



TCPDirect User Guide

SF-116303-CD , Issue 9

2019/09/26 12:43:34

Xilinx, Inc

TCPDirect User Guide

The information disclosed to you hereunder (the “Materials”) is provided solely for the selection and use of Xilinx products. To the maximum extent permitted by applicable law: (1) Materials are made available “AS IS” and with all faults, Xilinx hereby DISCLAIMS ALL WARRANTIES AND CONDITIONS, EXPRESS, IMPLIED, OR STATUTORY, INCLUDING BUT NOT LIMITED TO WARRANTIES OF MERCHANTABILITY, NON-INFRINGEMENT, OR FITNESS FOR ANY PARTICULAR PURPOSE; and (2) Xilinx shall not be liable (whether in contract or tort, including negligence, or under any other theory of liability) for any loss or damage of any kind or nature related to, arising under, or in connection with, the Materials (including your use of the Materials), including for any direct, indirect, special, incidental, or consequential loss or damage (including loss of data, profits, goodwill, or any type of loss or damage suffered as a result of any action brought by a third party) even if such damage or loss was reasonably foreseeable or Xilinx had been advised of the possibility of the same. Xilinx assumes no obligation to correct any errors contained in the Materials or to notify you of updates to the Materials or to product specifications. You may not reproduce, modify, distribute, or publicly display the Materials without prior written consent. Certain products are subject to the terms and conditions of Xilinx’s limited warranty, please refer to Xilinx’s Terms of Sale which can be viewed at <https://www.xilinx.com/legal.htm#tos>; IP cores may be subject to warranty and support terms contained in a license issued to you by Xilinx. Xilinx products are not designed or intended to be fail-safe or for use in any application requiring fail-safe performance; you assume sole risk and liability for use of Xilinx products in such critical applications, please refer to Xilinx’s Terms of Sale which can be viewed at <https://www.xilinx.com/legal.htm#tos>.

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.

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

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

AUTOMOTIVE APPLICATIONS DISCLAIMER

AUTOMOTIVE PRODUCTS (IDENTIFIED AS “XA” IN THE PART NUMBER) ARE NOT WARRANTED FOR USE IN THE DEPLOYMENT OF AIRBAGS OR FOR USE IN APPLICATIONS THAT AFFECT CONTROL OF A VEHICLE (“SAFETY APPLICATION”) UNLESS THERE IS A SAFETY CONCEPT OR REDUNDANCY FEATURE CONSISTENT WITH THE ISO 26262 AUTOMOTIVE SAFETY STANDARD (“SAFETY DESIGN”). CUSTOMER SHALL, PRIOR TO USING OR DISTRIBUTING ANY SYSTEMS THAT INCORPORATE PRODUCTS, THOROUGHLY TEST SUCH SYSTEMS FOR SAFETY PURPOSES. USE OF PRODUCTS IN A SAFETY APPLICATION WITHOUT A SAFETY DESIGN IS FULLY AT THE RISK OF CUSTOMER, SUBJECT ONLY TO APPLICABLE LAWS AND REGULATIONS GOVERNING LIMITATIONS ON PRODUCT LIABILITY.

Copyright

© Copyright 2019 Xilinx, Inc. Xilinx, the Xilinx logo, Solarflare, Onload, TCPDirect, and other designated brands included herein are trademarks of Xilinx in the United States and other countries. All other trademarks are the property of their respective owners.

SF-116303-CD

Last Revised: September 2019

Issue 9

Contents

1	TCPDirect	1
1.1	Introduction	1
2	What's New	3
2.1	Bug fixes	3
3	Overview	5
3.1	Platforms	5
3.2	Components	5
3.3	Capabilities and Restrictions	6
3.3.1	Protocols	6
3.3.2	OS	6
3.3.3	Network Interface Configuration	7
3.4	How TCPDirect Increases Performance	7
3.4.1	Overhead	7
3.4.2	Latency	8
3.4.3	Bandwidth	8
3.4.4	Scalability	8
3.5	Requirements	8
3.5.1	Adapter	8
3.5.2	License	9
3.5.3	Onload	9
3.5.4	Huge Pages	9
3.5.5	PIO	9

4	Concepts	11
4.1	Stacks	11
4.2	Zockets	11
4.2.1	TCP zockets	11
4.2.2	UDP zockets	12
4.2.3	Waitables	12
4.3	Multiplexers	12
4.4	TX alternatives	12
4.5	Cut-through PIO	12
4.6	Delegated sends	13
5	Example Applications	15
5.1	zfdppingpong	15
5.1.1	Usage	15
5.2	zftcpingpong	16
5.2.1	Usage	16
5.3	zfaltpingpong	16
5.4	zfsink	16
5.4.1	Usage	16
5.5	zftcpmtpong	17
5.5.1	Usage	17
5.6	exchange	17
5.7	trader_tcpdirect_ds_efvi	17
5.7.1	Usage	17
5.8	Building the Example Applications	18

6	Using TCPDirect	19
6.1	Components	19
6.2	Compiling and Linking	19
6.2.1	Header files	19
6.2.2	Linking	19
6.2.3	Debugging	20
6.3	General	20
6.4	Using stacks	21
6.5	Using zockets	21
6.6	UDP receive	21
6.7	UDP send	22
6.8	TCP listening	22
6.9	TCP send and receive	23
6.10	Alternative Tx queues	24
6.11	Epoll – muxer.h	26
6.12	Stack polling	27
6.13	Cut-through PIO	28
6.13.1	Underrun, poisoning and fallback:	28
6.13.2	CTPIO diagnostics	28
6.14	Delegated sends	29
6.15	Timestamps	29
6.16	VLANs	30
6.17	Miscellaneous	30
6.18	Errors issued by newer C++ compilers	30
6.19	zf_stackdump	30
6.19.1	Usage	31
6.19.2	stackdump output: stack	31
6.19.3	stackdump output: UDP RX	32
6.19.4	stackdump output: UDP TX	32
6.19.5	stackdump output: TCP TX/RX	32

7	Worked Examples	35
7.1	UDP ping pong example	35
7.2	TCP ping pong example	36
8	Attributes	39
8.1	alt_buf_size Attribute Reference	40
8.2	alt_count Attribute Reference	40
8.3	arp_reply_timeout Attribute Reference	41
8.4	ctpio Attribute Reference	41
8.5	ctpio_mode Attribute Reference	42
8.6	interface Attribute Reference	43
8.7	log_file Attribute Reference	43
8.8	log_format Attribute Reference	44
8.9	log_level Attribute Reference	44
8.10	max_tcp_endpoints Attribute Reference	45
8.11	max_tcp_listen_endpoints Attribute Reference	46
8.12	max_tcp_syn_backlog Attribute Reference	46
8.13	max_udp_rx_endpoints Attribute Reference	47
8.14	max_udp_tx_endpoints Attribute Reference	47
8.15	n_bufs Attribute Reference	48
8.16	name Attribute Reference	48
8.17	pio Attribute Reference	49
8.18	reactor_spin_count Attribute Reference	50
8.19	rx_ring_max Attribute Reference	50
8.20	rx_ring_refill_batch_size Attribute Reference	51
8.21	rx_ring_refill_interval Attribute Reference	51
8.22	rx_timestamping Attribute Reference	52
8.23	tcp_alt_ack_rewind Attribute Reference	52
8.24	tcp_delayed_ack Attribute Reference	53
8.25	tcp_finwait_ms Attribute Reference	53
8.26	tcp_initial_cwnd Attribute Reference	54
8.27	tcp_retries Attribute Reference	54
8.28	tcp_syn_retries Attribute Reference	55
8.29	tcp_synack_retries Attribute Reference	55
8.30	tcp_timewait_ms Attribute Reference	56
8.31	tcp_wait_for_time_wait Attribute Reference	56
8.32	tx_ring_max Attribute Reference	57
8.33	tx_timestamping Attribute Reference	57

