debugging . . . . .	16
5.3	General . . . . .	16
5.4	Using stacks . . . . .	17
5.5	Using zockets . . . . .	17
5.6	UDP receive . . . . .	17
5.7	UDP send . . . . .	18
5.8	TCP listening . . . . .	18
5.9	TCP send and receive . . . . .	19
5.10	Alternative Tx queues . . . . .	20
5.11	Epoll – muxer.h . . . . .	21
5.12	Stack polling . . . . .	22
5.13	Miscellaneous . . . . .	23
5.14	zf_stackdump . . . . .	23
5.14.1	Usage . . . . .	23
5.14.2	stackdump output: stack . . . . .	23
5.14.3	stackdump output: UDP RX . . . . .	24
5.14.4	stackdump output: UDP TX . . . . .	24
5.14.5	stackdump output: TCP TX/RX . . . . .	24
<b>6</b>	<b>Worked Examples</b>	<b>27</b>
6.1	UDP ping pong example . . . . .	27
6.2	TCP pong example . . . . .	28

<b>7</b>	<b>Attributes</b>	<b>31</b>
7.1	alt_buf_size Attribute Reference . . . . .	32
7.2	alt_count Attribute Reference . . . . .	32
7.3	arp_reply_timeout Attribute Reference . . . . .	33
7.4	interface Attribute Reference . . . . .	33
7.5	log_format Attribute Reference . . . . .	34
7.6	log_level Attribute Reference . . . . .	34
7.7	max_tcp_endpoints Attribute Reference . . . . .	35
7.8	max_tcp_listen_endpoints Attribute Reference . . . . .	36
7.9	max_tcp_syn_backlog Attribute Reference . . . . .	36
7.10	max_udp_rx_endpoints Attribute Reference . . . . .	36
7.11	max_udp_tx_endpoints Attribute Reference . . . . .	37
7.12	n_bufs Attribute Reference . . . . .	37
7.13	name Attribute Reference . . . . .	38
7.14	reactor_spin_count Attribute Reference . . . . .	39
7.15	rx_ring_max Attribute Reference . . . . .	39
7.16	rx_ring_refill_batch_size Attribute Reference . . . . .	40
7.17	rx_ring_refill_interval Attribute Reference . . . . .	40
7.18	tcp_alt_ack_rewind Attribute Reference . . . . .	41
7.19	tcp_delayed_ack Attribute Reference . . . . .	41
7.20	tcp_finwait_ms Attribute Reference . . . . .	42
7.21	tcp_initial_cwnd Attribute Reference . . . . .	42
7.22	tcp_retries Attribute Reference . . . . .	43
7.23	tcp_syn_retries Attribute Reference . . . . .	43
7.24	tcp_synack_retries Attribute Reference . . . . .	44
7.25	tcp_timewait_ms Attribute Reference . . . . .	44
7.26	tcp_wait_for_time_wait Attribute Reference . . . . .	45
7.27	tx_ring_max Attribute Reference . . . . .	45



<b>8</b>	<b>Data Structure Index</b>	<b>47</b>
8.1	Data Structures . . . . .	47
<b>9</b>	<b>File Index</b>	<b>49</b>
9.1	File List . . . . .	49
<b>10</b>	<b>Data Structure Documentation</b>	<b>51</b>
10.1	zf_attr Struct Reference . . . . .	51
10.1.1	Detailed Description . . . . .	51
10.2	zf_muxer_set Struct Reference . . . . .	52
10.2.1	Detailed Description . . . . .	52
10.3	zf_stack Struct Reference . . . . .	52
10.3.1	Detailed Description . . . . .	52
10.4	zf_waitable Struct Reference . . . . .	52
10.4.1	Detailed Description . . . . .	53
10.5	zft Struct Reference . . . . .	53
10.5.1	Detailed Description . . . . .	53
10.6	zft_handle Struct Reference . . . . .	53
10.6.1	Detailed Description . . . . .	53
10.7	zft_msg Struct Reference . . . . .	54
10.7.1	Detailed Description . . . . .	54
10.7.2	Field Documentation . . . . .	54
10.7.2.1	flags . . . . .	54
10.7.2.2	iov . . . . .	54
10.7.2.3	iovcnt . . . . .	54
10.7.2.4	pkts_left . . . . .	55
10.7.2.5	reserved . . . . .	55
10.8	zftl Struct Reference . . . . .	55
10.8.1	Detailed Description . . . . .	55
10.9	zfur Struct Reference . . . . .	55
10.9.1	Detailed Description . . . . .	55
10.10	zfur_msg Struct Reference . . . . .	56
10.10.1	Detailed Description . . . . .	56
10.10.2	Field Documentation . . . . .	56
10.10.2.1	dgrams_left . . . . .	56
10.10.2.2	flags . . . . .	56
10.10.2.3	iov . . . . .	56
10.10.2.4	iovcnt . . . . .	57
10.10.2.5	reserved . . . . .	57
10.11	zfut Struct Reference . . . . .	57
10.11.1	Detailed Description . . . . .	57

<b>11 File Documentation</b>	<b>59</b>
11.1 attr.h File Reference	59
11.1.1 Detailed Description	60
11.1.2 Function Documentation	60
11.1.2.1 zf_attr_alloc	60
11.1.2.2 zf_attr_doc	61
11.1.2.3 zf_attr_dup	61
11.1.2.4 zf_attr_free	61
11.1.2.5 zf_attr_get_int	62
11.1.2.6 zf_attr_get_str	62
11.1.2.7 zf_attr_reset	62
11.1.2.8 zf_attr_set_from_fmt	62
11.1.2.9 zf_attr_set_from_str	63
11.1.2.10 zf_attr_set_int	63
11.1.2.11 zf_attr_set_str	63
11.2 muxer.h File Reference	64
11.2.1 Detailed Description	64
11.2.2 Function Documentation	65
11.2.2.1 zf_muxer_add	65
11.2.2.2 zf_muxer_alloc	65
11.2.2.3 zf_muxer_del	65
11.2.2.4 zf_muxer_free	66
11.2.2.5 zf_muxer_mod	66
11.2.2.6 zf_muxer_wait	66
11.2.2.7 zf_waitable_event	67
11.2.2.8 zf_waitable_fd_get	67
11.2.2.9 zf_waitable_fd_prime	68
11.3 types.h File Reference	68
11.3.1 Detailed Description	68

11.4	x86.h File Reference	69
11.4.1	Detailed Description	69
11.5	zf.h File Reference	69
11.5.1	Detailed Description	69
11.6	zf_alts.h File Reference	69
11.6.1	Detailed Description	70
11.6.2	Function Documentation	70
11.6.2.1	zf_alternatives_alloc	70
11.6.2.2	zf_alternatives_cancel	70
11.6.2.3	zf_alternatives_free_space	70
11.6.2.4	zf_alternatives_query_overhead_tcp	71
11.6.2.5	zf_alternatives_release	71
11.6.2.6	zf_alternatives_send	72
11.6.2.7	zft_alternatives_queue	72
11.7	zf_platform.h File Reference	73
11.7.1	Detailed Description	73
11.8	zf_reactor.h File Reference	73
11.8.1	Detailed Description	74
11.8.2	Function Documentation	74
11.8.2.1	zf_reactor_perform	74
11.8.2.2	zf_stack_has_pending_work	74
11.9	zf_stack.h File Reference	75
11.9.1	Detailed Description	76
11.9.2	Macro Definition Documentation	76
11.9.2.1	EPOLLSTACKHUP	76
11.9.3	Function Documentation	76
11.9.3.1	zf_deinit	76
11.9.3.2	zf_init	76
11.9.3.3	zf_stack_alloc	76

11.9.3.4	zf_stack_free	77
11.9.3.5	zf_stack_is_quiescent	77
11.9.3.6	zf_stack_to_waitable	78
11.10	zf_tcp.h File Reference	78
11.10.1	Detailed Description	79
11.10.2	Function Documentation	79
11.10.2.1	zft_addr_bind	79
11.10.2.2	zft_alloc	80
11.10.2.3	zft_connect	80
11.10.2.4	zft_error	81
11.10.2.5	zft_free	82
11.10.2.6	zft_get_header_size	82
11.10.2.7	zft_get_mss	82
11.10.2.8	zft_getname	82
11.10.2.9	zft_handle_free	83
11.10.2.10	zft_handle_getname	83
11.10.2.11	zft_recv	83
11.10.2.12	zft_send	84
11.10.2.13	zft_send_single	85
11.10.2.14	zft_send_space	85
11.10.2.15	zft_shutdown_tx	86
11.10.2.16	zft_state	86
11.10.2.17	zft_to_waitable	86
11.10.2.18	zft_zc_recv	87
11.10.2.19	zft_zc_recv_done	87
11.10.2.20	zft_zc_recv_done_some	87
11.10.2.21	zftl_accept	88
11.10.2.22	zftl_free	88
11.10.2.23	zftl_getname	88

11.10.2.24	zftl_listen . . . . .	89
11.10.2.25	zftl_to_waitable . . . . .	89
11.11	zf_udp.h File Reference . . . . .	90
11.11.1	Detailed Description . . . . .	91
11.11.2	Function Documentation . . . . .	91
11.11.2.1	zfur_addr_bind . . . . .	91
11.11.2.2	zfur_addr_unbind . . . . .	92
11.11.2.3	zfur_alloc . . . . .	93
11.11.2.4	zfur_free . . . . .	93
11.11.2.5	zfur_pkt_get_header . . . . .	93
11.11.2.6	zfur_to_waitable . . . . .	94
11.11.2.7	zfur_zc_rcv . . . . .	94
11.11.2.8	zfur_zc_rcv_done . . . . .	94
11.11.2.9	zfut_alloc . . . . .	95
11.11.2.10	zfut_free . . . . .	95
11.11.2.11	zfut_get_header_size . . . . .	95
11.11.2.12	zfut_get_mss . . . . .	96
11.11.2.13	zfut_send . . . . .	96
11.11.2.14	zfut_send_single . . . . .	96
11.11.2.15	zfut_to_waitable . . . . .	97
	<b>Index</b>	<b>99</b>



## Chapter 1

# TCPDirect

Solarflare's TCPDirect is highly accelerated network middleware. It uses similar techniques to Onload, but delivers lower latency. In order to achieve this, TCPDirect supports a reduced feature set and uses a proprietary API.

### 1.1 Introduction

The TCPDirect API provides an interface to an implementation of TCP and UDP over IP. This is dynamically linked into the address space of user-mode applications, and granted direct (but safe) access to the network-adaptor hardware. The result is that data can be transmitted to and received from the network directly by the application, without involvement of the operating system. This technique is known as 'kernel bypass'.

Kernel bypass avoids disruptive events such as system calls, context switches and interrupts and so increases the efficiency with which a processor can execute application code. This also directly reduces the host processing overhead, typically by a factor of two, leaving more CPU time available for application processing. This effect is most pronounced for applications which are network intensive.

The key features of TCPDirect are:

- **User-space:** TCPDirect can be used by unprivileged user-space applications.
- **Kernel bypass:** Data path operations do not require system calls.
- **Low CPU overhead:** Data path operations consume very few CPU cycles.
- **Low latency:** Suitable for low latency applications.
- **High packet rates:** Supports millions of packets per second per core.
- **Zero-copy:** Particularly efficient for filtering and forwarding applications.
- **Flexibility:** Supports many use cases.





## Chapter 2

# Overview

This part of the documentation gives an overview of TCPDirect and how it is often used.

### 2.1 Platforms

TCPDirect can be run on 8000-series Solarflare adapters with a suitable license (e.g. the 'Plus' license).

TCPDirect can also be run on 7000-series adapters which require both the Onload license and a TCPDirect license.

Refer to the *Solarflare Server Adapter User Guide* 'Product Specifications' for adapter details; for licensing queries, please contact your sales representative.

TCPDirect is supported on:

- Red Hat Enterprise Linux 6.6 - 6.8
- Red Hat Enterprise Linux 7.0 - 7.2
- SuSE Linux Enterprise Server 11 sp2, sp3, sp4
- SuSE Linux Enterprise Server 12 base and sp1
- Canonical Ubuntu Server LTS 14.04, 16.04
- Canonical Ubuntu Server 16.10
- Debian 7 "Wheezy"
- Debian 8 "Jessie"
- Linux kernels 2.6.32 - 4.8.

## 2.2 Components

TCPDirect is supplied as:

- header files containing the proprietary public API
- a binary library for linking into your application.

To use TCPDirect, you must have access to the source code for your application, and the toolchain required to build it. You must then replace the existing calls for network access with appropriate calls from the TCPDirect API. Typically this involves replacing calls to the BSD sockets API. Finally you must recompile your application, linking in the TCPDirect library.

For more details, see [Using TCPDirect](#).

If you do not have access to source code for your application, you can instead accelerate it with Onload.

## 2.3 Capabilities and Restrictions

TCPDirect supports a carefully selected feature set that allows it to run many real-world applications, without losing performance to resource-intensive features that are seldom used.

Before porting an application to TCPDirect, you should ensure that it supports the features that you require. The subsections below list the support for different features.

If your application requires features that are unsupported by TCPDirect, consider instead using Onload or ef\_vi:

- Onload has higher latency than TCPDirect, but a full feature set.

Onload supports all of the standard BSD sockets API, meaning that no modifications are required to POSIX-compliant socket-based applications being accelerated. Like TCPDirect, Onload uses kernel bypass for applications over TCP/IP and UDP/IP protocols.

- Ef\_vi has even lower latency than TCPDirect, but operates at a lower level.

Ef\_vi is a low level OSI level 2 interface which sends and receives raw Ethernet frames, and exposes many of the advanced capabilities of Solarflare network adapters. But because the ef\_vi API operates at this low level, any application using it must implement the higher layer protocols itself, and also deal with any exceptions or other unusual conditions.

### 2.3.1 Protocols

The table below shows the protocols that are supported by TCPDirect and (for comparison) by Onload:

Protocol	TCPDirect	Onload
IPv4	Yes	Yes
IPv6	No	No

UDP	Yes	Yes
TCP	Yes	Yes
TCP header options (e.g. timestamps)	No	Yes
VLANs	Yes	Yes
Multicast RX	Yes	Yes
Multicast TX	Yes	Yes
Multicast loopback	No	Yes

## 2.3.2 OS

The table below shows the OS features that are supported by TCPDirect and (for comparison) by Onload:

OS	TCPDirect	Onload
Preload	No	Yes
Static link	Yes	Yes
Dynamic link	Yes	Yes
Direct API	Yes	Yes
Bonding	No	Yes
Teaming	No	Yes
Send/receive via non-SFC interface	No	Yes
Multiple threads	Yes	Yes
Multiple processes	Yes	Yes
Sharing stacks between threads and processes	No	Yes
Multiple stacks	Yes	Yes
fork()	Yes, with limitations (no shared stacks or zockets)	Yes
dup()	N/A (no file descriptors)	Yes
User-level only	Yes	No
Interrupts	No	Yes
Huge pages	Yes	Yes

As mentioned in the table above, TCPDirect stacks and zockets cannot be shared between processes. The one partial exception is that a stack may be used in a `fork()` child if it is not used in the parent after calling `fork()`.

## 2.4 How TCPDirect Increases Performance

TCPDirect can significantly reduce the costs associated with networking by reducing CPU overheads and improving performance for latency, bandwidth and application scalability.

### 2.4.1 Overhead

Transitioning into and out of the kernel from a user-space application is a relatively expensive operation: the equivalent of hundreds or thousands of instructions. With conventional networking such a transition is required every time the application sends and receives data. With TCPDirect, the TCP/IP processing can be done entirely within the user-process, eliminating expensive application/kernel transitions, i.e. system calls. In addition, the TCPDirect TCP/IP stack is highly tuned, offering further overhead savings.

The overhead savings of TCPDirect mean more of the CPU's computing power is available to the application to do useful work.

## 2.4.2 Latency

Conventionally, when a server application is ready to process a transaction it calls into the OS kernel to perform a 'receive' operation, where the kernel puts the calling thread 'to sleep' until a request arrives from the network. When such a request arrives, the network hardware 'interrupts' the kernel, which receives the request and 'wakes' the application.

All of this overhead takes CPU cycles as well as increasing cache and translation lookaside-buffer (TLB) footprint. With TCPDirect, the application can remain at user level waiting for requests to arrive at the network adapter and process them directly. The elimination of a kernel-to-user transition, an interrupt, and a subsequent user-to-kernel transition can significantly reduce latency. In short, reduced overheads mean reduced latency.

## 2.4.3 Bandwidth

Because TCPDirect imposes less overhead, it can process more bytes of network traffic every second. Along with specially tuned buffering and algorithms designed for high speed networks, TCPDirect allows applications to achieve significantly improved bandwidth.

## 2.4.4 Scalability

Modern multi-core systems are capable of running many applications simultaneously. However, the advantages can be quickly lost when the multiple cores contend on a single resource, such as locks in a kernel network stack or device driver. These problems are compounded on modern systems with multiple caches across many CPU cores and Non-Uniform Memory Architectures.

TCPDirect results in the network adapter being partitioned and each partition being accessed by an independent copy of the TCP/IP stack. The result is that with TCPDirect, doubling the cores really can result in doubled throughput.

## 2.5 Requirements

### 2.5.1 Adapter

The following list identifies the minimum driver and firmware requirements for adapters running TCPDirect applications:

- Net driver minimum: 4.10.1011 - this is available in the openonload-201606-u1 and Enterprise Onload 5.0 distributions.
- Firmware minimum: 6.2.3.1000
- Firmware variant: The adapter must be configured to use the ultra-low-latency firmware variant. The firmware variant can be identified/set using the Solarflare sfboot utility from the Solarflare Linux Utilities package (SF-107601-LS).

## 2.5.2 License

SFN8000 series adapters - require an Onload license and TCPDirect license. The 'Plus' license will include both required licenses.

SFN7000 series adapters - require an Onload license and TCPDirect license.

Installed licenses can be identified using the Solarflare sfkey utility from the Solarflare Linux Utilities package (SF-107601-LS).

```
# sfkey
enp4s0f0,enp4s0f1: 712200205071234567890123 (Flareon)
  Product name      Solarflare SFN7122F SFP+ Server Adapter
  Serial number     712200205071234567890123
  Installed keys    Onload, TCP Direct

enp5s0f0,enp5s0f1: 000F5341C700 (8xxx)
  Product name      Solarflare Flareon Ultra 8000 Series 10G Adapter
  Serial number     852200203001234567890123
  Installed keys    Plus
```

## 2.5.3 Onload

Openonload from version 201606-u1 or EnterpriseOnload 5.0 must be installed on the server.

## 2.5.4 Huge Pages

TCPDirect requires the allocation of huge pages. Huge pages are needed for each stack created - approximate minimum 10 huge pages per stack, the number of zockets created and number of packet buffers required. Some experimentation is needed to identify suitable page allocation needs for an application, but a general recommendation would be to allocate at least 40 huge pages per stack and then to use `zf_stackdump` to identify packet buffer usage.

Further information including allocation commands for huge pages is available in the Onload User Guide (SF-104474-CD).

## 2.5.5 PIO

TCPDirect uses PIO packet buffers and these are available by default from the adapter driver. Users should be aware that PIO buffers are a limited resource used by the driver for non-accelerated sockets, by Onload stacks which require 1 PIO buffer per VI created and by the TCPDirect application which require 1 PIO buffer per stack.

To ensure there are sufficient PIO buffers available, it may be necessary to restrict or prevent the driver and Onload non-critical sockets from using PIO.

Further information about PIO including configuration commands can be found in the Onload User Guide (SF-104474-CD).



## Chapter 3

# Concepts

This part of the documentation describes the concepts involved in TCPDirect.

### 3.1 Stacks

TCPDirect can have multiple network stacks. Each stack accesses a separate partition of the network adaptor, which improves security. Access to a given stack is not thread-safe, so typically each thread uses a separate TCP-Direct stack. This model avoids problems of lock-contention and cache-line bouncing and allows for scalability in a natural way.

### 3.2 Zockets

TCPDirect endpoints are represented by *zockets*. Zockets are similar to BSD or POSIX sockets, but are categorized by protocol and state into different zocket types.

Each type of zocket is represented by its own distinct data structure, and is handled by its own API calls. Because the type of zocket is known, code can be more efficient.

Zockets cannot be converted from one type to another.

#### 3.2.1 TCP zockets

For TCP, a new zocket can be listening or non-listening:

- a listening zocket is represented by a *TCP listening zocket* data structure
- a non-listening zocket is represented by a *TCP zocket handle*.

If the zocket is later connected, a *TCP zocket* data structure is created and returned. This occurs if:

- a connection to a TCP listening zocket is accepted
- a TCP zocket handle is connected.

### 3.2.2 UDP zockets

For UDP, a zocket can be used to receive or transmit:

- a receive zocket is represented by a *UDP receive zocket* data structure
- a transmit zocket is represented by a *UDP transmit zocket* data structure

### 3.2.3 Waitables

*Waitables* are handles that represent zockets and are used with the [multiplexer](#) interface.

Each type of zocket has an API call to return a waitable representing the zocket.

## 3.3 Multiplexers

A *multiplexer* (or *muxer*) allows multiple zockets to be polled for activity through a single call. The interface and behaviour are similar to the standard Linux *epoll* mechanism.

Each multiplexer is associated with a [stack](#) and can only be used to poll zockets from that stack. Zockets (represented by [waitables](#)) can be added to a multiplexer together with the set of events that the application is interested in. A zocket can only be a member of one multiplexer at a time.

When a multiplexer is polled the corresponding stack is polled for network events, and a list of zockets that are *ready* (readable, writable etc.) is returned.

## 3.4 TX alternatives

*TX alternatives* provide multiple alternative queues for transmission, that can be used to minimize latency. Different possible responses can be pushed through the TX path on the NIC, and held in different queues ready to transmit. When it is decided which response to transmit, the appropriate alternative queue is selected, and the queued packets are sent. Because the packets are already prepared, and are held close to the wire, latency is greatly reduced.

Due to differences in hardware architecture, the TX alternatives feature is not available on SFN7000 series adapters.



## Chapter 4

# Example Applications

Solarflare TCPDirect comes with a range of example applications - including source code and make files. This is a quick guide to using them, both for testing TCPDirect's effectiveness in an environment, and as starting points for developing applications.

Application	Description
<a href="#">zfudppingpong</a>	Measure round-trip latency with UDP.
<a href="#">zftcppingpong</a>	Measure round-trip latency with TCP.
<a href="#">zfaltppingpong</a>	Measure round-trip latency with TCP TX Alternatives.
<a href="#">zfsink</a>	Receive stream of UDP datagrams and demonstrate muxer.
<a href="#">zftcpmtpong</a>	Use of TCPDirect in multi-threaded applications.

### 4.1 zfudppingpong

The zfudppingpong application passes messages back and forth between two hosts using UDP, and uses this to measure the average round-trip latency.

#### 4.1.1 Usage

##### Server:

```
export ZF_ATTR=interface=ethX
zfudppingpong pong serverhost:serverport clienthost:clientport
```

##### Client:

```
export ZF_ATTR=interface=ethX
zfudppingpong ping clienthost:clientport serverhost:serverport
```

where:

- *ethX* is the name of the network interface to use,
- *serverhost* and *clienthost* identify the server and client machines (e.g. *hostname* or *192.168.0.10*), and
- *serverport* and *clientport* are port numbers of your choosing on the server and client machines

There are various additional options. See the help text for details.

## 4.2 zftcpingpong

The `zftcpingpong` application passes messages back and forth between two hosts using TCP, and uses this to measure the average round-trip latency. It illustrates actively and passively opened TCP connections, and has an option to use a *muxer*.

### 4.2.1 Usage

#### Server:

```
export ZF_ATTR=interface=ethX
zftcpingpong pong serverhost:serverport
```

#### Client:

```
export ZF_ATTR=interface=ethX
zftcpingpong ping serverhost:serverport
```

## 4.3 zfaltpingpong

The `zfaltpingpong` application illustrates use of the [TX alternatives](#) feature, which supports lower latency sends with TCP.

Usage is as for [zftcpingpong](#).

## 4.4 zfsink

The `zfsink` application demonstrates how to receive UDP datagrams, how to use the muxer, and the "waitable fd" mechanism for integration with other I/O and blocking.

By default it traces the calls it makes, and this can be suppressed with the `-q` option.

### 4.4.1 Usage

```
export ZF_ATTR=interface=ethX
zfsink localaddr:port
```

*localaddr* should be an IP address on interface *ethX*, or a multicast address. There are various additional options – run "`zfsink -h`" for details.

## 4.5 zftcpmtpong

The `zftcpmtpong` application demonstrates how to use TCPDirect in an application that does sends and receives on TCP sockets in separate threads.

By default it traces the calls it makes, and this can be suppressed with the `-q` option.

### 4.5.1 Usage

```
export ZF_ATTR=interface=ethX
zftcpmtpong localaddr:port
```

*localaddr* should be an IP address on interface *ethX*. This application accepts incoming TCP connections and waits for messages to arrive. It sends on each connection an equal number of bytes as are received (although not with the same contents).

## 4.6 Building the Example Applications

The TCPDirect example applications are built along with the Onload installation and should be present in the `openonload/build/gnu_x86_64/tests/zf_apps` subdirectory.

Source code for the example applications is in the `src/tests/zf_apps` subdirectory.

To rebuild the example applications use the following procedure:

```
cd openonload/scripts/
export PATH="$PWD:$PATH"
cd ../build/gnu_x86_64/tests/zf_apps/
make clean
make
```



## Chapter 5

# Using TCPDirect

This part of the documentation gives information on using TCPDirect to write and build applications.

### 5.1 Components

All components required to build and link a user application with the Solarflare TCPDirect API are distributed with Onload. When Onload is installed all required directories/files are located under the Onload distribution directory.

### 5.2 Compiling and Linking

#### 5.2.1 Header files

Applications or libraries using TCPDirect include the `zf.h` header which is installed into the system include directory. For example:

```
#include <zf/zf.h>
```

#### 5.2.2 Linking

The application will need to be linked either:

- with `libonload_zf_static.a` and `libciull.a`, to link statically, or
- with `libonload_zf.so`, to link dynamically.

All of the above libraries are deployed to the system library directory by `onload_install`.

If compiling your application against TCPDirect libraries distributed with one version of Onload, and running on a system with a different version of Onload, some care is required. TCPDirect currently preserves compatibility and provides a stable API between the user-space components and the kernel drivers, so that applications compiled against an older TCPDirect library will work when run with newer drivers. Compatibility in the other direction (newer TCPDirect libraries running with older drivers) is not guaranteed. Finally, TCPDirect does not support linking dynamically against one version of the libraries, and then running against another.

The simplest approach is to link statically to `libonload_zf_static.a`, as this ensures that the version of the library used will match the one you have compiled against.

For those wishing to use TCPDirect in combination with Onload, it is possible to link either statically or dynamically to TCPDirect and then to run the application with the `onload` wrapper in the usual way to allow the Onload intercepts to take effect.

### 5.2.3 Debugging

By default, the TCPDirect libraries are optimized for performance, and in particular perform only a minimum of logging and parameter-validation. To aid testing, debug versions of the TCPDirect libraries are provided, which *do* offer such validation and logging. As with the production libraries, these are available both as static and as shared libraries.

To use the static debug library, an application must be linked against it explicitly, rather than being linked against the production library. The debug library is not installed to the linker's default search path, and so the full path to the library must be passed to the linker. The debug library is named `libonload_zf_static.a`, as is the production library, but is installed to the `zf/debug` subdirectory of the system library directory (typically `/usr/lib64`).

To use the shared debug library, the application should link as normal against the shared library as described in the [Linking](#) section above, but when run should be invoked via the `zf_debug` wrapper. For example, an application called `app` linked against the shared TCPDirect library will use the production library when invoked as

```
app
```

and will use the debug library when invoked as

```
zf_debug app
```

By default, the debug libraries emit the same logging messages as do the production libraries, but the `log_level` attribute described in the [Attributes](#) chapter can be used to enable additional logging selectively. As a convenience, the `-l` option to the `zf_debug` wrapper will set this attribute to the specified value.

## 5.3 General

The majority of the functions in this API will return 0 on success, or a negative error code on failure. These are negated values of standard Linux error codes as defined in the system's `errno.h`. `errno` itself is not used.

Most of the API is non-blocking. The cases where this is not the case (e.g. [zf\\_muxer\\_wait\(\)](#)) are highlighted in the rest of this document.

The public API is defined by the headers in the `zf` subdirectory of the system include directory (typically `/usr/include`).

Attributes (defined by struct [zf\\_attr](#)) are used to pass configuration details through the API. This is similar to the existing SolarCapture attribute system.

The following sections discuss the most common operations. Zocket shutdown, obtaining addresses, and some other details are generally omitted for clarity – please refer to the suggested headers and example code for full details.

## 5.4 Using stacks

Before zockets can be created, the calling application must first create a stack using the following functions:

```
int zf_stack_alloc(struct zf_attr* attr, struct zf_stack** stack_out);
int zf_stack_free(struct zf_stack* stack);
```

The `attr` parameter to `zf_stack_alloc()` configures various aspects of the stack's behavior. In particular, the `interface` attribute specifies which network interface the stack should use, and the `n_bufs` attribute determines the total number of packet buffers allocated by the stack. Packet buffers are required to send and receive packets to and from the network, and also to queue packets on zockets for sending and receiving. A value of `n_bufs` that is too small can result in dropped packets and in various API calls failing with `ENOMEM`. Please see the [Attributes](#) chapter and the documentation for each API call for more details.

## 5.5 Using zockets

TCPDirect supports both TCP and UDP, but in contrast to the BSD sockets API the type of these zockets is explicit through the API types and function calls and UDP zockets are separated into receive (RX) and transmit (TX) parts.

## 5.6 UDP receive

First allocate a UDP receive zocket:

```
int zfur_alloc(struct zfur** us_out,
              struct zf_stack* st,
              const struct zf_attr* attr);
```

Then bind to associate the zocket with an address, port, and add filters:

```
int zfur_addr_bind(struct zfur* us,
                  struct sockaddr* laddr,
                  socklen_t laddrlen,
                  const struct sockaddr* raddr,
                  socklen_t raddrlen,
                  int flags);
```

Then receive packets:

```
int zfur_zc_recv(struct zfur* us,
                struct zfur_msg* msg,
                int flags);
```

`zfur_zc_recv()` will perform a zero-copy read of a single UDP datagram. The struct `zfur_msg` is completed to point to the buffers used by this message. Because it is zero-copy, the buffers used are locked (preventing re-use by the stack) until `zfur_zc_recv_done()` is called:

```
int zfur_zc_recv_done(struct zfur* us,
                     struct zfur_msg* msg);
```

### Note

These functions can all be found in [zf\\_udp.h](#).

## 5.7 UDP send

First allocate a UDP TX zocket, using the supplied addresses and ports:

```
int zfut_alloc(struct zfut** us_out,
              struct zf_stack* st,
              const struct sockaddr* laddr,
              socklen_t laddrlen,
              const struct sockaddr* raddr,
              socklen_t raddrlen,
              int flags,
              const struct zf_attr* attr);
```

Then perform a copy-based send (potentially using PIO) of a single datagram:

```
int zfut_send(struct zfut* us,
             const struct iovec* iov,
             int iov_cnt,
             int flags);
```

### Note

These functions can all be found in [zf\\_udp.h](#).

## 5.8 TCP listening

A TCP listening zocket can be created:

```
int zftl_listen(struct zf_stack* st,
              const struct sockaddr* laddr,
              socklen_t laddrlen,
              const struct zf_attr* attr,
              struct zftl** tl_out);
```

And a passively opened zocket accepted:

```
int zftl_accept(struct zftl* tl,
              struct zft** ts_out);
```

Listening zockets can be closed and freed:

```
int zftl_free(struct zftl* ts);
```

### Note

These functions can all be found in [zf\\_tcp.h](#).



## 5.9 TCP send and receive

Allocate a TCP (non-listening) zocket. Unlike UDP, this can be used for both send and receive:

```
int zft_alloc(struct zf_stack* st,
             const struct zf_attr* attr,
             struct zft_handle** handle_out);
```

Bind the zocket to a local address/port:

```
int zft_addr_bind(struct zft_handle* handle,
                 const struct sockaddr* laddr,
                 socklen_t laddrlen,
                 int flags);
```

Then connect the zocket to a remote address/port. Note that the supplied zocket handle is replaced with a different type as part of this operation. This function does not block (subsequent operations will return an error until it has completed).

```
int zft_connect(struct zft_handle* handle,
               const struct sockaddr* raddr,
               socklen_t raddrlen,
               struct zft** ts_out);
```

Perform a zero-copy receive on the connected TCP zocket:

```
int zft_zc_recv(struct zft* ts,
               struct zft_msg* msg,
               int flags);
```

The struct `zft_msg` is completed to point to the received message. Because it is zero-copy, this will lock the buffers used until the caller indicates that it has finished with them by calling:

```
void zft_zc_recv_done(struct zft* ts,
                    struct zft_msg* msg);
```

Alternatively a copy-based receive call can be made:

```
int zft_recv(struct zft* ts,
            struct iovec* iov_out,
            int iovcnt,
            int flags);
```

A copy-based send call can be made, and the supplied buffers reused immediately after this call returns:

```
int zft_send(struct zft* ts,
            const struct iovec* iov,
            int iov_cnt,
            int flags);
```

### Note

These functions can all be found in [zf\\_tcp.h](#).

## 5.10 Alternative Tx queues

Finally, for lowest latency on the fast path, a special API based around different alternative queues of data can be used. The TX alternative API is used to minimise latency on send, by pushing packets through the TX path on the NIC before a decision can be made whether they are needed.

```
int zf_alternatives_alloc(struct zf_stack* stack,
                        const struct zf_attr* attr,
                        zf_althandle* alt_out);
int zf_alternatives_release(struct zf_stack* stack,
                          zf_althandle alt);
int zf_alternatives_send(struct zf_stack* stack,
                       zf_althandle alt);
int zf_alternatives_cancel(struct zf_stack* stack,
                        zf_althandle alt);
int zft_alternatives_queue(struct zft* ts,
                        zf_althandle alt,
                        const struct iovec* iov,
                        int iov_cnt,
                        int flags);
unsigned zf_alternatives_free_space(struct zf_stack* stack,
                                  zf_althandle alt);
```

At the point when the decision to send is made the packet has already nearly reached the wire, minimising latency on the critical path.

Multiple queues are available for this, allowing alternative packets to be queued. Then when it is known what needs to be sent the appropriate alternative queue is selected. Packets queued on this are then sent to the wire.

When a packet is queued a handle is provided to allow future updates to the packet data. However, packet data update requires requeuing all packets on the affected alternative, so incurs a time penalty.

Here is an example, where there are 2 things that need updates, A and B, but it's not yet known which will be needed. The application has allocated 3 alternative queues, allowing them to queue updates for either A only, B only, or both:

```
zf_alternatives_alloc(ts, attr, &queue_a);
zft_alternatives_queue(ts, queue_a, <UpdateA_data>, flags);
zf_alternatives_alloc(ts, attr, &queue_b);
zft_alternatives_queue(ts, queue_b, <UpdateB_data>, flags);
zf_alternatives_alloc(ts, attr, &queue_ab);
zft_alternatives_queue(ts, queue_ab, <UpdateA_data>, flags);
zft_alternatives_queue(ts, queue_ab, <UpdateB_data>, flags);
```

After running the above code, the queues are as follows:

- queue\_a: <UpdateA\_data>
- queue\_b: <UpdateB\_data>
- queue\_ab: <UpdateA\_data><UpdateB\_data>

A single packet can only be queued on one alternative. In the example above each instance of an update is a separate buffer.

When it is known which update is required the application can select the appropriate alternative. The `zf_alternatives_send()` function is used to do this. This will send out the packets on the selected alternative. If other alternatives have queued packets, you must flush them without sending them, as the TCP headers will then be incorrect on these packets. The `zf_alternatives_cancel()` function is used to do this.

```
zf_alternatives_send(ts, queue_a);  
zf_alternatives_cancel(ts, queue_b);  
zf_alternatives_cancel(ts, queue_ab);
```

After running the above code, the packet containing <UpdateA\_data> has been sent from queue\_a, and all three queues are empty and available for re-use.

Packet data cannot be edited in place once a packet has been queued on an alternative. If a queued packet needs to be updated it must be requeued, together with all other packets currently queued on the alternative. The [zf\\_alternatives\\_cancel\(\)](#) and [zft\\_alternatives\\_queue\(\)](#) functions are used to do this.

To avoid having to wait for the original alternative to be canceled before re-use a replacement alternative can be supplied. The unwanted alternative could then be freed:

```
zf_alternatives_alloc(ts, attr, &queue_new_ab);  
zft_alternatives_queue(ts, queue_new_ab, <UpdateA_edited_data>, flags);  
zft_alternatives_queue(ts, queue_new_ab, <UpdateB_data>, flags);  
zf_alternatives_release(ts, queue_ab);
```

Before running the above code, queue\_ab contains unwanted data for editing:

- queue\_ab: <UpdateA\_data><UpdateB\_data>

After running the above code, queue\_new\_ab contains the new edited data, and queue\_ab has been freed:

- queue\_new\_ab: <UpdateA\_edited\_data><UpdateB\_data>

To determine the maximum packet size you can queue on an alternative, use the [zf\\_alternatives\\_free\\_space\(\)](#) function.

```
fs = zf_alternatives_free_space(ts, queue_ab);
```

## Note

These functions can all be found in [zf\\_alts.h](#).

## 5.11 Epoll – muxer.h

The multiplexer allows multiple zockets to be polled in a single operation. The multiplexer owes much of its design (and some of its datatypes) to epoll.

The basic unit of functionality is the multiplexer set implemented by [zf\\_muxer\\_set](#). Each type of zocket (e.g. UDP receive, UDP transmit, TCP listening, TCP) that can be multiplexed is equipped with a method for obtaining a [zf\\_waitable](#) that represents a given zocket:

```
struct zf_waitable* zfur_to_waitable(struct zfur* us);  
struct zf_waitable* zfut_to_waitable(struct zfut* us);  
struct zf_waitable* zftl_to_waitable(struct zftl* tl);  
struct zf_waitable* zft_to_waitable(struct zft* ts);
```

This `zf_waitable` can then be added to a multiplexer set by calling `zf_muxer_add()`. Each waitable can only exist in a single multiplexer set at once. Each multiplexer set can only contain waitables from a single stack.

```
int zf_muxer_add(struct zf_muxer_set*,
                struct zf_waitable* w,
                const struct epoll_event* event);
```

Having added all of the desired zockets to a set, the set can be polled using `zf_muxer_wait()`.

```
int zf_muxer_wait(struct zf_muxer_set*,
                  struct epoll_event* events,
                  int maxevents,
                  int64_t timeout);
```

This function polls a multiplexer set and populates an array of event descriptors representing the zockets in that set that are ready. The events member of each descriptor specifies the events for which the zocket is actually ready, and the data member is set to the user-data associated with that descriptor, as specified in the call to `zf_muxer_add()` or `zf_muxer_mod()`.