9	Data Structure Index	59
9.1	Data Structures	59
10	File Index	61
10.1	File List	61
11	Data Structure Documentation	63
11.1	zf_attr Struct Reference	63
11.1.1	Detailed Description	63
11.2	zf_ds Struct Reference	64
11.2.1	Detailed Description	64
11.2.2	Field Documentation	64
11.2.2.1	cong_wnd	64
11.2.2.2	delegated_wnd	65
11.2.2.3	headers	65
11.2.2.4	headers_len	65
11.2.2.5	headers_size	65
11.2.2.6	ip_len_offset	65
11.2.2.7	ip_tcp_hdr_len	66
11.2.2.8	mss	66
11.2.2.9	reserved	66
11.2.2.10	send_wnd	66
11.2.2.11	tcp_seq_offset	66
11.3	zf_muxer_set Struct Reference	67
11.3.1	Detailed Description	67
11.4	zf_pkt_report Struct Reference	67
11.4.1	Detailed Description	67
11.4.2	Field Documentation	67
11.4.2.1	bytes	68
11.4.2.2	flags	68

11.4.2.3	start	68
11.4.2.4	timestamp	68
11.5	zf_stack Struct Reference	68
11.5.1	Detailed Description	69
11.6	zf_waitable Struct Reference	69
11.6.1	Detailed Description	69
11.7	zft Struct Reference	69
11.7.1	Detailed Description	70
11.8	zft_handle Struct Reference	70
11.8.1	Detailed Description	70
11.9	zft_msg Struct Reference	70
11.9.1	Detailed Description	70
11.9.2	Field Documentation	71
11.9.2.1	flags	71
11.9.2.2	iov	71
11.9.2.3	iovcnt	71
11.9.2.4	pkts_left	71
11.9.2.5	reserved	72
11.10	zftl Struct Reference	72
11.10.1	Detailed Description	72
11.11	zfur Struct Reference	72
11.11.1	Detailed Description	72
11.12	zfur_msg Struct Reference	73
11.12.1	Detailed Description	73
11.12.2	Field Documentation	73
11.12.2.1	dgrams_left	73
11.12.2.2	flags	73
11.12.2.3	iov	74
11.12.2.4	iovcnt	74
11.12.2.5	reserved	74
11.13	zfut Struct Reference	74
11.13.1	Detailed Description	74

12 File Documentation	75
12.1 attr.h File Reference	75
12.1.1 Detailed Description	75
12.1.2 Function Documentation	76
12.1.2.1 zf_attr_alloc()	76
12.1.2.2 zf_attr_doc()	76
12.1.2.3 zf_attr_dup()	77
12.1.2.4 zf_attr_free()	77
12.1.2.5 zf_attr_get_int()	77
12.1.2.6 zf_attr_get_str()	78
12.1.2.7 zf_attr_reset()	79
12.1.2.8 zf_attr_set_from_fmt()	79
12.1.2.9 zf_attr_set_from_str()	80
12.1.2.10 zf_attr_set_int()	80
12.1.2.11 zf_attr_set_str()	80
12.2 muxer.h File Reference	81
12.2.1 Detailed Description	82
12.2.2 Function Documentation	82
12.2.2.1 zf_muxer_add()	82
12.2.2.2 zf_muxer_alloc()	83
12.2.2.3 zf_muxer_del()	83
12.2.2.4 zf_muxer_free()	83
12.2.2.5 zf_muxer_mod()	84
12.2.2.6 zf_muxer_wait()	84
12.2.2.7 zf_waitable_event()	85
12.2.2.8 zf_waitable_fd_get()	86
12.2.2.9 zf_waitable_fd_prime()	86
12.3 types.h File Reference	87
12.3.1 Detailed Description	87

12.3.2	Macro Definition Documentation	87
12.3.2.1	ZF_PKT_REPORT_CLOCK_SET	87
12.3.2.2	ZF_PKT_REPORT_DROPPED	87
12.3.2.3	ZF_PKT_REPORT_IN_SYNC	88
12.3.2.4	ZF_PKT_REPORT_NO_TIMESTAMP	88
12.3.2.5	ZF_PKT_REPORT_TCP_FIN	88
12.3.2.6	ZF_PKT_REPORT_TCP_RETRANS	88
12.3.2.7	ZF_PKT_REPORT_TCP_SYN	88
12.4	x86.h File Reference	89
12.4.1	Detailed Description	89
12.5	zf.h File Reference	89
12.5.1	Detailed Description	89
12.6	zf_alts.h File Reference	89
12.6.1	Detailed Description	90
12.6.2	Function Documentation	90
12.6.2.1	zf_alternatives_alloc()	90
12.6.2.2	zf_alternatives_cancel()	90
12.6.2.3	zf_alternatives_free_space()	91
12.6.2.4	zf_alternatives_query_overhead_tcp()	91
12.6.2.5	zf_alternatives_release()	92
12.6.2.6	zf_alternatives_send()	92
12.6.2.7	zft_alternatives_queue()	93
12.7	zf_ds.h File Reference	94
12.7.1	Detailed Description	95
12.7.2	Enumeration Type Documentation	95
12.7.2.1	zf_delegated_send_rc	95
12.7.3	Function Documentation	95
12.7.3.1	zf_delegated_send_cancel()	95
12.7.3.2	zf_delegated_send_complete()	96

12.7.3.3	zf_delegated_send_prepare()	96
12.7.3.4	zf_delegated_send_tcp_advance()	97
12.7.3.5	zf_delegated_send_tcp_update()	98
12.8	zf_platform.h File Reference	98
12.8.1	Detailed Description	98
12.9	zf_reactor.h File Reference	98
12.9.1	Detailed Description	99
12.9.2	Function Documentation	99
12.9.2.1	zf_reactor_perform()	99
12.9.2.2	zf_reactor_perform_attr()	99
12.9.2.3	zf_stack_has_pending_events()	100
12.9.2.4	zf_stack_has_pending_work()	101
12.10	zf_stack.h File Reference	101
12.10.1	Detailed Description	102
12.10.2	Macro Definition Documentation	102
12.10.2.1	EPOLLSTACKHUP	102
12.10.3	Function Documentation	102
12.10.3.1	zf_deinit()	102
12.10.3.2	zf_init()	102
12.10.3.3	zf_stack_alloc()	103
12.10.3.4	zf_stack_free()	103
12.10.3.5	zf_stack_is_quiescent()	104
12.10.3.6	zf_stack_to_waitable()	104
12.11	zf_tcp.h File Reference	105
12.11.1	Detailed Description	106
12.11.2	Function Documentation	106
12.11.2.1	zft_addr_bind()	106
12.11.2.2	zft_alloc()	107
12.11.2.3	zft_connect()	108

12.11.2.4	zft_error()	108
12.11.2.5	zft_free()	109
12.11.2.6	zft_get_header_size()	109
12.11.2.7	zft_get_mss()	110
12.11.2.8	zft_get_tx_timestamps()	110
12.11.2.9	zft_getname()	110
12.11.2.10	zft_handle_free()	111
12.11.2.11	zft_handle_getname()	111
12.11.2.12	zft_pkt_get_timestamp()	112
12.11.2.13	zft_recv()	112
12.11.2.14	zft_send()	113
12.11.2.15	zft_send_single()	114
12.11.2.16	zft_send_single_warm()	115
12.11.2.17	zft_send_space()	116
12.11.2.18	zft_shutdown_tx()	116
12.11.2.19	zft_state()	117
12.11.2.20	zft_to_waitable()	117
12.11.2.21	zft_zc_recv()	117
12.11.2.22	zft_zc_recv_done()	118
12.11.2.23	zft_zc_recv_done_some()	118
12.11.2.24	zftl_accept()	119
12.11.2.25	zftl_free()	119
12.11.2.26	zftl_getname()	120
12.11.2.27	zftl_listen()	120
12.11.2.28	zftl_to_waitable()	121
12.12	zfd_udp.h File Reference	121
12.12.1	Detailed Description	123
12.12.2	Function Documentation	123
12.12.2.1	zfd_addr_bind()	123

12.12.2.2 zfur_addr_unbind()	123
12.12.2.3 zfur_alloc()	124
12.12.2.4 zfur_free()	124
12.12.2.5 zfur_pkt_get_header()	125
12.12.2.6 zfur_pkt_get_timestamp()	125
12.12.2.7 zfur_to_waitable()	126
12.12.2.8 zfur_zc_recv()	126
12.12.2.9 zfur_zc_recv_done()	127
12.12.2.10 zfut_alloc()	127
12.12.2.11 zfut_free()	128
12.12.2.12 zfut_get_header_size()	128
12.12.2.13 zfut_get_mss()	129
12.12.2.14 zfut_get_tx_timestamps()	129
12.12.2.15 zfut_send()	129
12.12.2.16 zfut_send_single()	130
12.12.2.17 zfut_send_single_warm()	131
12.12.2.18 zfut_to_waitable()	131
Index	133

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-adapter 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

What's New

This chapter tells you what's new in this release of TCPDirect.

If you are not a previous user of TCPDirect, go to [Overview](#).

2.1 Bug fixes

Various bugs have been fixed. For details, see the [ChangeLog](#) file.

Chapter 3

Overview

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

3.1 Platforms

TCPDirect can be run on X2000-series and 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 the same Linux distributions as Onload. For details, refer to the *Onload User Guide* that was supplied with the Onload package.

3.2 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.

3.3 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.

3.3.1 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

3.3.2 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	Yes	Yes
Teaming	Yes	Yes
Send/receive via non-SFC interface	No	Yes

OS	TCPDirect	Onload
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()`.

3.3.3 Network Interface Configuration

TCPDirect can accelerate traffic over Solarflare network interfaces, and will respect the operating system's address and routing configuration. It can also accelerate traffic over bonds whose slaves are all Solarflare interfaces. The bonding implementations provided by both the `bonding` and `teaming` drivers are supported. This functionality is subject to the following restrictions:

- Only LACP and active-backup modes are supported.
- With LACP, the only supported hashing modes are `layer2`, `layer2+3`, and `layer3+4`.
- There is no support for the addition of slaves to bonds during the lifetime of the stack, and nor is there support for the removal of slaves from bonds. Failover between interfaces in a bond is supported, however.
- [TX alternatives](#) are not supported on stacks created on bonded interfaces.

3.4 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.

3.4.1 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.

3.4.2 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.

3.4.3 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.

3.4.4 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.

3.5 Requirements

3.5.1 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 sboot utility from the Solarflare Linux Utilities package (SF-107601-LS).

3.5.2 License

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

X2000 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
```

3.5.3 Onload

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

3.5.4 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).

3.5.5 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 4

Concepts

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

4.1 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 TCPDirect stack. This model avoids problems of lock-contention and cache-line bouncing and allows for scalability in a natural way.

4.2 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.

4.2.1 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.

4.2.2 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

4.2.3 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.

4.3 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.

4.4 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. It is also not supported on stacks running on bonded interfaces.

4.5 Cut-through PIO

CTPIO (Cut-through PIO) improves send latency by moving packets from the PCIe bus to network port with minimal latency. It can be used in three modes:

1. Cut-through: The frame is transmitted onto the network as it is streamed across the PCIe bus. This mode offers the best latency.
2. Store-and-forward: The frame is buffered on the adapter before transmitting onto the network.
3. Store-and-forward with poison disabled: As for (2), except that it is guaranteed that frames are never poisoned. When this mode is enabled on any VI, all VIs are placed into store-and-forward mode.

Due to differences in hardware architecture, CTPIO is not available on SFN8000-series or SFN7000-series adapters.

4.6 Delegated sends

Delegated sends allow the user to have TCPDirect handle the TCP state machine while performing critical sends through another mechanism (such as `ef_vi`) to achieve lower latency. This is analogous both in principle and in the arrangement of the API to the delegated sends feature of Onload.

Chapter 5

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
zfudppingpong	Measure round-trip latency with UDP.
zftcpingpong	Measure round-trip latency with TCP.
zfaltppingpong	Measure round-trip latency with TCP TX Alternatives.
zfsink	Receive stream of UDP datagrams and demonstrate muxer.
zftcpmtpong	Use of TCPDirect in multi-threaded applications.
exchange	Simplified electronic trading exchange.
trader_tcpdirect_ds_efvi	Simplified electronic trader.

5.1 zfudppingpong

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

5.1.1 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.

5.2 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*.

5.2.1 Usage

Server:

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

Client:

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

5.3 zfaltpingpong

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

Usage is as for [zftcpingpong](#).

5.4 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.

5.4.1 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.

5.5 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.

5.5.1 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).

5.6 exchange

The `exchange` application plays the role of a simplified electronic trading exchange. It is to be used in conjunction with the `trader_tcpdirect_ds_efvi` application.

5.7 trader_tcpdirect_ds_efvi

The `trader_tcpdirect_ds_efvi` application demonstrates various techniques to reduce latency. These techniques are often useful in electronic trading applications, and so this example takes the form of an extremely simplified electronic trading system.

The `exchange` application provides a simplified electronic trading exchange, and this application provides a simplified electronic trader.

A `trader_onload_ds_efvi` application demonstrates similar techniques for Onload.

For full details, see the `README` file in the `tests/trade_sim` directory.

5.7.1 Usage

For normal socket-based sends, run as follows:

Server: `onload -p latency-best ./exchange mcast-intf`

Client: `./trader_tcpdirect_ds_efvi mcast-intf server`

For "delegated" sends, run as follows:

Server: `onload -p latency-best ./exchange mcast-intf`

Client: `./trader_tcpdirect_ds_efvi -d mcast-intf server`

where:

- *mcast-intf* is the multicast interface on the server or client machine (e.g. `eth0`)
- *server* is the IP address of the server machine (e.g. `192.168.0.10`)

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

5.8 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 6

Using TCPDirect

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

6.1 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.

6.2 Compiling and Linking

6.2.1 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>
```

6.2.2 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`.

TCPDirect provides a stable API and ABI between the application and TCPDirect library. An application that works with an older version of TCPDirect should also work with newer versions of TCPDirect. Exceptions to this are noted in the Release Notes.

The TCPDirect user-space library and kernel drivers must always match. The best way to ensure this is to link to TCPDirect dynamically. Then when the Onload packages are upgraded to a newer version both the user-space and kernel components are upgraded together.

Applications that link to TCPDirect statically are effectively tied to a single version of Onload, and must be re-linked when Onload packages are upgraded.

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.

6.2.3 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:

- Additional logging can be selectively enabled at application start-up by using the `ZF_ATTR` environment variable to set the `log_level` attribute, as described in the [Attributes](#) chapter.
- As a convenience, the `-l` option to the `zf_debug` wrapper will set the `log_level` attribute to the specified value.
- Changing the `log_level` attribute while the application is running has no effect.

6.3 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 API is not re-entrant, and so cannot be called from signal handlers.

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.

6.4 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.

6.5 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.

6.6 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](#).

6.7 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](#).

6.8 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](#).

6.9 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](#).

6.10 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.

The number of stacks that can use TX alternatives simultaneously is limited, and varies by adapter and port mode. Typical limitations are as follows:

- SFN8522, 2x10Gb: at least 6 stacks can use TX alternatives
- SFN8542, 2x40Gb: at least 6 stacks can use TX alternatives
- SFN8542, 1x40Gb + 2x10Gb: at least 3 stacks can use TX alternatives
- SFN8542, 4x10Gb: TX alternatives are *not* available.

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>

Mixing `zft_recv()` with calls to `zft_alternatives` functions is OK. But receiving more than `tcp_alt_ack_rewind` bytes of data will trigger an automatic rebuild of the alternative, which might add a bit of latency to any other sends which are happening at the time.

Mixing `zft_send()` or `zft_send_single()` with an alternative is OK, except if `zft_send()` or `zft_send_single()` is called after a `zft_alternatives_queue()` call for the same socket. In this situation:

- any subsequent call to `zf_alternatives_send()` for the same alternative will fail (returning -EINVAL)
- the alternative must be cancelled before it can be re-used.

So this is OK:

```
zft_send(ts, <data>, 0);
zft_alternatives_queue(ts, q, <data>, flags);
zft_alternatives_send(stack, q);
```

and this is also OK:

```
zft_alternatives_queue(ts, q, <data>, flags);
zft_alternatives_send(stack, q);
zft_send(ts, <data>, 0);
```

but this is not OK:

```
zft_alternatives_queue(ts, q, <data>, flags);
zft_send(ts, <data>, 0);
zft_alternatives_send(stack, q); // Will return -EINVAL
```

To determine the maximum packet size you can queue on an alternative, use the `zft_alternatives_free_space()` function.

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

Note

These functions can all be found in [zft_alts.h](#).

6.11 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 `zft_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 `zft_waitable` that represents a given zocket:

```
struct zft_waitable* zft_to_waitable(struct zft* us);
struct zft_waitable* zftl_to_waitable(struct zftl* us);
struct zft_waitable* zfttl_to_waitable(struct zfttl* tl);
struct zft_waitable* zftt_to_waitable(struct zftt* ts);
```

This `zft_waitable` can then be added to a multiplexer set by calling `zft_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 zft_muxer_add(struct zft_muxer_set*,
                 struct zft_waitable* w,
                 const struct epoll_event* event);
```

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

```
int zft_muxer_wait(struct zft_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](#).