Before checking for ready zockets, the function calls `zf_reactor_perform()` on the set's stack in order to process events from the hardware. In contrast to the rest of the API, `zf_muxer_wait()` can block. The maximum time to block is specified timeout, and a value of zero results in non-blocking behaviour. A negative value for timeout will allow the function to block indefinitely. If the function, blocks, it will call `zf_reactor_perform()` repeatedly in a tight loop.

The multiplexer supports only edge-triggered events: that is, if `zf_muxer_wait()` reports that a zocket is ready, it will not do so again until a new event occurs on that zocket, even if the zocket is in fact ready.

Waitables already in a set can be modified:

```
int zf_muxer_mod(struct zf_waitable* w,
                 const struct epoll_event* event);
```

and deleted from the set:

```
int zf_muxer_del(struct zf_waitable* w);
```

These functions can all be found in `muxer.h`.

## 5.12 Stack polling

The majority of the calls in the API are non-blocking and for performance reasons do not attempt to speculatively process events on a stack. The API provides the following function to allow the calling application to request the stack process events. It will return zero if nothing user-visible occurred as a result, or greater than zero if something potentially user-visible happened (e.g. received packet delivered to a zocket, zocket became writeable, etc). It may return false positives, i.e. report that something user-visible occurred, when in fact it did not.

```
int zf_reactor_perform(struct zf_stack* st);
```

Any calls which block (e.g. `zf_muxer_wait()`) will make this call internally. The code examples at the end of this document show how `zf_reactor_perform()` can be used.

The API also provides the following function to determine whether a stack has work pending. It will return non-zero if the stack has work pending, and therefore the application should call `zf_reactor_perform()` or `zf_muxer_wait()`.

```
int zf_stack_has_pending_work(const struct zf_stack* st);
```

These functions can all be found in `zf_reactor.h`.

## 5.13 Miscellaneous

For TCP zockets you can discover the local and/or remote IP addresses and ports in use:

```
void zft_getname(struct zft* ts, struct sockaddr* laddr_out,
                 struct sockaddr* raddr_out);
```

These functions can all be found in [zf\\_tcp.h](#) and [zf\\_udp.h](#).

## 5.14 zf\_stackdump

TCPDirect does not use the Onload bypass datapaths, it uses its own datapath therefore traffic sent/received by TCPDirect stacks and zockets is NOT visible using tcpdump or onload\_tcpdump or onload\_stackdump.

The TCPDirect zf\_stackdump feature can be used to analyse stacks/zockets created by the TCPDirect application.

### 5.14.1 Usage

```
# zf_stackdump -h
zf_stackdump [command [stack_ids...]]
```

Commands:

```
list      List stack(s)
dump      Show state of stack(s)
```

The default command is 'list'. Commands iterate over all stacks if no stacks are specified on the command line.

```
enp4s0f0/0f0      id=10      pid=8845
```

```
# zf_stackdump dump
=====
name=enp4s0f0/0f0
pool: pkt_bufs_n=17536 free=17025
config: tcp_timewait_ticks=666 tcp_finwait_ticks=666
config: tcp_initial_cwnd=0 ms_per_tcp_tick=90
alts: n_alts=0
stats: ring_refill_nomem=0
nic0: vi=240 flags=0 intf=enp4s0f0 index=6 hw=1A1
txq: pio_buf_size=2048
=====
UDP RX enp4s0f0/0f0:0
filter: lcl=172.16.130.252:8012 rmt=0.0.0.0:0
rx: unread=1 begin=0 process=0 end=1
udp rx: release_n=1 q_drops=0
```

### 5.14.2 stackdump output: stack

Parameter	Description
pkt_bufs_n	num of packet buffers allocated to the stack.
free	num of free (available) packet buffers.
tcp_timewait_ticks	length of the TIME-WAIT timer in ticks.

tcp_finwait_ticks	length of the FIN-WAIT-2 timer in ticks.
tcp_initial_cwnd	size of TCP congestion window.
ms_per_tcp_tick	granularity of TCP timer in milliseconds.
n_alts	total number of TX alternatives allocated to this stack.
ring_refill_nomem	num times there were no free packet buffers to refill rx ring (increase buffers with n_bufs attr).
vi	VI being used by this stack.
flags	stack flags.
intf	physical interface being used.
index	index of the physical interface.
pio_buf_size	size (bytes) of a PIO buffer.

### 5.14.3 stackdump output: UDP RX

Parameter	Description
filter	identifies filters installed on the adapter.
rx unread	num packets received, but still in zocket buffer rx queue
rx begin	zf_rx_ring: oldest pkt buffer not yet read.
rx process	zf_rx_ring: oldest pkt buffer not yet processed - so buffer not yet reaped.
rx end	zf_rx_ring: index of the last pkt in the queue.
release_n	num zero-copy packet awaiting release.
q_drops	num packets dropped from the zockets rx queue.

### 5.14.4 stackdump output: UDP TX

Parameter	Description
UDP TX	local interface local_ip:port remote_ip:port.
path	dst server src server.
tx posted	num tx descriptors posted to the NIC.
tx completed	num tx descriptors that have completed TX.

### 5.14.5 stackdump output: TCP TX/RX

Parameter	Description
tx posted	num tx descriptors posted to the NIC.
tx completed	num tx descriptors that have completed TX.
rx	unread, begin, process end, see UDP RX.
flags	zocket flags.
flags_ack_delay	ACK flags TF_ACK_DELAY 0x01   TF_ACK_NOW 0x02   TF_INTR 0x04 (in fast recovery)   TF_ACK_NEXT 0x08.
parent	identifies the listeing zocket from which a passive-open zocket was accepted.

refcount	when a zocket is used both from the TCP state machine and the application. This allows us to track when both have finished using it, and it can be freed.
snd_nxt	next sequence num to send.
lastack	last acknowledged sequence num.
snd_wnd_max	size of size window advertised by the peer.
snd_wl1	max sequence number advertised.
snd_wl2	last sequence number acknowledged.
snd_lbb	sequence num of next byte to be buffered.
snd_right_edge	sequence num of TCP send window.
send	num segments held in the send buffer.
inflight	num segments sent, but not yet acknowledged.
qbegin	TCP segment at sendq start.
qmiddle	TCP segment at sendq middle.
qend	TCP segment at sendq end.
sndbuf	zocket send buffer size (bytes).
cwnd	size of congestion avoidance window.
ssthresh	slow start threshold - num bytes that have to be sent before exiting slow start.
mss_lim	max segment size limit set by peer.
rcv_nxt	next expected sequence number.
rcv_ann_wnd	receiver window to announce.
rcv_ann_right_edge	announced right edge of window.
mss	max segment size.
seq	sequence number used for RTT estimation.
sa	rtt variable: smoothed round trip time.
sv	rtt variable: round trip time variance estimate.
nrtx	num of RTO retransmission attempts - reset to zero when a new ACK is received.
dupacks	num duplicate acks received.
persist_backoff	num of zero send win probes - sends probe pkt to keep connection alive.
timers	active timers.
ooo: added	out of order pkt added to sendq.
ooo: removed	count removals from overflow including segments that become in-order.
ooo: replaced	out of order pkt replaced in sendq.
ooo: handling deferred	count of deferred out-of-order pkts.
ooo: dropped_nomem	num of out of order pkts dropped when memory allocation fails.
ooo: drop_overfilled	num of out of order pkts dropped to prevent buffer overflowing.
msg_more_send_delayed	num of times there was no send because of MSG_MORE flag.
send_nomem	num of times there were no free packet buffers to perform a send.





## Chapter 6

# Worked Examples

This part of the documentation examines simplified versions of [zfudppingpong](#) and [zftcppingpong](#). These are small applications which listen for packets and replies, with as low latency as possible.

### Note

These examples do not set the values of attributes programmatically. Instead, they are left with the values set from the defaults and the `ZF_ATTR` environment variable. In particular, it is necessary to set the value of the `interface` attribute in `ZF_ATTR` in order to use these examples. Fully-fledged applications might prefer instead to set attributes using (for example) [zf\\_attr\\_set\\_str\(\)](#). Please see the [Attributes](#) chapter and the documentation for [attr.h](#) for more information.

## 6.1 UDP ping pong example

In the following example various boiler plate code has been omitted for clarity. For a full usable example see `src/tests/zf_apps/zfudppingpong.c`.

```
void ping(struct zf_stack* stack, struct zfut* ut, struct zfur* ur)
{
    unsigned char data[1500];
    struct iovec siov = { data, 1};
    struct {
        struct zfur_msg zcr;
        struct iovec iov[2];
    } rd = { { .iovcnt = 2 } };

    uint64_t* ping_last_word = cfg.size>=8 ? ((uint64_t*)&data[cfg.size]) -1 : 0;

    siov.iov_len = cfg.size;
    if( ping_last_word )
        *ping_last_word = 0x1122334455667788;
    ZF_TRY(zfut_send(ut, &siov, 1, 0));

    for(int it = 0; it < cfg.itercount; ) {
        while(zf_reactor_perform(stack) == 0);

        rd.zcr.iovcnt = 2;
        ZF_TRY(zfur_zc_recv(ur, &rd.zcr, 0));

        if( rd.zcr.iovcnt == 0 )
            continue;

        it += rd.zcr.iovcnt;

        zfur_zc_recv_done(ur, &rd.zcr);
    }
}
```

```

    if( ping_last_word )
        ++ping_last_word;
    ZF_TRY(zfut_send(ut, &siov, 1, 0));
}

void pong(struct zf_stack* stack, struct zfut* ut, struct zfur* ur)
{
    struct {
        struct zfur_msg zcr;
        struct iovec iov[2];
    } rd = { { .iovcnt = 2 } };

    for(int it = 0; it < cfg.itercount; ) {
        while(zf_reactor_perform(stack) == 0);

        rd.zcr.iovcnt = 2;
        ZF_TRY(zfur_zc_rcv(ur, &rd.zcr, 0));

        if( rd.zcr.iovcnt == 0 )
            continue;

        /* in pong we reply with the same data */
        for( int i = 0 ; i < rd.zcr.iovcnt; ++i ) {
            ZF_TRY(zfut_send(ut, &rd.zcr.iov[i], 1, 0));
        }

        it += rd.zcr.iovcnt;
        zfur_zc_rcv_done(ur, &rd.zcr);
    }
}

int main(int argc, char* argv[])
{
    ZF_TRY(zf_init());

    struct zf_attr* attr;
    ZF_TRY(zf_attr_alloc(&attr));

    struct zf_stack* stack;
    ZF_TRY(zf_stack_alloc(attr, &stack));

    struct zfur* ur;
    ZF_TRY(zfur_alloc(&ur, stack, attr));

    struct sockaddr_in laddr = parse_addr((char*)cfg.laddr);
    struct sockaddr_in raddr = parse_addr((char*)cfg.raddr);
    ZF_TRY(zfur_addr_bind(ur, &laddr, &raddr, 0));

    struct zfut* ut;
    ZF_TRY(zfut_alloc(&ut, stack, &laddr, &raddr, 0, attr));

    (cfg.ping ? &ping : pong)(stack, ut, ur);

    return 0;
}

```

## 6.2 TCP pong example

In the following example some boiler plate code has been omitted for clarity. For a full usable example see `src/tests/zf_apps/zftcppingpong.c`, which includes the “ping” side, how to use the multiplexer, and other details.

```

int main(int argc, const char* argv[])
{
    ZF_TRY(zf_init());

    struct zf_attr* attr;
    ZF_TRY(zf_attr_alloc(&attr));

    struct zf_stack* stack;
    ZF_TRY(zf_stack_alloc(attr, &stack));

```

```

struct zft_handle* tcp_handle;
ZF_TRY(zft_alloc(stack, attr, &tcp_handle));

struct sockaddr_in laddr = {
    .sin_addr = { inet_addr(argv[1]) },
    .sin_port = htons(2000),
};
struct sockaddr_in raddr = {
    .sin_addr = { inet_addr(argv[2]) },
    .sin_port = htons(2000),
};
ZF_TRY(zft_addr_bind(tcp_handle, &laddr));

struct zft* tcp;
ZF_TRY(zft_connect(tcp_handle, &raddr, &tcp));

int first_recv = 1;
unsigned char data[1500];
struct iovec siov = { data, 1};
struct {
    struct zft_msg zcr;
    struct iovec __iov_tail;
} rd = { { .iovcnt = 1 } };

while(1) {
    while(zf_reactor_perform(stack) == 0);

    rd.zcr.iovcnt = 1;
    ZF_TRY(zft_zc_recv(tcp, &rd.zcr, 0));

    if( rd.zcr.iovcnt == 0 )
        continue;

    if( first_recv ) {
        first_recv = 0;

        siov.iov_len = rd.zcr.iov[0].iov_len;
        memcpy(data, ((char*)rd.zcr.iov[0].iov_base), siov.iov_len);
    }
    for( int i = 0 ; i < rd.zcr.iovcnt; ++i ) {
        ZF_TRY(zft_send(tcp, &siov, 1, 0));
    }
    zft_zc_recv_done(tcp, &rd.zcr);
}

ZF_TRY(zft_shutdown_tx(tcp));
ZF_TRY(zft_free(tcp));

return 0;
}

```



## Chapter 7

# Attributes

Many TCPDirect API functions take an *attribute object* of type [zf\\_attr](#). Each attribute object specifies a set of attributes, which are key-value pairs. These attributes are documented in this section.

Attribute	Description
<a href="#">alt_buf_size</a>	Amount of NIC-side buffer space to allocate for use with TCP alternatives on this VI.
<a href="#">alt_count</a>	Number of TCP alternatives to allocate on this VI.
<a href="#">arp_reply_timeout</a>	Maximum time to wait for ARP replies, in microseconds (approx).
<a href="#">interface</a>	Use this interface name as <a href="#">zf_stack</a> interface.
<a href="#">log_format</a>	Combination of flags: ZF_LF_STACK_NAME (0x1), ZF_LF_FRC(0x2), ZF_LF_TCP_TIME(0x4), ZF_LF_PROCESS(0x8)
<a href="#">log_level</a>	Bitmask to enable different log message levels for each logging component.
<a href="#">max_tcp_endpoints</a>	Sets the maximum number of TCP endpoints (i.e. struct zftl).
<a href="#">max_tcp_listen_endpoints</a>	Sets the maximum number of TCP endpoints (i.e. struct zftl).
<a href="#">max_tcp_syn_backlog</a>	Sets the maximum number of half-open connections maintained in the stack.
<a href="#">max_udp_rx_endpoints</a>	Sets the maximum number of UDP RX endpoints (i.e. struct zfur).
<a href="#">max_udp_tx_endpoints</a>	Sets the maximum number of UDP TX endpoints (i.e. struct zfut).
<a href="#">n_bufs</a>	Number of packet buffers to allocate for the stack.
<a href="#">name</a>	The object name.
<a href="#">reactor_spin_count</a>	Sets how many iterations of the event processing loop <a href="#">zf_reactor_perform()</a> will make (in the absence of any events) before returning.
<a href="#">rx_ring_max</a>	Set the size and maximum fill level of the RX descriptor ring, which provides buffering between the network adapter and software.

<a href="#">rx_ring_refill_batch_size</a>	Sets the number of packet buffers rx ring is refilled with on each <a href="#">zf_reactor_perform</a> call.
<a href="#">rx_ring_refill_interval</a>	Sets the frequency of rx buffer ring refilling during inner <a href="#">zf_reactor_perform()</a> loop.
<a href="#">tcp_alt_ack_rewind</a>	The maximum number of bytes by which outgoing ACKs will be allowed to go backwards when sending an alternative queue.
<a href="#">tcp_delayed_ack</a>	Enable TCP delayed ACK ("on" by default).
<a href="#">tcp_finwait_ms</a>	Length of TCP FIN-WAIT-2 timer in ms, 0 - disabled.
<a href="#">tcp_initial_cwnd</a>	The initial congestion window for new TCP zockets.
<a href="#">tcp_retries</a>	The maximum number of TCP retransmits if data is not acknowledged by the network peer in general case.
<a href="#">tcp_syn_retries</a>	The maximum number of TCP SYN retransmits during <a href="#">zft_connect()</a> .
<a href="#">tcp_synack_retries</a>	The maximum number of TCP SYN-ACK retransmits before incoming connection is dropped.
<a href="#">tcp_timewait_ms</a>	Length of TCP TIME-WAIT timer in ms.
<a href="#">tcp_wait_for_time_wait</a>	Do not consider a stack to be quiescent if there are any TCP zockets in the TIME_WAIT state.
<a href="#">tx_ring_max</a>	Set the size of the TX descriptor ring, which provides buffering between the software and the network adaptor.

## 7.1 alt\_buf\_size Attribute Reference

Amount of NIC-side buffer space to allocate for use with TCP alternatives on this VI.

### Detailed Description

#### Type

Integer.

#### Default

40960.

#### Relevant components

[zf\\_vi](#).

## 7.2 alt\_count Attribute Reference

Number of TCP alternatives to allocate on this VI.

## Detailed Description

### Type

Integer.

### Default

0.

### Relevant components

zf\_vi.

## 7.3 arp\_reply\_timeout Attribute Reference

Maximum time to wait for ARP replies, in microseconds (approx).

## Detailed Description

### Type

Integer.

### Default

1000.

### Relevant components

[zf\\_stack](#).

## 7.4 interface Attribute Reference

Use this interface name as [zf\\_stack](#) interface.

## Detailed Description

### Type

String.

### Default

none.

### Relevant components

[zf\\_stack](#).

## 7.5 log\_format Attribute Reference

Combination of flags: ZF\_LF\_STACK\_NAME (0x1), ZF\_LF\_FRC(0x2), ZF\_LF\_TCP\_TIME(0x4), ZF\_LF\_PROCESS(0x8)

## Detailed Description

### Type

Integer.

### Default

stack name and tcp time.

### Relevant components

[zf\\_stack](#).

## 7.6 log\_level Attribute Reference

Bitmask to enable different log message levels for each logging component.



## Detailed Description

The log message level for each component is specified using a separate 4 bit nibble within the bitmask. The value of each nibble is a bitwise combination of: 0(none), 0x1(errors), 0x2(warnings), 0x4(info), 0x8(trace - debug build only). The following components are available: stack (bits 0-3), TCP-rx (4-7), TCP-tx (8-11), TCP-connection (12-15), UDP-rx (16-19), UDP-tx (20-23), UDP-connection (24-27), muxer (28-31), pool (32-35), fast-path (36-39), timers (40-43), filters (44-47), cplane (48-51). E.g. 0xfff0 will enable all TCP related logging and disable all other logging.

### Type

Bitmask.

### Default

ERR-level on all components.

### Relevant components

[zf\\_stack](#).

## 7.7 max\_tcp\_endpoints Attribute Reference

Sets the maximum number of TCP endpoints (i.e. struct zft).

## Detailed Description

### Type

Integer.

### Default

64.

### Relevant components

[zf\\_stack](#).

## 7.8 max\_tcp\_listen\_endpoints Attribute Reference

Sets the maximum number of TCP endpoints (i.e. struct zftl).

### Detailed Description

#### Type

Integer.

#### Default

16.

#### Relevant components

[zf\\_stack](#).

## 7.9 max\_tcp\_syn\_backlog Attribute Reference

Sets the maximum number of half-open connections maintained in the stack.

### Detailed Description

#### Type

Integer.

#### Default

net.ipv4.tcp\_max\_syn\_backlog.

#### Relevant components

[zf\\_stack](#).