6.12 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](#).

6.13 Cut-through PIO

When using a suitable adapter CTPIO is enabled by default in "sf-np" mode. The operation of CTPIO can be controlled via the following TCPDirect attributes:

ctpio:

- 0: disable
- 1: enable (default)
- 2: enable, warn if not available
- 3: enable, fail if not available

ctpio_mode:

- ct: cut-through mode
- sf: store-and-forward mode
- sf-np: store-and-forward with no poisoning

6.13.1 Underrun, poisoning and fallback:

When using cut-through mode, if the frame is not streamed to the adapter at at least line rate, then the frame is likely to be poisoned. This is most likely to happen if the application thread is interrupted while writing the frame to the adapter. In the underrun case, the frame is terminated with an invalid FCS – this is referred to as "poisoning" – and so will be discarded by the link partner. Cut-through mode is currently expected to perform well only on 10G links.

CTPIO may be aborted for other reasons, including timeout while writing a frame, contention between threads using CTPIO at the same time, and the CPU writing data to the adapter in the wrong order.

In all of the above failure cases the adapter falls-back to sending via the traditional DMA mechanism, which incurs a latency penalty. So a valid copy of the packet is always transmitted, whether the CTPIO operation succeeds or not.

Normally only an underrun in cut-through mode will result in a poisoned frame being transmitted. In rare cases it is also possible for a poisoned frame to be emitted in store-and-forward mode. If it is necessary to strictly prevent poisoned packets from reaching the network, then poisoning can be disabled globally.

6.13.2 CTPIO diagnostics

The adapter maintains counters that show whether CTPIO is being used, and any reasons for CTPIO sends failing. These can be inspected as follows:

```
ethtool -S ethX | grep ctpio
```

Note that some of these counters are maintained on a per-adapter basis, whereas others are per network port.

6.14 Delegated sends

TCPDirect supports delegated sends.

Note

This is a new API that is offered as a technology preview: it is functionally complete but has had limited testing, and the API may change in future releases. Customers are invited to experiment with the feature and to direct feedback to support@solarflare.com.

Zockets are created through the TCPDirect API as normal. The user can then request that TCPDirect delegate sending to their application. Once the application has completed a send through (for example) `ef_vi`, it updates TCPDirect and TCPDirect will handle the TCP state machinery, retransmissions, and so on.

There is an example application at

```
<onload_install_dir>/src/tests/zf_apps/zfdelegated_client.c
```

which is intended to be used in conjunction with the existing `f_vi` server test app at

```
<onload_install_dir>/src/tests/ef_vi/efdelegated_server.c
```

This API is intended to be used by servers that make sporadic TCP sends on a zocket rather than large amounts of bi-directional traffic. It should be used carefully as there are small windows of time (while the send has been delegated to the application) where either the application or TCPDirect could be using out of date sequence or acknowledgement numbers. It has been designed such that this should be harmless, but may still have the potential to confuse other TCP implementations.

6.15 Timestamps

To enable RX timestamping, use the `rx_timestamping` attribute. To get the RX timestamps:

```
int zft_pkt_get_timestamp(struct zft* ts, const struct
    zft_msg* msg,
                        struct timespec* ts_out, int pktind,
                        unsigned* flags);

int zfur_pkt_get_timestamp(struct zfur* us, const struct
    zfur_msg* msg,
                        struct timespec *ts_out, int pktind,
                        unsigned* flags);
```

To enable TX timestamping, use the `tx_timestamping` attribute. To get the TX timestamps:

```
int zft_get_tx_timestamps(struct zft* ts, struct
    zf_pkt_report* reports,
                        int* count_in_out);

int zfut_get_tx_timestamps(struct zfut* us, struct
    zf_pkt_report* reports_out,
                        int* count_in_out);
```

These functions can all be found in `zf_tcp.h` and `zf_udp.h`.

6.16 VLANs

TCPDirect only supports a single interface per stack, and this restriction also applies to VLAN interfaces. Therefore, a separate TCPDirect stack is required for each separate VLAN tagged interface used (even if the underlying interface is the same).

The name of the VLAN tagged interface (e.g. enp4s0f0.100) must be specified in the interface attribute.

TCPDirect does not check VLAN tags when traffic is received. If traffic that arrives has a VLAN tag that does not match that of the specified interface, it will still be received. This includes untagged traffic to a tagged interface and vice-versa.

6.17 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](#).

6.18 Errors issued by newer C++ compilers

Applications using TCPDirect may fail to build with g++ version 6 and above with the message "error: flexible array member 'zft_msg::iov' not at end of 'struct my_msg'". To work around this issue, application code may be modified from:

```
struct my_msg {
    zft_msg msg;
    iovec iov[1];
};
```

to:

```
typedef struct {
    zft_msg msg;
    iovec iov[1];
} my_msg;
```

6.19 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.

6.19.1 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 discard_csum_bad=0 discard_mcast_mismatch=0
       discard_crc_bad=0 discard_trunc=0 discard_rights=0
       discard_ev_error=0 discard_other=0 discard_inner_csum_bad=0
       cplane_alien_ifindex=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
```

6.19.2 stackdump output: stack

Title	Parameter	Description
-	name	name of the stack.
pool	pkt_bufs_n	num of packet buffers allocated to the stack.
pool	free	num of free (available) packet buffers.
config	tcp_timewait_ticks	length of the TIME-WAIT timer in ticks.
config	tcp_finwait_ticks	length of the FIN-WAIT-2 timer in ticks.
config	ctpio_threshold	the cut-through threshold for CTPIO transmits.
config	tcp_initial_cwnd	size of TCP congestion window.
config	ms_per_tcp_tick	granularity of TCP timer in milliseconds.
alts	n_alts	total number of TX alternatives allocated to this stack.
stats	ring_refill_nomem	num times there were no free packet buffers to refill rx ring (increase buffers with n_bufs attr).
stats	discard_csum_bad	num of packets discarded due to IP, UDP or TCP checksum error.
stats	discard_mcast_mismatch	num of packets discarded due to hash mismatch in a multicast packet.
stats	discard_crc_bad	num of packets discarded due to Ethernet CRC error.
stats	discard_trunc	num of packets discarded due to a truncated frame
stats	discard_rights	num of packets discarded due to non-ownership rights to the packet.
stats	discard_ev_error	num of events discarded due to event queue overrun.
stats	discard_other	num of packets discarded due to other unspecified reason.
stats	discard_inner_csum_bad	num of packets discarded due to inner IP, UDP or TCP checksum error.
nic	vi	vi_i instance id of the VI being used by this stack.
nic	flags	ef_vi_flags of the VI being used by this stack.
nic	intf	name of the physical interface being used.
nic	index	ifindex of the physical interface being used.

nic	hw	type of the physical interface being used.
txq	pio_buf_size	size (bytes) of a PIO buffer.

6.19.3 stackdump output: UDP RX

Title	Parameter	Description
-	-	local interface.
filter	lcl	local_ip:port.
filter	rmt	remote_ip:port.
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.
udp rx	release_n	num zero-copy packet awaiting release.
udp rx	q_drops	num packets dropped from the zockets rx queue.

6.19.4 stackdump output: UDP TX

Title	Parameter	Description
-	-	local interface.
-	lcl	local_ip:port.
-	rmt	remote_ip:port.
path	dst	destination server.
path	src	source server.

6.19.5 stackdump output: TCP TX/RX

Title	Parameter	Description
path	dst	destination server.
path	src	source server.
rx	rx unread	num packets received, but still in zocket buffer rx queue
rx	rx begin	zf_rx_ring: oldest pkt buffer not yet read.
rx	rx process	zf_rx_ring: oldest pkt buffer not yet processed - so buffer not yet reaped.
rx	rx end	zf_rx_ring: index of the last pkt in the queue.
tcp	flags	zocket flags.
tcp	flags_ack_delay	ACK flags TF_ACK_DELAY 0x01 TF_ACK_NOW 0x02 TF_INTR 0x04 (in fast recovery) TF_ACK_NEXT 0x08.
tcp	error	error count.
tcp	parent	identifies the listing zocket from which a passive-open zocket was accepted.
tcp	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	snd_next	next sequence num to send.

Title	Parameter	Description
snd	lastack	last acknowledged sequence num.
snd	wnd	sender window.
snd	snd_wnd_max	maximum sender window advertised by the peer.
snd	snd_wl1	max sequence number advertised.
snd	snd_wl2	last sequence number acknowledged.
snd	snd_lbb	sequence num of next byte to be buffered.
snd	snd_right_edge	sequence num of TCP send window.
snd	delegated	bytes reserved by user of delegated send API
snd	send	num segments held in the send buffer.
snd	inflight	num segments sent, but not yet acknowledged.
snd	qbegin	TCP segment at sendq start.
snd	qmiddle	TCP segment at sendq middle.
snd	qend	TCP segment at sendq end.
snd	sndbuf	socket send buffer size (bytes).
snd	cwnd	size of congestion avoidance window.
snd	ssthresh	slow start threshold - num bytes that have to be sent before exiting slow start.
snd	mss_lim	max segment size limit set by peer, in bytes.
rcv	rcv_nxt	next expected sequence number.
rcv	rcv_ann_wnd	receiver window to announce.
rcv	rcv_ann_right_edge	announced right edge of window.
rcv	mss	max segment size, in bytes.
rtt	est	RTT estimate in ticks.
rtt	seq	sequence number used for RTT estimation.
rtt	sa	smoothed round trip time.
rtt	sv	round trip time variance estimate.
cong	nrtx	num of RTO retransmission attempts - reset to zero when a new ACK is received.
cong	dupacks	num duplicate acks received.
cong	persist_backoff	num of zero send win probes - sends probe pkt to keep connection alive.
timers	-	types of active timers (e.g. RTO, DACK, ZWIN, TIMEWAIT).
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.
stats	msg_more_send_delayed	num of times there was no send because of MSG_MORE flag.
stats	send_nomem	num of times there were no free packet buffers to perform a send.

Chapter 7

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.

7.1 UDP ping pong example

In the following example various boiler plate code has been omitted for clarity. For a full usable version of this example see `src/tests/zf_apps/zfudppingpong.c`, which includes how to use timestamping, and other details.

```
#define SIZE 12
#define ITERATIONS 1000000
int ping;

static void ping_pongs(struct zf_stack* stack, struct zfur* ur, struct
    zfut* ut)
{
    char send_buf[SIZE];
    int sends_left = ITERATIONS;
    int recvs_left = ITERATIONS;

    struct {
        /* The iovec used by zfur_msg must be immediately afterwards. */
        struct zfur_msg msg;
        struct iovec iov[2];
    } msg;
    const int max_iov = sizeof(msg.iov) / sizeof(msg.iov[0]);

    if( ping ) {
        ZF_TEST(zfut_send_single(ut, send_buf, SIZE) == SIZE);
        --sends_left;
    }

    do {
        /* Poll the stack until something happens. */
        while( zf_reactor_perform(stack) == 0 )
```

```

        ;
        msg.msg.iovcnt = max_iov;
        zfur_zc_recv(ur, &msg.msg, 0);
        if( msg.msg.iovcnt ) {
            if( sends_left ) {
                ZF_TEST(zfut_send_single(ut, send_buf, SIZE) == SIZE);
                --sends_left;
            }
            zfur_zc_recv_done(ur, &msg.msg);
            --recvs_left;
        }
    } while( recvs_left );
}

int main(int argc, char* argv[])
{
    if( argc != 3 )
        usage_err();

    if( ! strcmp(argv[0], "ping") )
        ping = true;
    else if( ! strcmp(argv[0], "pong") )
        ping = false;
    else
        usage_err();

    struct addrinfo *ai_local, *ai_remote;
    if( getaddrinfo_hostport(argv[1], NULL, &ai_local) != 0 ) {
        fprintf(stderr, "ERROR: failed to lookup address '%s'\n", argv[1]);
        exit(2);
    }
    if( getaddrinfo_hostport(argv[2], NULL, &ai_remote) != 0 ) {
        fprintf(stderr, "ERROR: failed to lookup address '%s'\n", argv[2]);
        exit(2);
    }

    /* Initialise the TCPDirect library and allocate a stack. */
    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));

    /* Allocate sockets and bind them to the given addresses. TCPDirect has
       separate objects for sending and receiving UDP datagrams.
    */
    struct zfur* ur;
    ZF_TRY(zfur_alloc(&ur, stack, attr));
    ZF_TRY(zfur_addr_bind(ur, ai_local->ai_addr, ai_local->ai_addrlen,
                        ai_remote->ai_addr, ai_remote->ai_addrlen, 0));

    struct zfut* ut;
    ZF_TRY(zfut_alloc(&ut, stack, ai_local->ai_addr, ai_local->ai_addrlen,
                    ai_remote->ai_addr, ai_remote->ai_addrlen, 0, attr));

    ping_pongs(stack, ur, ut);

    return 0;
}

```

7.2 TCP ping pong example

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

```

#define SIZE 12
#define ITERATIONS 1000000
int ping;

struct rx_msg {

```

```

/* The iovect used by zft_msg must be immediately afterwards */
struct zft_msg msg;
struct iovect iov[1];
};

static void ping_pongs(struct zf_stack* stack, struct zft* zock)
{
    char send_buf[SIZE];
    struct rx_msg msg;
    const int max_iov = sizeof(msg.iov) / sizeof(msg.iov[0]);
    int sends_left = ITERATIONS;
    int recvs_left = ITERATIONS;
    bool zock_has_rx_data = false;

    if( ping ) {
        ZF_TEST(zft_send_single(zock, send_buf, SIZE, 0) == SIZE);
        --sends_left;
    }

    do {
        size_t bytes_left = SIZE;
        do {
            if( ! zock_has_rx_data )
                /* Poll the stack until something happens. */
                while( zf_reactor_perform(stack) == 0 )
                    ;
            msg.msg.iovcnt = max_iov;
            zft_zc_recv(zock, &msg.msg, 0);
            if( msg.msg.iovcnt ) {
                /* NB. msg.iov[0].iov_len==0 indicates we're not going to get any
                 * more data (ie. the other end has shutdown or connection has
                 * died). We don't check for that here...instead it will be
                 * detected if zft_zc_recv_done()!=1.
                 */
                ZF_TEST(msg.iov[0].iov_len <= bytes_left);
                bytes_left -= msg.iov[0].iov_len;
                if( bytes_left == 0 )
                    /* Break out to do send before zft_zc_recv_done() to save a few
                     * nanoseconds.
                     */
                    break;
                ZF_TEST(zft_zc_recv_done(zock, &msg.msg) == 1);
            }
            zock_has_rx_data = msg.msg.pkts_left != 0;
        } while( bytes_left );

        if( sends_left ) {
            ZF_TEST(zft_send_single(zock, send_buf, SIZE, 0) == SIZE);
            --sends_left;
        }
        ZF_TEST(zft_zc_recv_done(zock, &msg.msg) == 1);
        --recvs_left;
    } while( recvs_left );
}

int main(int argc, char* argv[])
{
    if( argc != 2 )
        usage_err();

    if( ! strcmp(argv[0], "ping") )
        ping = true;
    else if( ! strcmp(argv[0], "pong") )
        ping = false;
    else
        usage_err();

    struct addrinfo* ai;
    if( getaddrinfo_hostport(argv[1], NULL, &ai) != 0 ) {
        fprintf(stderr, "ERROR: failed to lookup address '%s'\n", argv[1]);
        exit(2);
    }

    /* Initialise the TCPDirect library and allocate a stack. */
    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* zock;

if( ping ) {
    /* In 'ping' mode, connect to the specified remote address. */
    struct zft_handle* tcp_handle;
    ZF_TRY(zft_alloc(stack, attr, &tcp_handle));
    printf("Connecting to ponger\n");
    ZF_TRY(zft_connect(tcp_handle, ai->ai_addr, ai->ai_addrlen, &zock));
    /* The zft_connect() call is non-blocking, so the zocket is not yet
       connected. Wait until the connect completes or fails...
    */
    while( zft_state(zock) == TCP_SYN_SENT )
        zf_reactor_perform(stack);
    ZF_TEST( zft_state(zock) == TCP_ESTABLISHED );
}
else {
    /* In 'pong' mode, create a listening zocket and wait until we've
       accepted a connection.
    */
    struct zftl* listener;
    int rc;
    ZF_TRY(zftl_listen(stack, ai->ai_addr, ai->ai_addrlen, attr, &listener));
    printf("Waiting for incoming connection\n");
    do {
        while( zf_reactor_perform(stack) == 0 );
    } while( (rc = zftl_accept(listener, &zock) == -EAGAIN );
    ZF_TRY(rc);
    ZF_TRY(zftl_free(listener));
}
printf("Connection established\n");

ping_pongs(stack, zock);