## 7.10 max\_udp\_rx\_endpoints Attribute Reference

Sets the maximum number of UDP RX endpoints (i.e. struct zfur).

## Detailed Description

### Type

Integer.

### Default

64.

### Relevant components

[zf\\_stack](#).

## 7.11 max\_udp\_tx\_endpoints Attribute Reference

Sets the maximum number of UDP TX endpoints (i.e. struct zfut).

## Detailed Description

### Type

Integer.

### Default

64.

### Relevant components

[zf\\_stack](#).

## 7.12 n\_bufs Attribute Reference

Number of packet buffers to allocate for the stack.

## Detailed Description

The optimal value for this parameter depends on the size of the RX and TX queues, the total number of zockets in the stack, the number of alternatives in use and the frequency at which the application polls the stack and reads pending data from zockets. 0 - use maximum the stack with given parameters can use.

### Type

Integer.

### Default

0.

### Relevant components

zf\_pool.

## 7.13 name Attribute Reference

The object name.

## Detailed Description

Object names are visible in log messages, but have no other effect.

### Type

String.

### Default

(none).

### Relevant components

[zf\\_stack](#), [zf\\_pool](#), [zf\\_vi](#).

## 7.14 reactor\_spin\_count Attribute Reference

Sets how many iterations of the event processing loop [zf\\_reactor\\_perform\(\)](#) will make (in the absence of any events) before returning.

### Detailed Description

The default value makes [zf\\_reactor\\_perform\(\)](#) briefly spin if there are no new events present. A higher number can give better latency, however [zf\\_reactor\\_perform\(\)](#) will take more time to return when no new events are present. The minimum value is 1, which disables spinning. This attribute also affects the cost of [zf\\_muxer\\_wait\(\)](#) when invoked with `timeout_ns=0`.

### Type

Integer.

### Default

128.

### Relevant components

[zf\\_stack](#).

## 7.15 rx\_ring\_max Attribute Reference

Set the size and maximum fill level of the RX descriptor ring, which provides buffering between the network adapter and software.

### Detailed Description

The RX ring sizes supported are 512, 1024, 2048 and 4096. The `n_bufs` attribute may need to be increased when changing this value.

### Type

Integer.

### Default

512.

**Relevant components**

[zf\\_vi](#).

## 7.16 `rx_ring_refill_batch_size` Attribute Reference

Sets the number of packet buffers rx ring is refilled with on each `zf_reactor_perform` call.

**Detailed Description**

Must be multiple of 8.

**Type**

Integer.

**Default**

16.

**Relevant components**

[zf\\_stack](#).

## 7.17 `rx_ring_refill_interval` Attribute Reference

Sets the frequency of rx buffer ring refilling during inner `zf_reactor_perform()` loop.

**Detailed Description**

Set to 1 to have the ring refilled at each iteration.

**Type**

Integer.

#### Default

1.

#### Relevant components

[zf\\_stack](#).

## 7.18 tcp\_alt\_ack\_rewind Attribute Reference

The maximum number of bytes by which outgoing ACKs will be allowed to go backwards when sending an alternative queue.

### Detailed Description

#### Type

Integer.

#### Default

64K.

#### Relevant components

[zf\\_stack](#).

## 7.19 tcp\_delayed\_ack Attribute Reference

Enable TCP delayed ACK ("on" by default).

### Detailed Description

#### Type

Integer.

**Default**

1.

**Relevant components**

[zf\\_stack](#).

## 7.20 tcp\_finwait\_ms Attribute Reference

Length of TCP FIN-WAIT-2 timer in ms, 0 - disabled.

**Detailed Description****Type**

Integer.

**Default**

net.ipv4.tcp\_fin\_timeout.

**Relevant components**

[zf\\_stack](#).

## 7.21 tcp\_initial\_cwnd Attribute Reference

The initial congestion window for new TCP sockets.

**Detailed Description****Type**

Integer.

**Default**

10 \* MSS.



#### Relevant components

[zf\\_stack](#).

## 7.22 tcp\_retries Attribute Reference

The maximum number of TCP retransmits if data is not acknowledged by the network peer in general case.

### Detailed Description

See also [tcp\\_synack\\_retries](#), [tcp\\_syn\\_retries](#).

#### Type

Integer.

#### Default

[net.ipv4.tcp\\_retries2](#).

#### Relevant components

[zf\\_stack](#).

## 7.23 tcp\_syn\_retries Attribute Reference

The maximum number of TCP SYN retransmits during [zft\\_connect\(\)](#).

### Detailed Description

#### Type

Integer.

#### Default

[net.ipv4.tcp\\_syn\\_retries](#).

**Relevant components**

[zf\\_stack](#).

## 7.24 tcp\_synack\_retries Attribute Reference

The maximum number of TCP SYN-ACK retransmits before incoming connection is dropped.

**Detailed Description****Type**

Integer.

**Default**

net.ipv4.tcp\_synack\_retries.

**Relevant components**

[zf\\_stack](#).

## 7.25 tcp\_timewait\_ms Attribute Reference

Length of TCP TIME-WAIT timer in ms.

**Detailed Description****Type**

Integer.

**Default**

net.ipv4.tcp\_fin\_timeout.

**Relevant components**

[zf\\_stack](#).

## 7.26 tcp\_wait\_for\_time\_wait Attribute Reference

Do not consider a stack to be quiescent if there are any TCP zockets in the TIME\_WAIT state.

### Detailed Description

("off" by default).

#### Type

Integer.

#### Default

0.

#### Relevant components

[zf\\_stack](#).

## 7.27 tx\_ring\_max Attribute Reference

Set the size of the TX descriptor ring, which provides buffering between the software and the network adaptor.

### Detailed Description

The requested value is rounded up to the next size supported by the adapter. At time of writing the ring sizes supported are 512, 1024 and 2048. The `n_bufs` attribute may need to be increased when changing this value.

#### Type

Integer.

#### Default

512.

#### Relevant components

[zf\\_vi](#).



## Chapter 8

# Data Structure Index

### 8.1 Data Structures

Here are the data structures with brief descriptions:

<a href="#">zf_attr</a>	Attribute object . . . . .	51
<a href="#">zf_muxer_set</a>	Multiplexer set . . . . .	52
<a href="#">zf_stack</a>	. . . . .	52
<a href="#">zf_waitable</a>	Abstract multiplexable object . . . . .	52
<a href="#">zft</a>	Opaque structure describing a TCP zocket that is connected . . . . .	53
<a href="#">zft_handle</a>	Opaque structure describing a TCP zocket that is passive and not connected . . . . .	53
<a href="#">zft_msg</a>	TCP zero-copy RX message structure . . . . .	54
<a href="#">zftl</a>	Opaque structure describing a TCP listening zocket . . . . .	55
<a href="#">zfur</a>	Opaque structure describing a UDP-receive zocket . . . . .	55
<a href="#">zfur_msg</a>	UDP zero-copy RX message structure . . . . .	56
<a href="#">zfut</a>	Opaque structure describing a UDP-transmit zocket . . . . .	57



## Chapter 9

# File Index

### 9.1 File List

Here is a list of all documented files with brief descriptions:

<a href="#">attr.h</a>	TCPDirect API for attribute objects . . . . .	59
<a href="#">muxer.h</a>	TCPDirect multiplexer . . . . .	64
<a href="#">types.h</a>	TCPDirect types . . . . .	68
<a href="#">x86.h</a>	TCPDirect x86-specific definitions . . . . .	69
<a href="#">zf.h</a>	TCPDirect top-level API . . . . .	69
<a href="#">zf_alts.h</a>	TCPDirect Alternative Sends API . . . . .	69
<a href="#">zf_platform.h</a>	TCPDirect platform API . . . . .	73
<a href="#">zf_reactor.h</a>	TCPDirect reactor API for processing stack events . . . . .	73
<a href="#">zf_stack.h</a>	TCPDirect stack API . . . . .	75
<a href="#">zf_tcp.h</a>	TCPDirect TCP API . . . . .	78
<a href="#">zf_udp.h</a>	TCPDirect UDP API . . . . .	90





## Chapter 10

# Data Structure Documentation

### 10.1 `zf_attr` Struct Reference

Attribute object.

```
#include <zf/attr.h>
```

#### 10.1.1 Detailed Description

Attribute object.

Attributes are used to specify optional behaviours and parameters, usually when allocating objects. Each attribute object defines a complete set of the attributes that the stack understands.

For example, the "max\_udp\_rx\_endpoints" attribute controls how many UDP-receive sockets can be created per [zf\\_stack](#).

The default values for attributes may be overridden by setting the environment variable ZF\_ATTR. For example:

```
** ZF_ATTR="interface=enp4s0f0;log_level=3"  
**
```

Each function that takes an attribute argument will only be interested in a subset of the attributes specified by an [zf\\_attr](#) instance. Other attributes are ignored.

The set of attributes supported may change between releases, so applications should where possible tolerate failures when setting attributes.

The documentation for this struct was generated from the following file:

- [attr.h](#)

## 10.2 `zf_muxer_set` Struct Reference

Multiplexer set.

```
#include <zf/muxer.h>
```

### 10.2.1 Detailed Description

Multiplexer set.

Represents multiple objects (including sockets) that can be polled simultaneously.

The documentation for this struct was generated from the following file:

- [muxer.h](#)

## 10.3 `zf_stack` Struct Reference

```
#include <zf/zf_stack.h>
```

### 10.3.1 Detailed Description

A stack encapsulates hardware and protocol state. It is the fundamental object used to drive TCPDirect. Individual objects for handling TCP and UDP traffic — *sockets* — are created within a stack.

#### See Also

[zf\\_stack\\_alloc\(\)](#)  
[zf\\_stack\\_free\(\)](#)  
[zf\\_reactor\\_perform\(\)](#)

The documentation for this struct was generated from the following file:

- [zf\\_stack.h](#)

## 10.4 `zf_waitable` Struct Reference

Abstract multiplexable object.

```
#include <zf/muxer.h>
```

### 10.4.1 Detailed Description

Abstract multiplexable object.

Zockets that can be added to a multiplexer set can be represented by a pointer of this type, which can be obtained by making the appropriate API call for the given zocket.

A waitable can also be retrieved for a stack by calling [zf\\_stack\\_to\\_waitable\(\)](#). Such waitables indicate whether a stack has quiesced, in the sense documented at [zf\\_stack\\_is\\_quiescent\(\)](#).

Definition at line 43 of file muxer.h.

The documentation for this struct was generated from the following file:

- [muxer.h](#)

## 10.5 zft Struct Reference

Opaque structure describing a TCP zocket that is connected.

```
#include <zf/zf_tcp.h>
```

### 10.5.1 Detailed Description

Opaque structure describing a TCP zocket that is connected.

Definition at line 142 of file zf\_tcp.h.

The documentation for this struct was generated from the following file:

- [zf\\_tcp.h](#)

## 10.6 zft\_handle Struct Reference

Opaque structure describing a TCP zocket that is passive and not connected.

```
#include <zf/zf_tcp.h>
```

### 10.6.1 Detailed Description

Opaque structure describing a TCP zocket that is passive and not connected.

The documentation for this struct was generated from the following file:

- [zf\\_tcp.h](#)

## 10.7 zft\_msg Struct Reference

TCP zero-copy RX message structure.

```
#include <zf/zf_tcp.h>
```

### Data Fields

- int [reserved](#) [4]
- int [pkts\\_left](#)
- int [flags](#)
- int [iovcnt](#)
- struct iovec [iov](#) [ZF\_FLEXIBLE\_ARRAY\_COUNT]

### 10.7.1 Detailed Description

TCP zero-copy RX message structure.

This structure is passed to [zft\\_zc\\_rcv\(\)](#), which will populate it with pointers to received packets.

Definition at line 379 of file [zf\\_tcp.h](#).

### 10.7.2 Field Documentation

#### 10.7.2.1 int flags

Reserved.

Definition at line 385 of file [zf\\_tcp.h](#).

#### 10.7.2.2 struct iovec iov[ZF\_FLEXIBLE\_ARRAY\_COUNT]

In: base of iovec array; out: filled with iovecs pointing to the payload of the received packets.

Definition at line 391 of file [zf\\_tcp.h](#).

#### 10.7.2.3 int iovcnt

In: Length of [iov](#) array expressed as a count of iovecs; out: number of entries of [iov](#) populated with pointers to packets.

Definition at line 388 of file [zf\\_tcp.h](#).

#### 10.7.2.4 `int pkts_left`

Out: Number of outstanding packets in the queue after this read.

Definition at line 383 of file `zf_tcp.h`.

#### 10.7.2.5 `int reserved[4]`

Reserved.

Definition at line 381 of file `zf_tcp.h`.

The documentation for this struct was generated from the following file:

- [zf\\_tcp.h](#)

## 10.8 `zftl` Struct Reference

Opaque structure describing a TCP listening zocket.

```
#include <zf/zf_tcp.h>
```

### 10.8.1 Detailed Description

Opaque structure describing a TCP listening zocket.

Definition at line 34 of file `zf_tcp.h`.

The documentation for this struct was generated from the following file:

- [zf\\_tcp.h](#)

## 10.9 `zfur` Struct Reference

Opaque structure describing a UDP-receive zocket.

```
#include <zf/zf_udp.h>
```

### 10.9.1 Detailed Description

Opaque structure describing a UDP-receive zocket.

Definition at line 24 of file `zf_udp.h`.

The documentation for this struct was generated from the following file:

- [zf\\_udp.h](#)

## 10.10 zfur\_msg Struct Reference

UDP zero-copy RX message structure.

```
#include <zf/zf_udp.h>
```

### Data Fields

- int [reserved](#) [4]
- int [dgrams\\_left](#)
- int [flags](#)
- int [iovcnt](#)
- struct iovec [iov](#) [ZF\_FLEXIBLE\_ARRAY\_COUNT]

### 10.10.1 Detailed Description

UDP zero-copy RX message structure.

This structure is passed to [zfur\\_zc\\_recv\(\)](#), which will populate it with pointers to received packets.

Definition at line 121 of file `zf_udp.h`.

### 10.10.2 Field Documentation

#### 10.10.2.1 int dgrams\_left

Out: Number of outstanding datagrams in the queue after this read.

Definition at line 125 of file `zf_udp.h`.

#### 10.10.2.2 int flags

Reserved.

Definition at line 127 of file `zf_udp.h`.

#### 10.10.2.3 struct iovec iov[ZF\_FLEXIBLE\_ARRAY\_COUNT]

In: base of iovec array; out: filled with iovecs pointing to the payload of the received packets.

Definition at line 133 of file `zf_udp.h`.

#### 10.10.2.4 int iovcnt

In: Length of [iov](#) array expressed as a count of iovecs; out: number of entries of [iov](#) populated with pointers to packets.

Definition at line 130 of file [zf\\_udp.h](#).

#### 10.10.2.5 int reserved[4]

Reserved.

Definition at line 123 of file [zf\\_udp.h](#).

The documentation for this struct was generated from the following file:

- [zf\\_udp.h](#)

## 10.11 zfut Struct Reference

Opaque structure describing a UDP-transmit zocket.

```
#include <zf/zf_udp.h>
```

### 10.11.1 Detailed Description

Opaque structure describing a UDP-transmit zocket.

A UDP-transmit zocket encapsulates the state required to send UDP datagrams. Each such zocket supports only a single destination address.

Definition at line 226 of file [zf\\_udp.h](#).

The documentation for this struct was generated from the following file:

- [zf\\_udp.h](#)





## Chapter 11

# File Documentation

### 11.1 attr.h File Reference

TCPDirect API for attribute objects.

#### Functions

- `int zf_attr_alloc (struct zf_attr **attr_out)`  
*Allocate an attribute object.*
- `void zf_attr_free (struct zf_attr *attr)`  
*Free an attribute object.*
- `void zf_attr_reset (struct zf_attr *attr)`  
*Return attributes to their default values.*
- `int zf_attr_set_int (struct zf_attr *attr, const char *name, int64_t val)`  
*Set an attribute to an integer value.*
- `int zf_attr_get_int (struct zf_attr *attr, const char *name, int64_t *val)`  
*Get an integer-valued attribute.*
- `int zf_attr_set_str (struct zf_attr *attr, const char *name, const char *val)`  
*Set an attribute to a string value.*
- `int zf_attr_get_str (struct zf_attr *attr, const char *name, char **val)`  
*Get a string-valued attribute.*
- `int zf_attr_set_from_str (struct zf_attr *attr, const char *name, const char *val)`  
*Set an attribute from a string value.*
- `int zf_attr_set_from_fmt (struct zf_attr *attr, const char *name, const char *fmt,...)`  
*Set an attribute to a string value (with formatting).*
- `struct zf_attr * zf_attr_dup (const struct zf_attr *attr)`  
*Duplicate an attribute object.*
- `int zf_attr_doc (const char *attr_name_opt, const char ***docs_out, int *docs_len_out)`  
*Returns documentation for an attribute.*

### 11.1.1 Detailed Description

TCPDirect API for attribute objects.

Definition in file [attr.h](#).

### 11.1.2 Function Documentation

#### 11.1.2.1 `int zf_attr_alloc ( struct zf_attr ** attr_out )`

Allocate an attribute object.

#### Parameters

<i>attr_out</i>	The attribute object is returned here.
-----------------	--

#### Returns

0 on success, or a negative error code:  
-ENOMEM if memory could not be allocated  
-EINVAL if the ZF\_ATTR environment variable is malformed.

#### 11.1.2.2 `int zf_attr_doc ( const char * attr_name_opt, const char *** docs_out, int * docs_len_out )`

Returns documentation for an attribute.

#### Parameters

<i>attr_name_opt</i>	The attribute name.
<i>docs_out</i>	On success, the resulting doc string output.
<i>docs_len_out</i>	On success, the length of the doc string output.

#### Returns

0 on success, or a negative error code.

#### 11.1.2.3 `struct zf_attr* zf_attr_dup ( const struct zf_attr * attr )`

Duplicate an attribute object.

#### Parameters

<i>attr</i>	The attribute object.
-------------	-----------------------

#### Returns

A new attribute object.

This function is useful when you want to make non-destructive changes to an existing attribute object.

#### 11.1.2.4 `void zf_attr_free ( struct zf_attr * attr )`

Free an attribute object.

#### Parameters

<i>attr</i>	The attribute object.
-------------	-----------------------

#### 11.1.2.5 int zf\_attr\_get\_int ( struct zf\_attr \* *attr*, const char \* *name*, int64\_t \* *val* )

Get an integer-valued attribute.

##### Parameters

<i>attr</i>	The attribute object.
<i>name</i>	Name of the attribute.
<i>val</i>	Value of the attribute (output).

##### Returns

0 on success, or a negative error code: -ENOENT if *name* is not a valid attribute name -EINVAL if *name* does not have an integer type

#### 11.1.2.6 int zf\_attr\_get\_str ( struct zf\_attr \* *attr*, const char \* *name*, char \*\* *val* )

Get a string-valued attribute.

##### Parameters

<i>attr</i>	The attribute object.
<i>name</i>	Name of the attribute.
<i>val</i>	Value of the attribute (output). This is allocated with strdup() and must be free()ed by the caller.

##### Returns

0 on success, or a negative error code: -ENOENT if *name* is not a valid attribute name -EINVAL if *name* does not have a string type

#### 11.1.2.7 void zf\_attr\_reset ( struct zf\_attr \* *attr* )

Return attributes to their default values.

##### Parameters

<i>attr</i>	The attribute object.
-------------	-----------------------

#### 11.1.2.8 int zf\_attr\_set\_from\_fmt ( struct zf\_attr \* *attr*, const char \* *name*, const char \* *fmt*, ... )

Set an attribute to a string value (with formatting).