/* Do a clean shutdown and free all resources. */
printf("Completed\n");
while( zft_shutdown_tx(zock) == -EAGAIN )
    zf_reactor_perform(stack);

while( ! zf_stack_is_quiescent(stack) )
    zf_reactor_perform(stack);

ZF_TRY(zft_free(zock));
ZF_TRY(zf_stack_free(stack));
ZF_TRY(zf_deinit());
return 0;
}

```

Chapter 8

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
<code>alt_buf_size</code>	Amount of NIC-side buffer space to allocate for use with TCP alternatives on this VI.
<code>alt_count</code>	Number of TCP alternatives to allocate on this VI.
<code>arp_reply_timeout</code>	Maximum time to wait for ARP replies, in microseconds (approx).
<code>ctpio</code>	Enable/Disable CTPIO.
<code>ctpio_mode</code>	Set the CTPIO mode to use.
<code>interface</code>	Use this interface name as <code>zf_stack</code> interface.
<code>log_file</code>	Use this file instead of <code>stderr</code> for log messages.
<code>log_format</code>	Bitmask to set the format of log messages.
<code>log_level</code>	Bitmask to enable different log message levels for each logging component.
<code>max_tcp_endpoints</code>	Sets the maximum number of TCP endpoints (i.e. struct <code>zft</code>).
<code>max_tcp_listen_endpoints</code>	Sets the maximum number of TCP listen endpoints (i.e. struct <code>zftl</code>).
<code>max_tcp_syn_backlog</code>	Sets the maximum number of half-open connections maintained in the stack.
<code>max_udp_rx_endpoints</code>	Sets the maximum number of UDP RX endpoints (i.e. struct <code>zfur</code>).
<code>max_udp_tx_endpoints</code>	Sets the maximum number of UDP TX endpoints (i.e. struct <code>zfut</code>).
<code>n_bufs</code>	Number of packet buffers to allocate for the stack.
<code>name</code>	The object name.
<code>pio</code>	Enable/Disable PIO buffers.
<code>reactor_spin_count</code>	Sets how many iterations of the event processing loop <code>zf_reactor_perform()</code> will make (in the absence of any events) before returning.
<code>rx_ring_max</code>	Set the size and maximum fill level of the RX descriptor ring, which provides buffering between the network adapter and software.
<code>rx_ring_refill_batch_size</code>	Sets the number of packet buffers rx ring is refilled with on each <code>zf_reactor_perform()</code> call.
<code>rx_ring_refill_interval</code>	Sets the frequency of rx buffer ring refilling during inner <code>zf_reactor_perform()</code> loop.
<code>rx_timestamping</code>	Add timestamps to received packets.
<code>tcp_alt_ack_rewind</code>	The maximum number of bytes by which outgoing ACKs will be allowed to go backwards when sending an alternative queue.
<code>tcp_delayed_ack</code>	Enable TCP delayed ACK ("on" by default).
<code>tcp_finwait_ms</code>	Length of TCP FIN-WAIT-2 timer in ms, 0 - disabled.

Attribute	Description
tcp_initial_cwnd	The initial congestion window for new TCP sockets.
tcp_retries	The maximum number of TCP retransmits if data is not acknowledged by the network peer in general case.
tcp_syn_retries	The maximum number of TCP SYN retransmits during zft_connect() .
tcp_synack_retries	The maximum number of TCP SYN-ACK retransmits before incoming connection is dropped.
tcp_timewait_ms	Length of TCP TIME-WAIT timer in ms.
tcp_wait_for_time_wait	Do not consider a stack to be quiescent if there are any TCP sockets in the TIME_WAIT state.
tx_ring_max	Set the size of the TX descriptor ring, which provides buffering between the software and the network adaptor.
tx_timestamping	Report timestamps for transmitted packets.

8.1 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

zfv_i.

8.2 alt_count Attribute Reference

Number of TCP alternatives to allocate on this VI.

Detailed Description

Not supported on stacks running on bonded network interfaces.

Type

Integer.

Default

0.

Relevant components

zf_vi.

8.3 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](#).

8.4 ctpio Attribute Reference

Enable/Disable CTPIO.

Detailed Description

- 0: don't use CTPIO,
- 1: CTPIO if available, no warning if not (default),
- 2: CTPIO if available, warn if not,
- 3: CTPIO else fail + error messages.

Note that warnings are disabled by default. If setting this attribute to 2, then bit 1 of the [log_level](#) attribute must also be set to enable warnings for the stack component.

Type

Integer.

Default

1.

Relevant components

zf_vi.

8.5 ctpio_mode Attribute Reference

Set the CTPIO mode to use.

Detailed Description

Set to:
'sf' for store-and-forward;
'ct' for cut-through;
'sf-np' to guarantee that poisoned frames are never emitted.

Type

String.

Default

sf-np.

Relevant components

[zf_vi](#).

8.6 interface Attribute Reference

Use this interface name as [zf_stack](#) interface.

Detailed Description

Type

String.

Default

none.

Relevant components

[zf_stack](#).

8.7 log_file Attribute Reference

Use this file instead of stderr for log messages.

Detailed Description

Type

String.

Default

(stderr).

Relevant components

[zf_stack](#).

8.8 log_format Attribute Reference

Bitmask to set the format of log messages.

Detailed Description

Combination of flags:

ZF_LF_STACK_NAME (0x1),

ZF_LF_FRC(0x2),

ZF_LF_TCP_TIME(0x4),

ZF_LF_PROCESS(0x8).

Type

Integer.

Default

stack name and tcp time.

Relevant components

[zf_stack](#).

8.9 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 (bits 4-7),

TCP-tx (bits 8-11),

TCP-connection (bits 12-15),

UDP-rx (bits 16-19),
UDP-tx (bits 20-23),
UDP-connection (bits 24-27),
muxer (bits 28-31),
pool (bits 32-35),
fast-path (bits 36-39),
timers (bits 40-43),
filters (bits 44-47),
cplane (bits 48-51).

E.g. 0xffff0 will enable all TCP related logging and disable all other logging.

Type

Bitmask.

Default

ERR-level on all components.

Relevant components

[zf_stack](#).

8.10 max_tcp_endpoints Attribute Reference

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

Detailed Description

This can be a value up to 64 (which is also the default).

Type

Integer.

Default

64.

Relevant components

[zf_stack](#).

8.11 max_tcp_listen_endpoints Attribute Reference

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

Detailed Description

This can be a value up to 64.

Type

Integer.

Default

16.

Relevant components

[zf_stack](#).

8.12 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](#).

8.13 max_udp_rx_endpoints Attribute Reference

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

Detailed Description

This can be a value up to 64 (which is also the default).

Type

Integer.

Default

64.

Relevant components

[zf_stack](#).

8.14 max_udp_tx_endpoints Attribute Reference

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

Detailed Description

This can be a value up to 64 (which is also the default).

Type

Integer.

Default

64.

Relevant components

[zf_stack](#).

8.15 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](#).

8.16 name Attribute Reference

The object name.

Detailed Description

The object name has a maximum length of 20 characters. Object names are visible in log messages, but have no other effect.

Type

String.

Default

(none).

Relevant components

[zf_stack](#), [zf_pool](#), [zf_vi](#).

8.17 pio Attribute Reference

Enable/Disable PIO buffers.

Detailed Description

0: don't use PIO,
1: PIO if available, no warning if not,
2: PIO if available, warn if not,
3: PIO else fail + error messages (default).

Note that warnings are disabled by default. If setting this attribute to 2, then bit 1 of the [log_level](#) attribute must also be set to enable warnings for the stack component.

Type

Integer.

Default

3.

Relevant components

[zf_vi](#).

8.18 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`.

8.19 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](#).

8.20 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](#).

8.21 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](#).

8.22 rx_timestamping Attribute Reference

Add timestamps to received packets.

Detailed Description

"off" by default.

Type

Integer.

Default

0.

Relevant components

[zf_vi](#).

8.23 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](#).

8.24 tcp_delayed_ack Attribute Reference

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

Detailed Description

Type

Integer.

Default

1.

Relevant components

[zf_stack](#).

8.25 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](#).

8.26 tcp_initial_cwnd Attribute Reference

The initial congestion window for new TCP sockets.

Detailed Description

Type

Integer.

Default

10 * MSS.

Relevant components

[zf_stack](#).

8.27 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](#).

8.28 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](#).

8.29 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](#).

8.30 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](#).

8.31 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](#).

8.32 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.

8.33 tx_timestamping Attribute Reference

Report timestamps for transmitted packets.

Detailed Description

"off" by default.

Type

Integer.

Default

0.

Relevant components

zf_vi.

Chapter 9

Data Structure Index

9.1 Data Structures

Here are the data structures with brief descriptions:

zf_attr	Attribute object	63
zf_ds	Structure used for delegated sends	64
zf_muxer_set	Multiplexer set	67
zf_pkt_report	Report structure providing timestamp and other packet information	67
zf_stack		68
zf_waitable	Abstract multiplexable object	69
zft	Opaque structure describing a TCP zocket that is connected	69
zft_handle	Opaque structure describing a TCP zocket that is passive and not connected	70
zft_msg	TCP zero-copy RX message structure	70
zftl	Opaque structure describing a TCP listening zocket	72
zfur	Opaque structure describing a UDP-receive zocket	72
zfur_msg	UDP zero-copy RX message structure	73
zfut	Opaque structure describing a UDP-transmit zocket	74

Chapter 10

File Index

10.1 File List

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

attr.h	TCPDirect API for attribute objects	75
muxer.h	TCPDirect multiplexer	81
types.h	TCPDirect types	87
x86.h	TCPDirect x86-specific definitions	89
zf.h	TCPDirect top-level API	89
zf_alts.h	TCPDirect Alternative Sends API	89
zf_ds.h	TCPDirect Delegated Sends API	94
zf_platform.h	TCPDirect platform API	98
zf_reactor.h	TCPDirect reactor API for processing stack events	98
zf_stack.h	TCPDirect stack API	101
zf_tcp.h	TCPDirect TCP API	105
zf_udp.h	TCPDirect UDP API	121

Chapter 11

Data Structure Documentation

11.1 zf_attr Struct Reference

Attribute object.

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

11.1.1 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](#)

11.2 zf_ds Struct Reference

Structure used for delegated sends.

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

Data Fields

- void * [headers](#)
- int [headers_size](#)
- int [headers_len](#)
- int [mss](#)
- int [send_wnd](#)
- int [cong_wnd](#)
- int [delegated_wnd](#)
- int [tcp_seq_offset](#)
- int [ip_len_offset](#)
- int [ip_tcp_hdr_len](#)
- int [reserved](#)

11.2.1 Detailed Description

Structure used for delegated sends.

This structure is used for delegated sends. Field usage varies:

- in: input
- out: output
- internal: internal use only.

Definition at line 24 of file zf_ds.h.

11.2.2 Field Documentation

11.2.2.1 cong_wnd

```
int cong_wnd
```

out: congestion window

Definition at line 37 of file zf_ds.h.

11.2.2.2 delegated_wnd

```
int delegated_wnd
```

out: max bytes application can send

Definition at line 40 of file zf_ds.h.

11.2.2.3 headers

```
void* headers
```

in: set to buffer to store headers to

Definition at line 26 of file zf_ds.h.

11.2.2.4 headers_len

```
int headers_len
```

out: length of headers

Definition at line 30 of file zf_ds.h.

11.2.2.5 headers_size

```
int headers_size
```

in: size of headers buffer

Definition at line 28 of file zf_ds.h.

11.2.2.6 ip_len_offset

```
int ip_len_offset
```

internal

Definition at line 45 of file zf_ds.h.

11.2.2.7 ip_tcp_hdr_len

```
int ip_tcp_hdr_len
```

internal

Definition at line 47 of file zf_ds.h.

11.2.2.8 mss

```
int mss
```

out: max segment size (max payload per packet)

Definition at line 33 of file zf_ds.h.

11.2.2.9 reserved

```
int reserved
```

internal

Definition at line 49 of file zf_ds.h.

11.2.2.10 send_wnd

```
int send_wnd
```

out: send window

Definition at line 35 of file zf_ds.h.

11.2.2.11 tcp_seq_offset

```
int tcp_seq_offset
```

internal

Definition at line 43 of file zf_ds.h.

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

- [zf_ds.h](#)

11.3 zf_muxer_set Struct Reference

Multiplexer set.

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

11.3.1 Detailed Description

Multiplexer set.

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

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

- [muxer.h](#)

11.4 zf_pkt_report Struct Reference

Report structure providing timestamp and other packet information.

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

Data Fields

- struct timespec [timestamp](#)
- uint32_t [start](#)
- uint16_t [bytes](#)
- uint16_t [flags](#)

11.4.1 Detailed Description

Report structure providing timestamp and other packet information.

This is provided by [zft_get_tx_timestamps\(\)](#) and [zft_get_tx_timestamps\(\)](#) to associate timestamps with packet data.

Definition at line 24 of file types.h.

11.4.2 Field Documentation

11.4.2.1 bytes

`uint16_t bytes`

Byte count for this packet

Definition at line 33 of file types.h.

11.4.2.2 flags

`uint16_t flags`

Flags set for this packet

Definition at line 49 of file types.h.

11.4.2.3 start

`uint32_t start`

Total count for the socket up to the start of this packet. For UDP, this is a count of datagrams. For TCP, this is a count of bytes. The counter will wrap when it reaches the end of its 32-bit range.

Definition at line 31 of file types.h.

11.4.2.4 timestamp

`struct timespec timestamp`

Hardware timestamp for packet transmission

Definition at line 26 of file types.h.

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

- [types.h](#)

11.5 zf_stack Struct Reference

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

11.5.1 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](#)

11.6 zf_waitable Struct Reference

Abstract multiplexable object.

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

11.6.1 Detailed Description

Abstract multiplexable object.

Sockets 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 socket.

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 45 of file muxer.h.

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

- [muxer.h](#)

11.7 zft Struct Reference

Opaque structure describing a TCP socket that is connected.

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

11.7.1 Detailed Description

Opaque structure describing a TCP zocket that is connected.

Definition at line 154 of file `zf_tcp.h`.

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

- [zf_tcp.h](#)

11.8 zft_handle Struct Reference

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

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

11.8.1 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](#)

11.9 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]

11.9.1 Detailed Description

TCP zero-copy RX message structure.

This structure is passed to [zft_zc_rcv\(\)](#), which will populate it and a referenced iovec array with pointers to received packets.

Definition at line 393 of file `zf_tcp.h`.

11.9.2 Field Documentation

11.9.2.1 flags

```
int flags
```

Reserved.

Definition at line 399 of file `zf_tcp.h`.

11.9.2.2 iov

```
struct iovec iov[ZF_FLEXIBLE_ARRAY_COUNT]
```

In: A separate iovec array, available for writing, with `iovcnt` entries, must immediately follow this structure. This structure and the iovec array are typically wrapped by a structure. For an example, see the [zftcpingpong](#) application.

Out: iovec array is filled with iovecs pointing to the payload of the received packets.

Definition at line 409 of file `zf_tcp.h`.

11.9.2.3 iovcnt

```
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 402 of file `zf_tcp.h`.

11.9.2.4 pkts_left

```
int pkts_left
```

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

Definition at line 397 of file `zf_tcp.h`.

11.9.2.5 reserved

```
int reserved[4]
```

Reserved.

Definition at line 395 of file zf_tcp.h.

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

- [zf_tcp.h](#)

11.10 zftl Struct Reference

Opaque structure describing a TCP listening zocket.

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

11.10.1 Detailed Description

Opaque structure describing a TCP listening zocket.

Definition at line 38 of file zf_tcp.h.

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

- [zf_tcp.h](#)

11.11 zfur Struct Reference

Opaque structure describing a UDP-receive zocket.

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

11.11.1 Detailed Description

Opaque structure describing a UDP-receive zocket.

Definition at line 27 of file zf_udp.h.

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

- [zf_udp.h](#)

11.12 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]

11.12.1 Detailed Description

UDP zero-copy RX message structure.

This structure is passed to [zfur_zc_rcv\(\)](#), which will populate it and a referenced [iovec](#) array with pointers to received packets.

Definition at line 124 of file [zf_udp.h](#).

11.12.2 Field Documentation

11.12.2.1 dgrams_left

```
int dgrams_left
```

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

Definition at line 128 of file [zf_udp.h](#).