#### Parameters

<i>attr</i>	The attribute object.
<i>name</i>	Name of the attribute.
<i>fmt</i>	Format string for the new attribute value.

#### Returns

- 0 on success, or a negative error code:
- ENOENT if *name* is not a valid attribute name
- EINVAL if it is not possible to convert *fmt* to a valid value for the attribute
- EOVERFLOW if *fmt* is not within the range of values this attribut can take.

This function behaves exactly as [zf\\_attr\\_set\\_from\\_str\(\)](#), except that the string value is generated from a printf()-style format string.

#### 11.1.2.9 int zf\_attr\_set\_from\_str ( struct zf\_attr \* *attr*, const char \* *name*, const char \* *val* )

Set an attribute from a string value.

#### Parameters

<i>attr</i>	The attribute object.
<i>name</i>	Name of the attribute.
<i>val</i>	New value for the attribute.

#### Returns

- 0 on success, or a negative error code:
- ENOENT if *name* is not a valid attribute name
- EINVAL if it is not possible to convert *val* to a valid value for the attribute
- EOVERFLOW if *val* is not within the range of values this attribut can take.

#### 11.1.2.10 int zf\_attr\_set\_int ( struct zf\_attr \* *attr*, const char \* *name*, int64\_t *val* )

Set an attribute to an integer value.

#### Parameters

<i>attr</i>	The attribute object.
<i>name</i>	Name of the attribute.
<i>val</i>	New value for the attribute.

#### Returns

- 0 on success, or a negative error code:
- ENOENT if *name* is not a valid attribute name
- EOVERFLOW if *val* is not within the range of values this attribute can take.

#### 11.1.2.11 int zf\_attr\_set\_str ( struct zf\_attr \* *attr*, const char \* *name*, const char \* *val* )

Set an attribute to a string value.

## Parameters

<i>attr</i>	The attribute object.
<i>name</i>	Name of the attribute.
<i>val</i>	New value for the attribute (may be NULL).

## Returns

- 0 on success, or a negative error code:
- ENOENT if *name* is not a valid attribute name
- # -ENOMSG if the attribute is not a string attribute.

## 11.2 muxer.h File Reference

TCPDirect multiplexer.

```
#include <sys/epoll.h>
```

## Functions

- int [zf\\_muxer\\_alloc](#) (struct [zf\\_stack](#) \*stack, struct [zf\\_muxer\\_set](#) \*\*muxer\_out)  
*Allocates a multiplexer set.*
- void [zf\\_muxer\\_free](#) (struct [zf\\_muxer\\_set](#) \*muxer)  
*Frees a multiplexer set.*
- int [zf\\_muxer\\_add](#) (struct [zf\\_muxer\\_set](#) \*muxer, struct [zf\\_waitable](#) \*w, const struct [epoll\\_event](#) \*event)  
*Adds a waitable object to a multiplexer set.*
- int [zf\\_muxer\\_mod](#) (struct [zf\\_waitable](#) \*w, const struct [epoll\\_event](#) \*event)  
*Modifies the event data for a waitable object in a multiplexer set.*
- int [zf\\_muxer\\_del](#) (struct [zf\\_waitable](#) \*w)  
*Removes a waitable object from a multiplexer set.*
- int [zf\\_muxer\\_wait](#) (struct [zf\\_muxer\\_set](#) \*muxer, struct [epoll\\_event](#) \*events, int maxevents, int64\_t timeout\_ns)  
*Polls a multiplexer set.*
- struct [epoll\\_event](#) \* [zf\\_waitable\\_event](#) (struct [zf\\_waitable](#) \*w)  
*Find out the [epoll\\_event](#) data in use with this waitable.*
- int [zf\\_waitable\\_fd\\_get](#) (struct [zf\\_stack](#) \*stack, int \*fd)  
*Create an fd that can be used within an [epoll](#) set or other standard muxer.*
- int [zf\\_waitable\\_fd\\_prime](#) (struct [zf\\_stack](#) \*stack)  
*Prime the fd before blocking.*

### 11.2.1 Detailed Description

TCPDirect multiplexer. The multiplexer, which allows multiple objects to be polled in a single operation.

The multiplexer allows multiple zockets to be polled in a single operation. The basic unit of functionality is the *multiplexer set* implemented by [zf\\_muxer\\_set](#). Each type of zocket that can be multiplexed is equipped with a method for obtaining a [zf\\_waitable](#) that represents a given zocket; this [zf\\_waitable](#) can then be added to a multiplexer set by calling [zf\\_muxer\\_add\(\)](#). Having added all of the desired zockets to a set, the set can be polled using [zf\\_muxer\\_wait\(\)](#).

The multiplexer owes much of its design (and some of its datatypes) to `epoll(7)`.

Definition in file [muxer.h](#).

## 11.2.2 Function Documentation

### 11.2.2.1 `int zf_muxer_add ( struct zf_muxer_set * muxer, struct zf_waitable * w, const struct epoll_event * event )`

Adds a waitable object to a multiplexer set.

#### Parameters

<i>muxer</i>	Multiplexer set.
<i>w</i>	Waitable to add.
<i>event</i>	Descriptor specifying the events that will be polled on the waitable, and the data to be returned when those events are detected.

#### Returns

- 0 on success, or a negative error code:
- EXDEV Waitable does not belong to the multiplexer set's stack.
- EALREADY Waitable is already in this multiplexer set.
- EBUSY Waitable is already in another multiplexer set.

Adds a waitable object to a multiplexer set. Each waitable may belong to at most one multiplexer set at a time. The events of interest are specified by `event.events`, which is a bitfield that should be populated from one or more of `EPOLLIN`, `EPOLLOUT`, `EPOLLHUP` and `EPOLLERR` as desired. `event.data` specifies the data to be returned to a caller of `zf_muxer_wait()` when that waitable is ready. Note that the waitable itself is not in general returned to such callers; if this is desired, then `event.data` must be set in such a way that the waitable can be determined.

#### Note

Unlike `epoll` functions in Linux, you have to explicitly set `EPOLLHUP` and `EPOLLERR` if you want to be notified about these events.

### 11.2.2.2 `int zf_muxer_alloc ( struct zf_stack * stack, struct zf_muxer_set ** muxer_out )`

Allocates a multiplexer set.

#### Parameters

<i>stack</i>	Stack to associate with multiplexer set.
<i>muxer_out</i>	Holds the address of the allocated multiplexer set on success.

#### Returns

- 0 on success, or a negative error code:
- ENOMEM Out of memory.

Allocates a multiplexer set, which allows multiple waitable objects to be polled in a single operation. Waitable objects, together with a mask of desired events, can be added to the set using `zf_muxer_add()`. The set can then be polled using `zf_muxer_wait()`.

### 11.2.2.3 `int zf_muxer_del ( struct zf_waitable * w )`

Removes a waitable object from a multiplexer set.

#### Parameters

<i>w</i>	Waitable to remove.
----------	---------------------

#### Returns

0 on success, or a negative error code:  
-EINVAL *w* has not been added to a multiplexer set.

#### Note

This operation should be avoided on fast paths.

#### 11.2.2.4 void zf\_muxer\_free ( struct zf\_muxer\_set \* *muxer* )

Frees a multiplexer set.

#### Parameters

<i>muxer</i>	The multiplexer set to free.
--------------	------------------------------

#### Note

If there are waitables in the set at the point at which it is freed, the underlying memory will not be freed until all of those waitables have been removed from the set. Nonetheless, the caller should never continue to use a pointer passed to this function.

#### 11.2.2.5 int zf\_muxer\_mod ( struct zf\_waitable \* *w*, const struct epoll\_event \* *event* )

Modifies the event data for a waitable object in a multiplexer set.

#### Parameters

<i>w</i>	Waitable to modify.
<i>event</i>	Descriptor specifying the events that will be polled on the waitable, and the data to be returned when those events are detected.

#### Returns

0 on success, or a negative error code:  
-EINVAL *w* has not been added to a multiplexer set.

#### See Also

[zf\\_muxer\\_add\(\)](#).

#### Note

This function can be used to re-arm waitable after it is returned by [zf\\_muxer\\_wait\(\)](#) if user likes something like level-triggered events:

```
**  zf_muxer_mod(w, zf_waitable_event(w));
**
```

#### 11.2.2.6 int zf\_muxer\_wait ( struct zf\_muxer\_set \* *muxer*, struct epoll\_event \* *events*, int *maxevents*, int64\_t *timeout\_ns* )

Polls a multiplexer set.



#### Parameters

<i>muxer</i>	Multiplexer set.
<i>events</i>	Array into which to return event descriptors.
<i>maxevents</i>	Maximum number of events to return.
<i>timeout_ns</i>	Maximum time in nanoseconds to block.

#### Returns

Number of events. Negative values are reserved for future use as error codes, but are not returned at present.

This function polls a multiplexer set and populates an array of event descriptors representing the waitables in that set that are ready. The *events* member of each descriptor specifies the events for which the waitable is actually ready, and the *data* member is set to the user-data associated with that descriptor, as specified in the call to [zf\\_muxer\\_add\(\)](#) or [zf\\_muxer\\_mod\(\)](#).

Before checking for ready objects, the function calls [zf\\_reactor\\_perform\(\)](#) on the set's stack in order to process events from the hardware. In contrast to the rest of the API, [zf\\_muxer\\_wait\(\)](#) can block. The maximum time to block is specified by *timeout\_ns*, and a value of zero results in non-blocking behaviour. A negative value for *timeout\_ns* will allow the function to block indefinitely. If the function blocks, it will call [zf\\_reactor\\_perform\(\)](#) repeatedly in a tight loop.

The multiplexer only supports edge-triggered events: that is, if [zf\\_muxer\\_wait\(\)](#) reports that a waitable is ready, it need not do so again until a *new* event occurs on tha waitable, even if the waitable is in fact ready. On the other hand, a waitable *may* be reported as ready even when a new event has not occurred, but only when the waitable is in fact ready. A transition from "not ready" to "ready" always constitutes an edge, and in particular, for *EPOLLIN*, the arrival of any new data constitutes an edge.

By default this function has relatively high CPU overhead when no events are ready to be processed and *timeout\_ns*==0, because it polls repeatedly for events. The amount of time spent polling is controlled by stack attribute *reactor\_spin\_count*. Setting *reactor\_spin\_count* to 1 disables polling and minimises the cost of [zf\\_muxer\\_wait\(timeout\\_ns=0\)](#).

#### 11.2.2.7 struct epoll\_event\* zf\_waitable\_event ( struct zf\_waitable \* w )

Find out the *epoll\_event* data in use with this waitable.

#### Parameters

<i>w</i>	Waitable to explore.
----------	----------------------

#### Returns

The event data.

#### Note

Function behaviour is undefined if the waitable is not a member of any multiplexer set.

#### 11.2.2.8 int zf\_waitable\_fd\_get ( struct zf\_stack \* stack, int \* fd )

Create an *fd* that can be used within an *epoll* set or other standard muxer.

**Parameters**

<i>stack</i>	Stack the fd should indicate activity for
<i>fd</i>	Updated on success to contain the fd to use

**Returns**

0 on success, or a negative error code. The possible error-codes are returned from system calls and are system-dependent.

This function creates a file descriptor that can be used within an epoll set (or other standard muxer such as poll or select) to be notified when there is activity on the corresponding stack.

The fd supplied may indicate readiness for a variety of reasons not directly related to the availability of data on a socket. For example, there is an event that needs processing, a timer has expired, or a connection has changed state. When this occurs the caller should ensure they call [zf\\_muxer\\_wait\(\)](#) to allow the required activity to take place, and discover if this affected any of the stack's sockets that the caller is interested in. This may or may not result in a socket within the stack becoming readable or writeable.

Freeing the [zf\\_stack](#) will release all the resources associated with this fd, so it must not be used afterwards. You do not need to call close() on the supplied fd, it will be closed when the stack is freed as part of the [zf\\_stack\\_free\(\)](#) call.

**11.2.2.9 int zf\_waitable\_fd\_prime ( struct zf\_stack \* *stack* )**

Prime the fd before blocking.

**Parameters**

<i>stack</i>	Stack that matches the fd
--------------	---------------------------

**Returns**

0 on success, or a negative error code. The possible error-codes are returned from system calls and are system-dependent.

This primes an fd previously allocated with [zf\\_waitable\\_fd\\_get\(\)](#) so it is ready for use with a standard muxer like epoll\_wait. The fd should be primed in this way each time the caller blocks waiting for activity.

## 11.3 types.h File Reference

TCPDirect types.

### 11.3.1 Detailed Description

TCPDirect types.

Definition in file [types.h](#).

## 11.4 x86.h File Reference

TCPDirect x86-specific definitions.

### 11.4.1 Detailed Description

TCPDirect x86-specific definitions. This file contains system-dependent code that is used that is used by the other header files. It has no end-user API.

Definition in file [x86.h](#).

## 11.5 zf.h File Reference

TCPDirect top-level API.

### 11.5.1 Detailed Description

TCPDirect top-level API. This file should be included in TCPDirect clients. It includes any other TCPDirect header files that are required.

Definition in file [zf.h](#).

## 11.6 zf\_alts.h File Reference

TCPDirect Alternative Sends API.

### Typedefs

- typedef uint64\_t [zf\\_althandle](#)  
*Opaque handle for an alternative.*

### Functions

- int [zf\\_alternatives\\_alloc](#) (struct [zf\\_stack](#) \*stack, const struct [zf\\_attr](#) \*attr, [zf\\_althandle](#) \*alt\_out)  
*Acquire an ID for an alternative queue.*
- int [zf\\_alternatives\\_release](#) (struct [zf\\_stack](#) \*stack, [zf\\_althandle](#) alt)  
*Release an ID for an alternative queue.*
- int [zf\\_alternatives\\_send](#) (struct [zf\\_stack](#) \*stack, [zf\\_althandle](#) alt)  
*Select an alternative and send those messages.*
- int [zf\\_alternatives\\_cancel](#) (struct [zf\\_stack](#) \*stack, [zf\\_althandle](#) alt)  
*Cancel an alternative.*
- int [zft\\_alternatives\\_queue](#) (struct [zft](#) \*ts, [zf\\_althandle](#) alt, const struct iovec \*iov, int iov\_cnt, int flags)  
*Queue a TCP message for sending.*
- unsigned [zf\\_alternatives\\_free\\_space](#) (struct [zf\\_stack](#) \*stack, [zf\\_althandle](#) alt)  
*Query the amount of free buffering on an alt.*
- int [zf\\_alternatives\\_query\\_overhead\\_tcp](#) (struct [zft](#) \*ts, struct [ef\\_vi\\_transmit\\_alt\\_overhead](#) \*out)  
*Query TCP per-packet overhead parameters.*

## 11.6.1 Detailed Description

TCPDirect Alternative Sends API.

Definition in file [zf\\_alts.h](#).

## 11.6.2 Function Documentation

### 11.6.2.1 `int zf_alternatives_alloc ( struct zf_stack * stack, const struct zf_attr * attr, zf_althandle * alt_out )`

Acquire an ID for an alternative queue.

#### Parameters

<i>stack</i>	Stack to allocate the alternative for
<i>attr</i>	Requested attributes for the alternative. At the present time, the attributes are unused. Refer to the attribute documentation in <a href="#">Attributes</a> for details.
<i>alt_out</i>	Handle for the allocated alternative

#### Returns

0 Success  
-ENOMEM No alternative queues available

The alternative queue is identified by opaque handles, and is only able to be used with zockets in the stack provided to this function.

The number of alternatives available to a stack is controlled by the value of the `alt_count` attribute used when creating the stack. This value defaults to zero.

#### See Also

[zf\\_alternatives\\_release\(\)](#)

### 11.6.2.2 `int zf_alternatives_cancel ( struct zf_stack * stack, zf_althandle alt )`

Cancel an alternative.

#### Parameters

<i>stack</i>	Stack the alternative was allocated on
<i>alt</i>	Selected alternative

#### Returns

0 Success

Drops messages queued on this alternative without sending.

You can reuse the alternative queue immediately for new messages (including messages on a different zocket from the previous use) but [zft\\_alternatives\\_queue\(\)](#) may return -EBUSY until the cancel operation is completed.

### 11.6.2.3 `unsigned zf_alternatives_free_space ( struct zf_stack * stack, zf_althandle alt )`

Query the amount of free buffering on an alt.

#### Parameters

<i>stack</i>	Stack the alternative was allocated on
<i>alt</i>	Selected alternative

#### Returns

Number of bytes available

The return value of this function is the payload size in bytes of the largest packet which can be sent into this alternative at this moment. Larger packets than this will cause -ENOMEM errors from functions which queue data on alternatives.

Due to per-packet and other overheads, this amount may be different on different alternatives, and is not guaranteed to rise and fall by exactly the sizes of packets queued and sent.

The returned value includes all packet headers. The maximum length of data accepted by [zft\\_alternatives\\_queue\(\)](#) will be lower than this by the size of the TCP+IP+Ethernet headers. To find a zocket's header size, use [zft\\_get\\_header\\_size\(\)](#) or [zfut\\_get\\_header\\_size\(\)](#).

**11.6.2.4** `int zf_alternatives_query_overhead_tcp ( struct zft * ts, struct ef_vi_transmit_alt_overhead * out )`

Query TCP per-packet overhead parameters.

#### Parameters

<i>ts</i>	TCP connection to be queried
<i>out</i>	Returned overhead parameters

#### Returns

0 on success or -EINVAL if this stack doesn't support alternatives.

This function returns a set of parameters which can be used with [ef\\_vi\\_transmit\\_alt\\_usage\(\)](#) to calculate the amount of buffer space used when sending data via TCP, taking into account the space taken up by headers, VLAN tags, IP options etc.

Use of this function in this way assumes that the transmitted data fits entirely into a single TCP packet.

See the documentation for [ef\\_vi\\_transmit\\_alt\\_usage\(\)](#) for more.

**11.6.2.5** `int zf_alternatives_release ( struct zf_stack * stack, zf_althandle alt )`

Release an ID for an alternative queue.

#### Parameters

<i>stack</i>	Stack to release the alternative for
<i>alt</i>	zf_alternative to release

#### Returns

0 Success

Releases allocated alternative queue. If any messages are queued on the specified queue they will be flushed without being sent.

#### See Also

[zf\\_alternatives\\_alloc\(\)](#)

#### 11.6.2.6 int zf\_alternatives\_send ( struct zf\_stack \* *stack*, zf\_althandle *alt* )

Select an alternative and send those messages.

#### Parameters

<i>stack</i>	Stack the alternative was allocated on
<i>alt</i>	Selected alternative

#### Returns

0 Success

-EBUSY Unable to send due to a transient state (e.g. the alternative queue is being refreshed in response to receiving data).

-EINVAL Unable to send due to inconsistent TCP state (e.g. the zocket is not connected, or has been used via the normal send path after queueing messages on this alternative queue)

On success messages queued on the selected alternative are sent. If other alternative queues have messages queued for the same zocket, their headers will now be out of date and you must call [zf\\_alternatives\\_cancel\(\)](#) on those queues. You are free to reuse this alternative queue, but until it has finished sending the current set of messages calls to [zft\\_alternatives\\_queue\(\)](#) will return -EBUSY.

#### 11.6.2.7 int zft\_alternatives\_queue ( struct zft \* *ts*, zf\_althandle *alt*, const struct iovec \* *iov*, int *iov\_cnt*, int *flags* )

Queue a TCP message for sending.

#### Parameters

<i>ts</i>	TCP zocket
-----------	------------

<i>alt</i>	ID of the queue to push this message to. Must have been allocated via <a href="#">zf_alternatives_alloc()</a>
<i>iov</i>	TCP payload data to send in this message.
<i>iov_cnt</i>	Number of iovecs to send. Currently must be 1.
<i>flags</i>	Reserved for future use; must be zero.

## Returns

0 Success

-EAGAIN Unable to queue due to a transient problem, e.g. the TCP send queue is not empty. These errors may remain present for many milliseconds; the caller should decide whether to retry immediately or to perform other work in the meantime.

-EBUSY Unable to queue due to a transient problem, e.g. the alternative queue is still draining from a previous operation. These errors are expected to clear quickly without outside intervention; the caller can react by calling [zf\\_reactor\\_perform\(\)](#) and retrying the operation.

-EMSGSIZE Enqueuing the message would exceed the total congestion window.

-ENOMEM Unable to queue due to all packet buffers being allocated already.

-ENOBUFFS Unable to queue due to a lack of available buffer space, either in TCP Direct or in the NIC hardware.

-EINVAL Invalid parameters. This includes the case where the alternative already has data queued on another socket.

This function behaves similarly to [zft\\_send\(\)](#), but doesn't actually put the data on the wire.

For now it is only possible to send a single buffer of data in each call to [zft\\_alternatives\\_queue\(\)](#); this function will return -EINVAL if 'iov\_cnt' is not equal to 1. Future releases may change this. Multiple messages can be queued for sending on a single alternative by calling [zft\\_alternatives\\_queue\(\)](#) for each message.

The current implementation limits all messages enqueued on an alternative to be from the same socket. This may change in future.

In some cases where an alternative is in the middle of an operation such as a send, cancel, etc. this function may return -EBUSY. In this case the caller should process some events and retry.

## 11.7 zf\_platform.h File Reference

TCPDirect platform API.

### 11.7.1 Detailed Description

TCPDirect platform API. This file contains platform-dependent code that is used by the other header files. It has no end-user API.

Definition in file [zf\\_platform.h](#).

## 11.8 zf\_reactor.h File Reference

TCPDirect reactor API for processing stack events.

## Functions

- `int zf_reactor_perform (struct zf_stack *st)`  
*Process events on a stack.*
- `int zf_stack_has_pending_work (const struct zf_stack *st)`  
*Determine whether a stack has work pending.*

### 11.8.1 Detailed Description

TCPDirect reactor API for processing stack events.

Definition in file [zf\\_reactor.h](#).

### 11.8.2 Function Documentation

#### 11.8.2.1 `int zf_reactor_perform ( struct zf_stack * st )`

Process events on a stack.

##### Parameters

<i>st</i>	Stack for which to process events.
-----------	------------------------------------

This function processes events on a stack and performs the necessary handling. These events include transmit and receive events raised by the hardware, and also software events such as TCP timers. Applications must call this function or [zf\\_muxer\\_wait\(\)](#) frequently for each stack that is in use. Please see [Stack polling](#) in the User Guide for further information.

By default this function has relatively high CPU overhead when no events are ready to be processed, because it polls repeatedly for events. The amount of time spent polling is controlled by stack attribute `reactor_spin_count`. Setting `reactor_spin_count` to 1 disables polling and minimises the cost of [zf\\_reactor\\_perform\(\)](#).

##### Returns

0 if nothing user-visible occurred as a result.

>0 if something user-visible might have occurred as a result.

Here, "something user-visible occurred" means that the event-processing just performed has had an effect that can be seen by another API call: for example, new data might have arrived on a socket, in which case that data can be retrieved by one of the receive functions. False positives are possible: a value greater than zero indicates to the application that it should process its sockets, but it does not guarantee that this will yield anything new. Finer-grained advertisement of interesting events can be achieved using the multiplexer.

##### See Also

[zf\\_muxer\\_wait\(\)](#)

#### 11.8.2.2 `int zf_stack_has_pending_work ( const struct zf_stack * st )`

Determine whether a stack has work pending.



#### Parameters

<code>st</code>	Stack to check for pending work.
-----------------	----------------------------------

This function returns non-zero if the stack has work pending, and therefore the application should call [zf\\_reactor\\_perform\(\)](#) or [zf\\_muxer\\_wait\(\)](#).

This function can be called concurrently with other calls on a stack, and so can be used to avoid taking a serialisation lock (and therefore avoid inducing lock contention) when there isn't any work to do.

#### Returns

- 0 if there is nothing to do.
- >0 if there is some work pending.

#### See Also

[zf\\_reactor\\_perform\(\)](#) [zf\\_muxer\\_wait\(\)](#)

## 11.9 zf\_stack.h File Reference

TCPDirect stack API.

#### Macros

- `#define EPOLLSTACKHUP EPOLLRDHUP`  
*Event indicating stack quiescence.*

#### Functions

- `int zf_init (void)`  
*Initialize zf library.*
- `int zf_deinit (void)`  
*Deinitialize zf library.*
- `int zf_stack_alloc (struct zf_attr *attr, struct zf_stack **stack_out)`  
*Allocate a stack with the supplied attributes.*
- `int zf_stack_free (struct zf_stack *stack)`  
*Free a stack previously allocated with [zf\\_stack\\_alloc\(\)](#).*
- `struct zf_waitable * zf_stack_to_waitable (struct zf_stack *)`  
*Returns a waitable object representing the quiescence of a stack.*
- `int zf_stack_is_quiescent (struct zf_stack *)`  
*Returns a boolean value indicating whether a stack is quiescent.*
- `void zf_version (void)`  
*Print library name and version to stderr.*

### 11.9.1 Detailed Description

TCPDirect stack API.

Definition in file [zf\\_stack.h](#).

### 11.9.2 Macro Definition Documentation

#### 11.9.2.1 `#define EPOLLSTACKHUP EPOLLRDHUP`

Event indicating stack quiescence.

#### See Also

[zf\\_stack\\_to\\_waitable\(\)](#)

Definition at line 85 of file [zf\\_stack.h](#).

### 11.9.3 Function Documentation

#### 11.9.3.1 `int zf_deinit ( void )`

Deinitialize zf library.

#### Returns

0. Negative values are reserved for future use as error returns.

#### 11.9.3.2 `int zf_init ( void )`

Initialize zf library.

#### Returns

0 on success, or a negative error code. This function uses attributes internally and can return any of the error codes returned by [zf\\_attr\\_alloc\(\)](#). Additionally, it can return the following:  
-ENOENT Failed to initialize control plane. A likely cause is that Onload drivers are not loaded.

#### 11.9.3.3 `int zf_stack_alloc ( struct zf_attr * attr, struct zf_stack ** stack_out )`

Allocate a stack with the supplied attributes.

#### Parameters

<i>attr</i>	A set of properties to apply to the stack.
<i>stack_out</i>	A pointer to the newly allocated stack.

A stack encapsulates hardware and protocol state. A stack binds to a single network interface, specified by the `interface` attribute in `attr`. To process events on a stack, call `zf_reactor_perform()` or `zf_muxer_wait()`.

Relevant attributes to set in `attr` are those in the `zf_stack`, `zf_pool` and `zf_vi` categories described in the attributes documentation in [Attributes](#).

#### Returns

0 on success, or a negative error code:

- EBUSY Out of VI instances or resources for alternatives.
- EINVAL Attribute out of range.
- ENODEV Interface was not specified or was invalid.
- ENOENT Failed to initialize `ef_vi` or Onload libraries. A likely cause is that Onload drivers are not loaded.
- ENOKEY Adapter is not licensed for TCPDirect.
- ENOMEM Out of memory. N.B. Huge pages are required.
- ENOSPC Out of PIO buffers.

Errors from system calls are also possible. Please consult your system's documentation for `errno(3)`.

#### 11.9.3.4 `int zf_stack_free ( struct zf_stack * stack )`

Free a stack previously allocated with `zf_stack_alloc()`.

#### Parameters

<i>stack</i>	Stack to free
--------------	---------------

#### Returns

When called with a valid stack, this function always returns zero. Results on invalid stacks are undefined.

#### 11.9.3.5 `int zf_stack_is_quiescent ( struct zf_stack * )`

Returns a boolean value indicating whether a stack is quiescent.

A stack is quiescent precisely when all of the following are true:

- the stack will not transmit any packets except in response to external stimuli (including relevant API calls),
- closing zockets will not result in the transmission of any packets, and
- (optionally, controlled by the `tcp_wait_for_time_wait` stack attribute) there are no TCP zockets in the `TIME_WAIT` state. In practice, this is equivalent altogether to the condition that there are no open TCP connections.

This can be used to ensure that all connections have been closed gracefully before destroying a stack (or exiting the application). Destroying a stack while it is not quiescent is permitted by the API, but when doing so there is no guarantee that sent data has been acknowledged by the peer or even transmitted, and there is the possibility that peers' connections will be reset.

#### See Also

[zf\\_stack\\_to\\_waitable\(\)](#)

#### Returns

Non-zero if the stack is quiescent, or zero otherwise.

### 11.9.3.6 struct zf\_waitable\* zf\_stack\_to\_waitable ( struct zf\_stack \* )

Returns a waitable object representing the quiescence of a stack.

The waitable will be ready for [EPOLLSTACKHUP](#) if the stack is quiescent.

#### See Also

[zf\\_stack\\_is\\_quiescent\(\)](#)

#### Returns

Waitable.

## 11.10 zf\_tcp.h File Reference

TCPDirect TCP API.

```
#include <netinet/in.h>
#include <sys/uio.h>
```

### Data Structures

- struct [zftl](#)  
*Opaque structure describing a TCP listening socket.*
- struct [zft](#)  
*Opaque structure describing a TCP socket that is connected.*
- struct [zft\\_msg](#)  
*TCP zero-copy RX message structure.*

### Functions

- int [zftl\\_listen](#) (struct [zf\\_stack](#) \*st, const struct sockaddr \*laddr, socklen\_t laddrlen, const struct [zf\\_attr](#) \*attr, struct [zftl](#) \*\*tl\_out)  
*Allocate TCP listening socket.*
- int [zftl\\_accept](#) (struct [zftl](#) \*tl, struct [zft](#) \*\*ts\_out)  
*Accept incoming TCP connection.*
- struct [zf\\_waitable](#) \* [zftl\\_to\\_waitable](#) (struct [zftl](#) \*tl)  
*Returns a [zf\\_waitable](#) representing tl.*
- void [zftl\\_getname](#) (struct [zftl](#) \*ts, struct sockaddr \*laddr\_out, socklen\_t \*laddrlen)  
*Retrieve the local address of the socket.*
- int [zftl\\_free](#) (struct [zftl](#) \*ts)  
*Release resources associated with a TCP listening socket.*
- struct [zf\\_waitable](#) \* [zft\\_to\\_waitable](#) (struct [zft](#) \*ts)  
*Returns a [zf\\_waitable](#) representing the given zft.*
- int [zft\\_alloc](#) (struct [zf\\_stack](#) \*st, const struct [zf\\_attr](#) \*attr, struct [zft\\_handle](#) \*\*handle\_out)

- Allocate active-open TCP socket.*
- int [zft\\_handle\\_free](#) (struct [zft\\_handle](#) \*handle)  
*Release a handle to a TCP socket.*
- void [zft\\_handle\\_getname](#) (struct [zft\\_handle](#) \*ts, struct sockaddr \*laddr\_out, socklen\_t \*laddrlen)  
*Retrieve the local address to which a [zft\\_handle](#) is bound.*
- int [zft\\_addr\\_bind](#) (struct [zft\\_handle](#) \*handle, const struct sockaddr \*laddr, socklen\_t laddrlen, int flags)  
*Bind to a specific local address.*
- int [zft\\_connect](#) (struct [zft\\_handle](#) \*handle, const struct sockaddr \*raddr, socklen\_t raddrlen, struct [zft](#) \*\*ts\_out)  
*Connect a TCP socket.*
- int [zft\\_shutdown\\_tx](#) (struct [zft](#) \*ts)  
*Shut down outgoing TCP connection.*
- int [zft\\_free](#) (struct [zft](#) \*ts)  
*Release resources associated with a TCP socket.*
- int [zft\\_state](#) (struct [zft](#) \*ts)  
*Return the TCP state of a TCP socket.*
- int [zft\\_error](#) (struct [zft](#) \*ts)  
*Find out the error type happened on the TCP socket.*
- void [zft\\_getname](#) (struct [zft](#) \*ts, struct sockaddr \*laddr\_out, socklen\_t \*laddrlen, struct sockaddr \*raddr\_out, socklen\_t \*raddrlen)  
*Retrieve the local address of the socket.*
- void [zft\\_zc\\_rcv](#) (struct [zft](#) \*ts, struct [zft\\_msg](#) \*msg, int flags)  
*Zero-copy read of available packets.*
- int [zft\\_zc\\_rcv\\_done](#) (struct [zft](#) \*ts, struct [zft\\_msg](#) \*msg)  
*Concludes pending [zc\\_rcv](#) operation as done.*
- int [zft\\_zc\\_rcv\\_done\\_some](#) (struct [zft](#) \*ts, struct [zft\\_msg](#) \*msg, size\_t len)  
*Concludes pending [zc\\_rcv](#) operation as done acknowledging all or some of the data to have been read.*
- int [zft\\_rcv](#) (struct [zft](#) \*ts, const struct iovec \*iov, int iovcnt, int flags)  
*Copy-based receive.*
- ssize\_t [zft\\_send](#) (struct [zft](#) \*ts, const struct iovec \*iov, int iov\_cnt, int flags)  
*Send data specified in iovec array.*
- ssize\_t [zft\\_send\\_single](#) (struct [zft](#) \*ts, const void \*buf, size\_t buflen, int flags)  
*Send data given in single buffer.*
- int [zft\\_send\\_space](#) (struct [zft](#) \*ts, size\_t \*space)  
*Query available space in the send queue.*
- int [zft\\_get\\_mss](#) (struct [zft](#) \*ts)  
*Retrieve the maximum segment size (MSS) for a TCP connection.*
- unsigned [zft\\_get\\_header\\_size](#) (struct [zft](#) \*ts)  
*Return protocol header size for this connection.*

### 11.10.1 Detailed Description

TCPDirect TCP API.

Definition in file [zf\\_tcp.h](#).

### 11.10.2 Function Documentation

#### 11.10.2.1 int [zft\\_addr\\_bind](#) ( struct [zft\\_handle](#) \* *handle*, const struct sockaddr \* *laddr*, socklen\_t *laddrlen*, int *flags* )

Bind to a specific local address.

#### Parameters

<i>handle</i>	TCP zocket handle.
<i>laddr</i>	Local address.
<i>laddrlen</i>	Length of structure pointed to by <i>laddr</i>
<i>flags</i>	Reserved. Must be zero.

#### Returns

- 0 Success.
- EADDRINUSE Local address already in use.
- EADDRNOTAVAIL *laddr* is not a local address.
- EAFNOSUPPORT *laddr* is not an AF\_INET address.
- EFAULT Invalid pointer.
- EINVAL Zocket is already bound, invalid *flags*, or invalid *laddrlen*.
- ENOMEM Out of memory.

#### 11.10.2.2 `int zft_alloc ( struct zf_stack * st, const struct zf_attr * attr, struct zft_handle ** handle_out )`

Allocate active-open TCP zocket.

#### Parameters

<i>st</i>	Initialized <a href="#">zf_stack</a> .
<i>attr</i>	Attributes required for this TCP zocket. Note that not all attributes are relevant; only those which apply to objects of type "zf_socket" are applicable here. Refer to the <a href="#">Attributes</a> documentation for details.
<i>handle_out</i>	On successful return filled with pointer to a zocket handle. This handle can be used to refer to the zocket before it is connected.

#### Returns

- 0 Success.
- ENOBUFS No zockets of this type available.

This function initialises the datastructures needed to make an outgoing TCP connection.

The returned handle can be used to refer to the zocket before it is connected.

The handle must be released either by explicit release with [zft\\_handle\\_free\(\)](#), or by conversion to a connected zocket via [zft\\_connect\(\)](#).

#### See Also

[zft\\_addr\\_bind\(\)](#) [zft\\_connect\(\)](#) [zft\\_handle\\_free\(\)](#)

#### 11.10.2.3 `int zft_connect ( struct zft_handle * handle, const struct sockaddr * raddr, socklen_t raddrlen, struct zft ** ts_out )`

Connect a TCP zocket.

#### Parameters

<i>handle</i>	TCP zocket handle, to be replaced by the returned zocket.
<i>raddr</i>	Remote address to connect to.
<i>raddrlen</i>	Length of structure pointed to by raddr.
<i>ts_out</i>	On successful return, a pointer to a TCP zocket.

This replaces the zocket handle with a TCP zocket. On successful return the zocket handle has been released and is no longer valid.

If a specific local address has not been set via [zft\\_addr\\_bind\(\)](#) then an appropriate one will be selected.

This function does not block. Functions that attempt to transfer data on the zocket between [zft\\_connect\(\)](#) and the successful establishment of the underlying TCP connection will return an error. Furthermore, failure of the remote host to accept the connection will not be reported by this function, but instead by any attempts to read from the zocket (or by [zft\\_error\(\)](#)). As such, after calling [zft\\_connect\(\)](#), either

- read calls that fail with `-ENOTCONN` should be repeated after calling [zf\\_reactor\\_perform\(\)](#), or
- the zocket should be polled for readiness using [zf\\_muxer\\_wait\(\)](#).

This is analogous to the non-blocking connection model for POSIX sockets.

#### Returns

- 0 Success.
- EAFNOSUPPORT raddr is not an AF\_INET address
- EADDRINUSE Address already in use.
- EBUSY Out of hardware resources.
- EFAULT Invalid pointer.
- EHOSTUNREACH No route to remote host.
- ENOMEM Out of memory.

#### See Also

[zft\\_addr\\_bind\(\)](#)

#### 11.10.2.4 int zft\_error ( struct zft \* ts )

Find out the error type happened on the TCP zocket.

#### Parameters

<i>ts</i>	TCP zocket.
-----------	-------------

#### Return values

<i>errno</i>	value, similar to SO_ERROR value for sockets.
--------------	---

<i>Error</i>	values are designed to be similar to Linux SO_ERROR:
<i>ECONNREFUSED</i>	The connection attempt was refused by server.
<i>ECONNRESET</i>	The connection was reset by the peer after it was established.
<i>ETIMEDOUT</i>	The connection was timed out, probably because of network failure.
<i>EPIPE</i>	The connection was closed gracefully by the peer (i.e. we've received all the data they've sent to us), but the peer refused to receive the data we've tried to send.

#### 11.10.2.5 int zft\_free ( struct zft \* ts )

Release resources associated with a TCP zocket.

##### Parameters

<i>ts</i>	TCP zocket.
-----------	-------------

This call shuts down the zocket if necessary. The application must not use `ts` after this call.

##### Returns

0 on success. Negative values are reserved for future use as error codes, but are not returned at present.

#### 11.10.2.6 unsigned zft\_get\_header\_size ( struct zft \* ts )

Return protocol header size for this connection.

##### Parameters

<i>ts</i>	The TCP zocket to query the header size for.
-----------	--

##### Returns

Protocol header size in bytes.

This function returns the total size of all protocol headers in bytes. An outgoing packet's size will be exactly the sum of this value and the number of payload data bytes it contains.

This function cannot fail.

#### 11.10.2.7 int zft\_get\_mss ( struct zft \* ts )

Retrieve the maximum segment size (MSS) for a TCP connection.

##### Parameters

<i>ts</i>	The TCP zocket to query.
-----------	--------------------------

##### Returns

- >= 0 The value of the MSS in bytes.
- ENOTCONN Zocket is not in a valid TCP state for sending.

#### 11.10.2.8 void zft\_getname ( struct zft \* ts, struct sockaddr \* laddr\_out, socklen\_t \* laddrlen, struct sockaddr \* raddr\_out, socklen\_t \* raddrlen )

Retrieve the local address of the zocket.



#### Parameters

<i>ts</i>	TCP zocket.
<i>laddr_out</i>	Return the local address of the zocket.
<i>laddrlen</i>	The length of the structure pointed to by <i>laddr_out</i>
<i>raddr_out</i>	Return the remote address of the zocket.
<i>raddrlen</i>	The length of the structure pointed to by <i>raddr_out</i>

This function returns local and/or remote IP address and TCP port of the given connection. Caller may pass NULL pointer for local or remote address if he is interested in the other address only.

If the supplied address structures are too small the result will be truncated and *addrlen* updated to a length greater than that supplied.

#### 11.10.2.9 int zft\_handle\_free ( struct zft\_handle \* *handle* )

Release a handle to a TCP zocket.

#### Parameters

<i>handle</i>	Handle to be released.
---------------	------------------------

This function releases resources associated with a [zft\\_handle](#).

#### Returns

0 Success.

#### 11.10.2.10 void zft\_handle\_getname ( struct zft\_handle \* *ts*, struct sockaddr \* *laddr\_out*, socklen\_t \* *laddrlen* )

Retrieve the local address to which a [zft\\_handle](#) is bound.

#### Parameters

<i>ts</i>	TCP zocket handle
<i>laddr_out</i>	Return the local address of the zocket
<i>laddrlen</i>	On entry, the size in bytes of the structure pointed to by <i>laddr_out</i> . Set on return to be the size in bytes of the result.

This function returns the local IP address and TCP port of the given listener. The behavior is undefined if the zocket is not bound.

If the supplied structure is too small the result will be truncated and *laddrlen* updated to a length greater than that supplied.

#### 11.10.2.11 int zft\_recv ( struct zft \* *ts*, const struct iovec \* *iov*, int *iovcnt*, int *flags* )

Copy-based receive.

#### Parameters

<i>ts</i>	TCP zocket
<i>iov</i>	Array with vectors pointing to buffers to fill with packet payloads.
<i>iovcnt</i>	The maximum number of buffers supplied (i.e. size of <i>iov</i> )
<i>flags</i>	None yet, must be zero.

#### Returns

>0 Number of bytes successfully received  
 0 End of File - other end has closed the connection  
 -EAGAIN No data available to read.  
 Other error codes are as for [zft\\_zc\\_rcv\\_done\(\)](#).

Copies received data on a zocket into buffers provided by the caller. The number of bytes received is returned. The caller's buffers will be filled as far as possible, and so a positive return value of less than the total space available in *iov* implies that no further data is available.

If no data is available, there are two possibilities: either the connection is still open, in which case `-EAGAIN` is returned, or else the connection has been closed by the peer, in which case the function succeeds and returns zero.

#### 11.10.2.12 `ssize_t zft_send ( struct zft * ts, const struct iovec * iov, int iov_cnt, int flags )`

Send data specified in iovec array.

#### Parameters

<i>ts</i>	The TCP zocket to send on.
<i>iov</i>	The iovec of data to send.
<i>iov_cnt</i>	The length of <i>iov</i> .
<i>flags</i>	Flags. 0 or <code>MSG_MORE</code> .

Sends the supplied data or its part. `MSG_MORE` flag will prevent sending packet that is not filled up to MSS.

Provided buffers may be re-used on successful return from this function.

#### Returns

Number of bytes sent on success.  
 -EINVAL Incorrect arguments supplied.  
 -ENOTCONN Zocket is not in a valid TCP state for sending.  
 -EAGAIN Not enough space (either bytes or buffers) in the send queue.  
 -ENOMEM Not enough packet buffers available.

#### Note

This function does not support sending zero-length data, and does not raise an error if you do so. Every iovec in the *iov* array must have length greater than 0, and *iov\_cnt* must also be greater than 0.

The *flags* argument must be set to 0 or `MSG_MORE`.

Notes on current implementation:

1. Currently, this function will return `-ENOMEM` without sending any data if it is unable to send the entire message due to shortage of packet buffers. This behaviour might change in future releases.

2. In case of partial send, the data is queued with MSG\_MORE flag set, and so may not go out immediately. See below for details of how to flush a MSG\_MORE send.
3. MSG\_MORE flag prevents the last partially filled segment from being sent immediately. The only guaranteed way to flush such a segment is to follow MSG\_MORE send with normal send - otherwise the segment might never get sent at all or it may take undefined amount of time. Some non-guaranteed triggers that might induce flush of a MSG\_MORE segment:
  - further MSG\_MORE send causes the segment to become full,
  - preceding normal send left partially filled segment in sendqueue, or
  - during stack polling TCP state machine intends to send ACK in response to incoming data.

#### 11.10.2.13 ssize\_t zft\_send\_single ( struct zft \* *ts*, const void \* *buf*, size\_t *buflen*, int *flags* )

Send data given in single buffer.

##### Parameters

<i>ts</i>	The TCP socket to send on.
<i>buf</i>	The buffer of data to send.
<i>buflen</i>	The length of buffer.
<i>flags</i>	Flags. 0 or MSG_MORE.

Sends the supplied data or its part. MSG\_MORE flag will prevent sending packet that is not filled up to MSS.

Provided buffer may be re-used on successful return from this function.

##### Returns

- Number of bytes sent on success.
- EINVAL Incorrect arguments supplied.
- ENOTCONN Socket is not in a valid TCP state for sending.
- EAGAIN Not enough space (either bytes or buffers) in the send queue.
- ENOMEM Not enough packet buffers available.

##### Note

This function does not support sending zero-length data, and does not raise an error if you do so. The flags argument must be set to 0 or MSG\_MORE.

Notes on current implementation:

1. Currently, this function will return -ENOMEM without sending any data if it is unable to send the entire message due to shortage of packet buffers. This behaviour might change in future releases.
2. MSG\_MORE flag prevents the last partially filled segment from being sent immediately. The only guaranteed way to flush such a segment is to follow MSG\_MORE send with normal send - otherwise the segment might never get sent at all or it may take undefined amount of time. Some non-guaranteed triggers that might induce flush of a MSG\_MORE segment:
  - further MSG\_MORE send causes the segment to become full,
  - preceding normal send left partially filled segment in sendqueue, or
  - during stack polling TCP state machine intends to send ACK in response to incoming data.

#### 11.10.2.14 int zft\_send\_space ( struct zft \* *ts*, size\_t \* *space* )

Query available space in the send queue.

**Parameters**

<i>ts</i>	The TCP zocket to query the send queue for.
<i>space</i>	On successful return, the available space in bytes.

This function will return the current space available in the send queue for the given zocket. This can be used to avoid [zft\\_send\(\)](#) returning `-EAGAIN`.

**Returns**

- 0 Success.
- ENOTCONN Zocket is not in a valid TCP state for sending.

**Note**

Available send queue space is a function of the number of the number of bytes queued, the number of internal buffers in the queue, and the MSS. Making many small sends can therefore consume more space than a single large send, and force [zft\\_send\(\)](#) to compress the send queue to avoid returning `-EAGAIN`.

**11.10.2.15 int zft\_shutdown\_tx ( struct zft \* ts )**

Shut down outgoing TCP connection.

**Parameters**

<i>ts</i>	A connected TCP zocket.
-----------	-------------------------

This function closes the TCP connection, preventing further data transmission except for already-queued data. This function does not prevent the connection from receiving more data.

**Returns**

- 0 on success, or a negative error code. Error codes returned are similar to [zft\\_send\(\)](#) ones:
- ENOTCONN Inappropriate TCP state: not connected or already shut down.
- EAGAIN Not enough space (either bytes or buffers) in the send queue.
- ENOMEM Not enough packet buffers available.

**11.10.2.16 int zft\_state ( struct zft \* ts )**

Return the TCP state of a TCP zocket.

**Parameters**

<i>ts</i>	TCP zocket.
-----------	-------------

**Returns**

Standard TCP\_\* state constant (e.g. TCP\_ESTABLISHED).

**11.10.2.17 struct zf\_waitable\* zft\_to\_waitable ( struct zft \* ts )**

Returns a [zf\\_waitable](#) representing the given [zft](#).

#### Parameters

<i>ts</i>	The <a href="#">zft</a> to return as a <a href="#">zf_waitable</a>
-----------	--

#### Returns

The [zf\\_waitable](#)

This function is necessary to use TCP zockets with the multiplexer.

#### 11.10.2.18 void zft\_zc\_recv ( struct zft \* *ts*, struct zft\_msg \* *msg*, int *flags* )

Zero-copy read of available packets.

#### Parameters

<i>ts</i>	TCP zocket.
<i>msg</i>	Message structure.
<i>flags</i>	Reserved. Must be zero.

This function completes the supplied `msg` structure with details of received packet buffers. In case of EOF a zero-length buffer is appended at the end of data stream and to identify reason of stream termination check result of [zft\\_zc\\_recv\\_done\(\)](#) or of [zft\\_zc\\_recv\\_done\\_some\(\)](#).

The function will only fill fewer iovecs in `msg` than are provided in the case where no further data is available.

Buffers are 'locked' until [zft\\_zc\\_recv\\_done\(\)](#) or [zft\\_zc\\_recv\\_done\\_some\(\)](#) is performed. The caller must not modify the contents of `msg` until after it has been passed to [zft\\_zc\\_recv\\_done\(\)](#) or to [zft\\_zc\\_recv\\_done\\_some\(\)](#).

#### 11.10.2.19 int zft\_zc\_recv\_done ( struct zft \* *ts*, struct zft\_msg \* *msg* )

Concludes pending `zc_recv` operation as done.

#### Parameters

<i>ts</i>	TCP zocket
<i>msg</i>	Message

#### Returns

- >= 1 Connection still receiving.
- 0 EOF.
- ECONNREFUSED Connection refused. This is possible as [zft\\_connect\(\)](#) is non-blocking.
- ECONNRESET Connection reset by peer.
- EPIPE Peer closed connection gracefully, but refused to receive some data sent on this zocket.
- ETIMEDOUT Connection timed out.

This function (or [zft\\_zc\\_recv\\_done\\_some\(\)](#)) must be called after each successful [zft\\_zc\\_recv\(\)](#) operation that returned at least one packet. It must not be called otherwise (in particular, when [zft\\_zc\\_recv\(\)](#) returned no packets). The function releases resources and enables subsequent calls to [zft\\_zc\\_recv\(\)](#) or [zft\\_recv\(\)](#). `msg` must be passed unmodified from the call to [zft\\_zc\\_recv\(\)](#).

#### 11.10.2.20 int zft\_zc\_recv\_done\_some ( struct zft \* *ts*, struct zft\_msg \* *msg*, size\_t *len* )

Concludes pending `zc_recv` operation as done acknowledging all or some of the data to have been read.

#### Parameters

<i>ts</i>	TCP zocket.
<i>msg</i>	Message.
<i>len</i>	Total number of bytes read by the client.

#### Returns

As for [zft\\_zc\\_recv\\_done\(\)](#).

Can be called after each successful [zft\\_zc\\_recv\(\)](#) operation as an alternative to [zft\\_zc\\_recv\\_done\(\)](#) or in cases where not all payload have been consumed. The restrictions on when it may be called are the same as for [zft\\_zc\\_recv\\_done\(\)](#). The function releases resources and enables subsequent calls to [zft\\_zc\\_recv\(\)](#) or [zft\\_recv\(\)](#). [zft\\_zc\\_recv\(\)](#) or [zft\\_recv\(\)](#) functions will return data indicated as non-read when they are called next time. *msg* must be passed unmodified from the call to [zft\\_zc\\_recv\(\)](#). *len* must not be greater than total payload returned by [zft\\_zc\\_recv\(\)](#).

#### 11.10.2.21 int zftl\_accept ( struct zftl \* *tl*, struct zft \*\* *ts\_out* )

Accept incoming TCP connection.

#### Parameters

<i>tl</i>	The listening zocket from which to accept the connection.
<i>ts_out</i>	On successful return filled with pointer to a TCP zocket for the new connection.

#### Returns

0 Success.  
-EAGAIN No incoming connections available.

#### 11.10.2.22 int zftl\_free ( struct zftl \* *ts* )

Release resources associated with a TCP listening zocket.

#### Parameters

<i>ts</i>	A TCP listening zocket.
-----------	-------------------------

This call shuts down the listening zocket, closing any connections waiting on the zocket that have not yet been accepted. The application must not use *ts* after this call.

#### Returns

0 Success.

#### 11.10.2.23 void zftl\_getname ( struct zftl \* *ts*, struct sockaddr \* *laddr\_out*, socklen\_t \* *laddrlen* )

Retrieve the local address of the zocket.

#### Parameters

<i>ts</i>	TCP zocket.
<i>laddr_out</i>	Set on return to the local address of the zocket.
<i>laddrlen</i>	On entry, the size in bytes of the structure pointed to by <i>laddr_out</i> . Set on return to be the size in bytes of the result.

This function returns the local IP address and TCP port of the listening zocket. If the supplied structure is too small the result will be truncated and *laddrlen* updated to a length greater than that supplied.

#### 11.10.2.24 `int zftl_listen ( struct zf_stack * st, const struct sockaddr * laddr, socklen_t laddrlen, const struct zf_attr * attr, struct zftl ** tl_out )`

Allocate TCP listening zocket.

#### Parameters

<i>st</i>	Initialized <a href="#">zf_stack</a> in which to created the listener.
<i>laddr</i>	Local address on which to listen. Must be non-null, and must be a single local address (not INADDR_ANY).
<i>laddrlen</i>	The size in bytes of the structure pointed to by <i>laddr</i>
<i>attr</i>	Attributes to apply to this zocket. Note that not all attributes are relevant; only those which apply to objects of type "zf_socket" are applicable here. Refer to the attribute documentation in <a href="#">Attributes</a> for details.
<i>tl_out</i>	On successful return filled with pointer to created TCP listening zocket.

#### Returns

- 0 Success.
- EFAULT Invalid *laddr* pointer.
- EADDRINUSE Local address already in use.
- EADDRNOTAVAIL *laddr* is not a local address.
- EAFNOSUPPORT *laddr* is not an AF\_INET address.
- EINVAL Zocket is already listening, or invalid *addr* length.
- ENOBUFFS No zockets of this type available.
- ENOMEM Out of memory.
- EOPNOTSUPP *laddr* is INADDR\_ANY.

#### 11.10.2.25 `struct zf_waitable* zftl_to_waitable ( struct zftl * tl )`

Returns a [zf\\_waitable](#) representing *tl*.

#### Parameters

<i>tl</i>	The <a href="#">zftl</a> to return as a <a href="#">zf_waitable</a>
-----------	---

#### Returns

The [zf\\_waitable](#)

This function is necessary to use TCP listening zockets with the multiplexer.

## 11.11 zf\_udp.h File Reference

TCPDirect UDP API.

```
#include <netinet/in.h>
#include <netinet/ip.h>
#include <netinet/udp.h>
#include <assert.h>
```

### Data Structures

- struct [zfur](#)  
*Opaque structure describing a UDP-receive zocket.*
- struct [zfur\\_msg](#)  
*UDP zero-copy RX message structure.*
- struct [zfut](#)  
*Opaque structure describing a UDP-transmit zocket.*

### Macros

- #define [ZFUT\\_FLAG\\_DONT\\_FRAGMENT](#) IP\_DF /\* 0x2000\*/  
*Flags for [zfut\\_send\(\)](#)*

### Functions

- int [zfur\\_alloc](#) (struct [zfur](#) \*\*us\_out, struct [zf\\_stack](#) \*st, const struct [zf\\_attr](#) \*attr)  
*Creates UDP-receive zocket.*
- int [zfur\\_free](#) (struct [zfur](#) \*us)  
*Release UDP-receive zocket previously created with [zfur\\_alloc\(\)](#).*
- int [zfur\\_addr\\_bind](#) (struct [zfur](#) \*us, struct sockaddr \*laddr, socklen\_t laddrlen, const struct sockaddr \*raddr, socklen\_t raddrlen, int flags)  
*Configures UDP-receive zocket to receive on a specified address.*
- int [zfur\\_addr\\_unbind](#) (struct [zfur](#) \*us, const struct sockaddr \*laddr, socklen\_t laddrlen, const struct sockaddr \*raddr, socklen\_t raddrlen, int flags)  
*Unbind UDP-receive zocket from address.*
- void [zfur\\_zc\\_rcv](#) (struct [zfur](#) \*us, struct [zfur\\_msg](#) \*msg, int flags)  
*Zero-copy read of single datagram.*
- void [zfur\\_zc\\_rcv\\_done](#) (struct [zfur](#) \*us, struct [zfur\\_msg](#) \*msg)  
*Concludes pending zero-copy receive operation as done.*
- int [zfur\\_pkt\\_get\\_header](#) (struct [zfur](#) \*us, const struct [zfur\\_msg](#) \*msg, const struct iphdr \*\*iphdr, const struct udphdr \*\*udphdr, int pktind)  
*Retrieves remote address from the header of a received packet.*
- struct [zf\\_waitable](#) \* [zfur\\_to\\_waitable](#) (struct [zfur](#) \*us)  
*Returns a [zf\\_waitable](#) representing the given [zfur](#).*
- int [zfut\\_alloc](#) (struct [zfut](#) \*\*us\_out, struct [zf\\_stack](#) \*st, const struct sockaddr \*laddr, socklen\_t laddrlen, const struct sockaddr \*raddr, socklen\_t raddrlen, int flags, const struct [zf\\_attr](#) \*attr)



- Allocate a UDP-transmit zocket.*
- int [zfut\\_free](#) (struct [zfut](#) \*us)
- Free UDP-transmit zocket.*
- int [zfut\\_get\\_mss](#) (struct [zfut](#) \*us)
- Get the maximum segment size which can be transmitted.*
- int [zfut\\_send\\_single](#) (struct [zfut](#) \*us, const void \*buf, size\_t buflen)
- Copy-based send of single non-fragmented UDP packet.*
- int [zfut\\_send](#) (struct [zfut](#) \*us, const struct iovec \*iov, int iov\_cnt, int flags)
- Copy-based send of single UDP packet (possibly fragmented).*
- struct [zf\\_waitable](#) \* [zfut\\_to\\_waitable](#) (struct [zfut](#) \*us)
- Returns a [zf\\_waitable](#) representing the given [zfut](#).*
- unsigned [zfut\\_get\\_header\\_size](#) (struct [zfut](#) \*us)
- Return protocol header size for this zocket.*

### 11.11.1 Detailed Description

TCPDirect UDP API.

Definition in file [zf\\_udp.h](#).

### 11.11.2 Function Documentation

**11.11.2.1** int [zfur\\_addr\\_bind](#) ( struct [zfur](#) \* *us*, struct [sockaddr](#) \* *laddr*, [socklen\\_t](#) *laddrlen*, const struct [sockaddr](#) \* *raddr*, [socklen\\_t](#) *raddrlen*, int *flags* )

Configures UDP-receive zocket to receive on a specified address.

#### Parameters

<i>us</i>	The zocket to bind
<i>laddr</i>	Local address. Cannot be NULL or INADDR_ANY, but the port may be zero, in which case an ephemeral port is allocated.
<i>laddrlen</i>	Length of the structure pointed to by <i>laddr</i> .
<i>raddr</i>	Remote address. If NULL, traffic will be accepted from all remote addresses.
<i>raddrlen</i>	Length of the structure pointed to by <i>raddr</i> .
<i>flags</i>	Flags. Must be zero.

#### Returns

- 0 Success.
- EADDRINUSE Address already in use.
- EAFNOSUPPORT *laddr* and/or *raddr* are not AF\_INET addresses
- EBUSY Out of hardware resources.
- EINVAL Invalid address length supplied.
- EFAULT Invalid address supplied.
- ENOMEM Out of memory.

The port number in *laddr* is updated if it was set to 0 by the caller.

If the specified local address is multicast then this has the effect of joining the multicast group as well as setting the filter. The group membership will persist until either the address is unbound (see [zfur\\_addr\\_unbind\(\)](#)), or the zocket is closed.

**11.11.2.2** `int zfur_addr_unbind ( struct zfur * us, const struct sockaddr * laddr, socklen_t laddrlen,  
const struct sockaddr * raddr, socklen_t raddrlen, int flags )`

Unbind UDP-receive zocket from address.

#### Parameters

<i>us</i>	The zocket to unbind.
<i>laddr</i>	Local address. Can be NULL to match any local address.
<i>laddrlen</i>	Length of the structure pointed to by <i>laddr</i> .
<i>raddr</i>	Remote address. Can be NULL to match any remote address.
<i>raddrlen</i>	Length of the structure pointed to by <i>raddr</i> .
<i>flags</i>	Flags. Must be zero.

#### Returns

- 0 Success.
- EINVAL The zocket is not bound to the specified address.

The addresses specified must match those used in [zfur\\_addr\\_bind\(\)](#).

#### 11.11.2.3 int zfur\_alloc ( struct zfur \*\* *us\_out*, struct zf\_stack \* *st*, const struct zf\_attr \* *attr* )

Creates UDP-receive zocket.

#### Parameters

<i>us_out</i>	Pointer to receive new UDP-receive zocket's address.
<i>st</i>	Initialized <a href="#">zf_stack</a> in which to create the zocket.
<i>attr</i>	Attributes to apply to this zocket. Note that not all attributes are relevant; only those which apply to objects of type "zf_socket" are applicable here. Refer to the attribute documentation in <a href="#">Attributes</a> for details.

#### Returns

- 0 Success.
- ENOBUFFS No zockets of this type available.

Associates UDP-receive zocket with semi-wild or full hardware filter. Creates software filter and initializes receive queue. The zocket becomes ready to receive packets after this call.

#### 11.11.2.4 int zfur\_free ( struct zfur \* *us* )

Release UDP-receive zocket previously created with [zfur\\_alloc\(\)](#).

#### Parameters

<i>us</i>	The UDP zocket to release.
-----------	----------------------------

#### Returns

- 0 on success. Negative values are reserved for future use as error codes, but are not returned at present.

#### 11.11.2.5 int zfur\_pkt\_get\_header ( struct zfur \* *us*, const struct zfur\_msg \* *msg*, const struct iphdr \*\* *iphdr*, const struct udphdr \*\* *udphdr*, int *pktind* )

Retrieves remote address from the header of a received packet.

#### Parameters

<i>us</i>	UDP zocket.
<i>msg</i>	Message.
<i>iphdr</i>	Location to receive IP header.
<i>udphdr</i>	Location to receive UDP header.
<i>pktind</i>	Index of packet within <code>msg-&gt;iov</code> .

This is useful for zockets that can receive from many remote addresses, i.e. those for which `zfur_addr_bind()` was called with `raddr == NULL`.

#### 11.11.2.6 `struct zf_waitable* zfur_to_waitable ( struct zfur * us )`

Returns a `zf_waitable` representing the given `zfur`.

#### Parameters

<i>us</i>	The <code>zfur</code> to return as a <code>zf_waitable</code>
-----------	---

#### Returns

The `zf_waitable`

This is necessary for use with the multiplexer.

#### 11.11.2.7 `void zfur_zc_rcv ( struct zfur * us, struct zfur_msg * msg, int flags )`

Zero-copy read of single datagram.

#### Parameters

<i>us</i>	UDP zocket.
<i>msg</i>	Message structure.
<i>flags</i>	Must be zero.

This function completes the supplied `msg` structure with details of a received UDP datagram.

The function may not fill all the supplied iovecs in `msg` even in the case where further data is available, but you can discover if there is more data available using the `dgrams_left` field in `zfur_msg` after making this call.

TCPDirect does not yet support fragmented datagrams, but in the future such datagrams will be represented in the `msg` iovec as a scatter-gather array of packet buffers. If the iovec is not long enough it may return a partial datagram.

Buffers are 'locked' until `zfur_zc_rcv_done()` is performed. The caller must not modify the contents of `msg` until after it has been passed to `zfur_zc_rcv_done()`.

#### 11.11.2.8 `void zfur_zc_rcv_done ( struct zfur * us, struct zfur_msg * msg )`

Concludes pending zero-copy receive operation as done.

#### Parameters

<i>us</i>	UDP zocket.
<i>msg</i>	Message.

Must be called after each successful [zfur\\_zc\\_recv\(\)](#) operation that returns at least one packet. It must not be called otherwise (in particular, when [zfur\\_zc\\_recv\(\)](#) returned no packets). The function releases resources and enables subsequent calls to [zfur\\_zc\\_recv\(\)](#). *msg* must be passed unmodified from the call to [zfur\\_zc\\_recv\(\)](#).

**11.11.2.9** `int zfut_alloc ( struct zfut ** us_out, struct zf_stack * st, const struct sockaddr * laddr, socklen_t laddrlen, const struct sockaddr * raddr, socklen_t raddrlen, int flags, const struct zf_attr * attr )`

Allocate a UDP-transmit zocket.

#### Parameters

<i>us_out</i>	On success contains pointer to newly created UDP transmit zocket
<i>st</i>	Stack in which to create zocket
<i>laddr</i>	Local address. If INADDR_ANY is specified, the local address will be selected according to the route to <i>raddr</i> , but the port must be non-zero.
<i>laddrlen</i>	Length of the structure pointed to by <i>laddr</i> .
<i>raddr</i>	Remote address.
<i>raddrlen</i>	Length of the structure pointed to by <i>raddr</i> .
<i>flags</i>	Must be zero.
<i>attr</i>	Attributes to apply to the zocket. Note that not all attributes are relevant; only those which apply to objects of type "zf_socket" are applicable here. Refer to the attribute documentation in <a href="#">Attributes</a> for details.

#### Returns

- 0 Success.
- EFAULT Invalid pointer.
- EHOSTUNREACH No route to remote host.
- EINVAL Invalid local or remote address, or address lengths.
- ENOBUFS No zockets of this type available.

#### Note

Once the zocket is created, neither the local address nor the remote address can be changed.

**11.11.2.10** `int zfut_free ( struct zfut * us )`

Free UDP-transmit zocket.

#### Parameters

<i>us</i>	UDP-transmit zocket to free.
-----------	------------------------------

#### Returns

0 on success. Negative values are reserved for future use as error codes, but are not returned at present.

**11.11.2.11** `unsigned zfut_get_header_size ( struct zfut * us )`

Return protocol header size for this zocket.

#### Parameters

<i>us</i>	The UDP-TX zocket to query the header size for.
-----------	---

#### Returns

Protocol header size in bytes.

This function returns the total size of all protocol headers in bytes. An outgoing packet's size will be exactly the sum of this value and the number of payload data bytes it contains.

This function cannot fail.

#### 11.11.2.12 int zfut\_get\_mss ( struct zfut \* *us* )

Get the maximum segment size which can be transmitted.

#### Returns

Maximum buflen parameter which can be passed to [zfut\\_send\\_single\(\)](#). This value is constant for a given zocket.

#### 11.11.2.13 int zfut\_send ( struct zfut \* *us*, const struct iovec \* *iov*, int *iov\_cnt*, int *flags* )

Copy-based send of single UDP packet (possibly fragmented).

#### Parameters

<i>us</i>	The UDP zocket to send on.
<i>iov</i>	The iovec of data to send.
<i>iov_cnt</i>	The length of iov.
<i>flags</i>	Flags.

#### Returns

Payload bytes sent (i.e. `buflen`) on success.  
-AGAIN Hardware queue full. Call [zf\\_reactor\\_perform\(\)](#) until it returns non-zero and try again.  
-MSGSIZE Message too large.  
-NOBUFS Out of packet buffers.

For a small packet in a plain buffer with the ZFUT\_FLAG\_DONT\_FRAGMENT flag set, this function just calls [zfut\\_send\\_single\(\)](#). Otherwise it handles IO vector and fragments a UDP packet into multiple IP fragments as needed.

If ZFUT\_FLAG\_DONT\_FRAGMENT flag is specified, then the datagram should fit to the MSS value (see [zfut\\_get\\_mss\(\)](#) above), and the DontFragment bit in the IP header will be set.

#### See Also

[zfut\\_send\\_single\(\)](#)

#### 11.11.2.14 int zfut\_send\_single ( struct zfut \* *us*, const void \* *buf*, size\_t *buflen* )

Copy-based send of single non-fragmented UDP packet.

#### Parameters

<i>us</i>	The UDP zocket to send on.
<i>buf</i>	A buffer of the data to send.
<i>buflen</i>	The length of the buffer, in bytes.

#### Returns

Payload bytes sent (i.e. *buflen*) on success.

-EAGAIN Hardware queue full. Call [zf\\_reactor\\_perform\(\)](#) until it returns non-zero and try again.

-ENOBUFS Out of packet buffers.

The function uses PIO when possible (i.e. for small datagrams), and always sets the DontFragment bit in the IP header. *buflen* must be no larger than the value returned by [zfut\\_get\\_mss\(\)](#).

#### See Also

[zfut\\_get\\_mss\(\)](#) [zfut\\_send\(\)](#)

#### 11.11.2.15 struct zf\_waitable\* zfut\_to\_waitable ( struct zfut \* *us* )

Returns a [zf\\_waitable](#) representing the given [zfut](#).

#### Parameters

<i>us</i>	The <a href="#">zfut</a> to return as a <a href="#">zf_waitable</a> .
-----------	---

#### Returns

The [zf\\_waitable](#).

This function is necessary to use UDP-transmit zockets with the multiplexer.





# Index

- attr.h, [59](#)
  - zf\_attr\_alloc, [60](#)
  - zf\_attr\_doc, [61](#)
  - zf\_attr\_dup, [61](#)
  - zf\_attr\_free, [61](#)
  - zf\_attr\_get\_int, [62](#)
  - zf\_attr\_get\_str, [62](#)
  - zf\_attr\_reset, [62](#)
  - zf\_attr\_set\_from\_fmt, [62](#)
  - zf\_attr\_set\_from\_str, [63](#)
  - zf\_attr\_set\_int, [63](#)
  - zf\_attr\_set\_str, [63](#)
- dgrams\_left
  - zfur\_msg, [56](#)
- EPOLLSTACKHUP
  - zf\_stack.h, [76](#)
- flags
  - zft\_msg, [54](#)
  - zfur\_msg, [56](#)
- iov
  - zft\_msg, [54](#)
  - zfur\_msg, [56](#)
- iovcnt
  - zft\_msg, [54](#)
  - zfur\_msg, [56](#)
- muxer.h, [64](#)
  - zf\_muxer\_add, [65](#)
  - zf\_muxer\_alloc, [65](#)
  - zf\_muxer\_del, [65](#)
  - zf\_muxer\_free, [66](#)
  - zf\_muxer\_mod, [66](#)
  - zf\_muxer\_wait, [66](#)
  - zf\_waitable\_event, [67](#)
  - zf\_waitable\_fd\_get, [67](#)
  - zf\_waitable\_fd\_prime, [68](#)
- pkts\_left
  - zft\_msg, [54](#)
- reserved
  - zft\_msg, [55](#)
  - zfur\_msg, [57](#)
- types.h, [68](#)
- x86.h, [69](#)
- zf.h, [69](#)
  - zf\_alternatives\_alloc
    - zf\_alts.h, [70](#)
  - zf\_alternatives\_cancel
    - zf\_alts.h, [70](#)
  - zf\_alternatives\_free\_space
    - zf\_alts.h, [70](#)
  - zf\_alternatives\_query\_overhead\_tcp
    - zf\_alts.h, [71](#)
  - zf\_alternatives\_release
    - zf\_alts.h, [71](#)
  - zf\_alternatives\_send
    - zf\_alts.h, [72](#)
  - zf\_alts.h, [69](#)
    - zf\_alternatives\_alloc, [70](#)
    - zf\_alternatives\_cancel, [70](#)
    - zf\_alternatives\_free\_space, [70](#)
    - zf\_alternatives\_query\_overhead\_tcp, [71](#)
    - zf\_alternatives\_release, [71](#)
    - zf\_alternatives\_send, [72](#)
    - zft\_alternatives\_queue, [72](#)
  - zf\_attr, [51](#)
  - zf\_attr\_alloc
    - attr.h, [60](#)
  - zf\_attr\_doc
    - attr.h, [61](#)
  - zf\_attr\_dup
    - attr.h, [61](#)
  - zf\_attr\_free
    - attr.h, [61](#)
  - zf\_attr\_get\_int
    - attr.h, [62](#)
  - zf\_attr\_get\_str
    - attr.h, [62](#)
  - zf\_attr\_reset
    - attr.h, [62](#)
  - zf\_attr\_set\_from\_fmt
    - attr.h, [62](#)
  - zf\_attr\_set\_from\_str
    - attr.h, [63](#)
  - zf\_attr\_set\_int
    - attr.h, [63](#)
  - zf\_attr\_set\_str
    - attr.h, [63](#)
  - zf\_deinit

- zf\_stack.h, 76
- zf\_init
  - zf\_stack.h, 76
- zf\_muxer\_add
  - muxer.h, 65
- zf\_muxer\_alloc
  - muxer.h, 65
- zf\_muxer\_del
  - muxer.h, 65
- zf\_muxer\_free
  - muxer.h, 66
- zf\_muxer\_mod
  - muxer.h, 66
- zf\_muxer\_set, 52
- zf\_muxer\_wait
  - muxer.h, 66
- zf\_platform.h, 73
- zf\_reactor.h, 73
  - zf\_reactor\_perform, 74
  - zf\_stack\_has\_pending\_work, 74
- zf\_reactor\_perform
  - zf\_reactor.h, 74
- zf\_stack, 52
- zf\_stack.h, 75
  - EPOLLSTACKHUP, 76
  - zf\_deinit, 76
  - zf\_init, 76
  - zf\_stack\_alloc, 76
  - zf\_stack\_free, 77
  - zf\_stack\_is\_quiescent, 77
  - zf\_stack\_to\_waitable, 77
- zf\_stack\_alloc
  - zf\_stack.h, 76
- zf\_stack\_free
  - zf\_stack.h, 77
- zf\_stack\_has\_pending\_work
  - zf\_reactor.h, 74
- zf\_stack\_is\_quiescent
  - zf\_stack.h, 77
- zf\_stack\_to\_waitable
  - zf\_stack.h, 77
- zf\_tcp.h, 78
  - zft\_addr\_bind, 79
  - zft\_alloc, 80
  - zft\_connect, 80
  - zft\_error, 81
  - zft\_free, 82
  - zft\_get\_header\_size, 82
  - zft\_get\_mss, 82
  - zft\_getname, 82
  - zft\_handle\_free, 83
  - zft\_handle\_getname, 83
  - zft\_rcv, 83
  - zft\_send, 84
  - zft\_send\_single, 85
  - zft\_send\_space, 85
  - zft\_shutdown\_tx, 86
  - zft\_state, 86
  - zft\_to\_waitable, 86
  - zft\_zc\_rcv, 87
  - zft\_zc\_rcv\_done, 87
  - zft\_zc\_rcv\_done\_some, 87
  - zftl\_accept, 88
  - zftl\_free, 88
  - zftl\_getname, 88
  - zftl\_listen, 89
  - zftl\_to\_waitable, 89
- zf\_udp.h, 90
  - zfur\_addr\_bind, 91
  - zfur\_addr\_unbind, 91
  - zfur\_alloc, 93
  - zfur\_free, 93
  - zfur\_pkt\_get\_header, 93
  - zfur\_to\_waitable, 94
  - zfur\_zc\_rcv, 94
  - zfur\_zc\_rcv\_done, 94
  - zfut\_alloc, 95
  - zfut\_free, 95
  - zfut\_get\_header\_size, 95
  - zfut\_get\_mss, 96
  - zfut\_send, 96
  - zfut\_send\_single, 96
  - zfut\_to\_waitable, 97
- zf\_waitable, 52
- zf\_waitable\_event
  - muxer.h, 67
- zf\_waitable\_fd\_get
  - muxer.h, 67
- zf\_waitable\_fd\_prime
  - muxer.h, 68
- zft, 53
- zft\_addr\_bind
  - zf\_tcp.h, 79
- zft\_alloc
  - zf\_tcp.h, 80
- zft\_alternatives\_queue
  - zf\_alts.h, 72
- zft\_connect
  - zf\_tcp.h, 80
- zft\_error
  - zf\_tcp.h, 81
- zft\_free
  - zf\_tcp.h, 82
- zft\_get\_header\_size
  - zf\_tcp.h, 82
- zft\_get\_mss
  - zf\_tcp.h, 82
- zft\_getname
  - zf\_tcp.h, 82
- zft\_handle, 53
- zft\_handle\_free
  - zf\_tcp.h, 83

zft\_handle\_getname  
    zf\_tcp.h, 83

zft\_msg, 54  
    flags, 54  
    iov, 54  
    iovcnt, 54  
    pkts\_left, 54  
    reserved, 55

zft\_rcv  
    zf\_tcp.h, 83

zft\_send  
    zf\_tcp.h, 84

zft\_send\_single  
    zf\_tcp.h, 85

zft\_send\_space  
    zf\_tcp.h, 85

zft\_shutdown\_tx  
    zf\_tcp.h, 86

zft\_state  
    zf\_tcp.h, 86

zft\_to\_waitable  
    zf\_tcp.h, 86

zft\_zc\_rcv  
    zf\_tcp.h, 87

zft\_zc\_rcv\_done  
    zf\_tcp.h, 87

zft\_zc\_rcv\_done\_some  
    zf\_tcp.h, 87

zftl, 55

zftl\_accept  
    zf\_tcp.h, 88

zftl\_free  
    zf\_tcp.h, 88

zftl\_getname  
    zf\_tcp.h, 88

zftl\_listen  
    zf\_tcp.h, 89

zftl\_to\_waitable  
    zf\_tcp.h, 89

zfur, 55

zfur\_addr\_bind  
    zf\_udp.h, 91

zfur\_addr\_unbind  
    zf\_udp.h, 91

zfur\_alloc  
    zf\_udp.h, 93

zfur\_free  
    zf\_udp.h, 93

zfur\_msg, 56  
    dgrams\_left, 56  
    flags, 56  
    iov, 56  
    iovcnt, 56  
    reserved, 57

zfur\_pkt\_get\_header  
    zf\_udp.h, 93

zfur\_to\_waitable  
    zf\_udp.h, 94

zfur\_zc\_rcv  
    zf\_udp.h, 94

zfur\_zc\_rcv\_done  
    zf\_udp.h, 94

zfut, 57

zfut\_alloc  
    zf\_udp.h, 95

zfut\_free  
    zf\_udp.h, 95

zfut\_get\_header\_size  
    zf\_udp.h, 95

zfut\_get\_mss  
    zf\_udp.h, 96

zfut\_send  
    zf\_udp.h, 96

zfut\_send\_single  
    zf\_udp.h, 96

zfut\_to\_waitable  
    zf\_udp.h, 97