11.12.2.2 flags

```
int flags
```

Reserved.

Definition at line 130 of file [zf_udp.h](#).

11.12.2.3 iov

```
struct iovec iov[ZF_FLEXIBLE_ARRAY_COUNT]
```

In: A separate iovec array, available for writing, with `iovcnt` entries, must immediately follow this structure. This structure and the iovec array are typically wrapped by a structure. For an example, see the [zfudpppingpong](#) application.

Out: iovec array is filled with iovecs pointing to the payload of the received packets.

Definition at line 140 of file `zf_udp.h`.

11.12.2.4 iovcnt

```
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 133 of file `zf_udp.h`.

11.12.2.5 reserved

```
int reserved[4]
```

Reserved.

Definition at line 126 of file `zf_udp.h`.

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

- [zf_udp.h](#)

11.13 zfut Struct Reference

Opaque structure describing a UDP-transmit socket.

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

11.13.1 Detailed Description

Opaque structure describing a UDP-transmit socket.

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

Definition at line 271 of file `zf_udp.h`.

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

- [zf_udp.h](#)

Chapter 12

File Documentation

12.1 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, or names of all attributes.

12.1.1 Detailed Description

TCPDirect API for attribute objects.

12.1.2 Function Documentation

12.1.2.1 zf_attr_alloc()

```
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.

12.1.2.2 zf_attr_doc()

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

Returns documentation for an attribute, or names of all attributes.

Parameters

<i>attr_name_opt</i>	The attribute name.
<i>docs_out</i>	On success, a pointer to an array of pointers to doc strings.
<i>docs_len_out</i>	On success, the number of entries in <i>docs_out</i> .

Returns

- 0 on success, or a negative error code.

This function returns pointers to the documentation for an attribute, or to the names of all attributes.

If *attr_name_opt* is the name of an attribute, the array referenced by *docs_out* contains the following strings in order:

- the name of the attribute
- the type of the attribute (e.g. "int", "str")
- the status of the attribute (e.g. "stable", "hidden", "beta")
- a description of the default value of the attribute
- the type(s) of objects the attribute applies to (e.g. "zf_stack", "zf_pool", "zf_vi")
- a description of the attribute.

If *attr_name_opt* is NULL or the empty string, the array referenced by *docs_out* instead contains all attribute names.

After calling this function, the memory it allocates for pointers must be freed by calling `free(docs_out)`.

12.1.2.3 zf_attr_dup()

```
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.

12.1.2.4 zf_attr_free()

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

Free an attribute object.

Parameters

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

12.1.2.5 zf_attr_get_int()

```
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

12.1.2.6 zf_attr_get_str()

```
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 <code>strdup()</code> and must be <code>free()</code> 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

12.1.2.7 `zf_attr_reset()`

```
void zf_attr_reset (
    struct zf_attr * attr )
```

Return attributes to their default values.

Parameters

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

12.1.2.8 `zf_attr_set_from_fmt()`

```
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.

12.1.2.9 zf_attr_set_from_str()

```
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.

12.1.2.10 zf_attr_set_int()

```
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.

12.1.2.11 zf_attr_set_str()

```
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.

12.2 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.
- const 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.

12.2.1 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)`.

12.2.2 Function Documentation

12.2.2.1 `zf_muxer_add()`

```
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.

12.2.2.2 zf_muxer_alloc()

```
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\(\)](#).

12.2.2.3 zf_muxer_del()

```
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.

12.2.2.4 zf_muxer_free()

```
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.

12.2.2.5 zf_muxer_mod()

```
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));
```

12.2.2.6 zf_muxer_wait()

```
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 the 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\(timeout_ns=0\)](#).

12.2.2.7 zf_waitable_event()

```
const 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.

12.2.2.8 zf_waitable_fd_get()

```
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 zocket. 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 zockets that the caller is interested in. This may or may not result in a zocket within the stack becoming readable or writeable.

Using a waitable FD with the [zf_waitable_fd...\(\)](#) family of functions enables interrupts. For latency-critical applications, you should instead manually poll each reactor in turn, after first setting the [reactor_spin_count](#) attribute to 1.

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.

12.2.2.9 zf_waitable_fd_prime()

```
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.

12.3 types.h File Reference

TCPDirect types.

Data Structures

- struct [zf_pkt_report](#)
Report structure providing timestamp and other packet information.

Macros

- #define [ZF_PKT_REPORT_CLOCK_SET](#) 0x0001
- #define [ZF_PKT_REPORT_IN_SYNC](#) 0x0002
- #define [ZF_PKT_REPORT_NO_TIMESTAMP](#) 0x0004
- #define [ZF_PKT_REPORT_DROPPED](#) 0x0008
- #define [ZF_PKT_REPORT_TCP_RETRANS](#) 0x2000
- #define [ZF_PKT_REPORT_TCP_SYN](#) 0x4000
- #define [ZF_PKT_REPORT_TCP_FIN](#) 0x8000

12.3.1 Detailed Description

TCPDirect types.

12.3.2 Macro Definition Documentation

12.3.2.1 ZF_PKT_REPORT_CLOCK_SET

```
#define ZF_PKT_REPORT_CLOCK_SET 0x0001
```

Adapter clock has been set

Definition at line 35 of file types.h.

12.3.2.2 ZF_PKT_REPORT_DROPPED

```
#define ZF_PKT_REPORT_DROPPED 0x0008
```

Dropped reports before this

Definition at line 41 of file types.h.

12.3.2.3 ZF_PKT_REPORT_IN_SYNC

```
#define ZF_PKT_REPORT_IN_SYNC 0x0002
```

Adapter clock is in sync

Definition at line 37 of file types.h.

12.3.2.4 ZF_PKT_REPORT_NO_TIMESTAMP

```
#define ZF_PKT_REPORT_NO_TIMESTAMP 0x0004
```

No timestamp available

Definition at line 39 of file types.h.

12.3.2.5 ZF_PKT_REPORT_TCP_FIN

```
#define ZF_PKT_REPORT_TCP_FIN 0x8000
```

Final TCP FIN packet

Definition at line 47 of file types.h.

12.3.2.6 ZF_PKT_REPORT_TCP_RETRANS

```
#define ZF_PKT_REPORT_TCP_RETRANS 0x2000
```

Retransmitted TCP packet

Definition at line 43 of file types.h.

12.3.2.7 ZF_PKT_REPORT_TCP_SYN

```
#define ZF_PKT_REPORT_TCP_SYN 0x4000
```

Initial TCP SYN packet

Definition at line 45 of file types.h.

12.4 x86.h File Reference

TCPDirect x86-specific definitions.

12.4.1 Detailed Description

TCPDirect x86-specific definitions.

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

12.5 zf.h File Reference

TCPDirect top-level API.

12.5.1 Detailed Description

TCPDirect top-level API.

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

12.6 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.

12.6.1 Detailed Description

TCPDirect Alternative Sends API.

12.6.2 Function Documentation

12.6.2.1 zf_alternatives_alloc()

```
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 Attributes 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.

Note

TX alternatives are not supported on stacks running on bonded network interfaces.

See also

[zf_alternatives_release\(\)](#)

12.6.2.2 zf_alternatives_cancel()

```
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.

12.6.2.3 [zf_alternatives_free_space\(\)](#)

```
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\(\)](#).

12.6.2.4 [zf_alternatives_query_overhead_tcp\(\)](#)

```
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.

12.6.2.5 `zf_alternatives_release()`

```
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\(\)](#)

12.6.2.6 `zf_alternatives_send()`

```
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 queuing 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.

12.6.2.7 zft_alternatives_queue()

```
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 zf_alternatives_alloc()
<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.

-ENOBUFS 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 zocket.

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.

12.7 zf_ds.h File Reference

TCPDirect Delegated Sends API.

Data Structures

- struct `zf_ds`
Structure used for delegated sends.

Macros

- #define `ZF_DELEGATED_SEND_RC_FATAL` 0x80
Mask to test for fatal errors in the Delegated Sends API.

Enumerations

- enum `zf_delegated_send_rc` {
`ZF_DELEGATED_SEND_RC_OK` = 0x00, `ZF_DELEGATED_SEND_RC_NOCWIN` = 0x01,
`ZF_DELEGATED_SEND_RC_NOWIN` = 0x02, `ZF_DELEGATED_SEND_RC_BAD_SOCKET` = 0x83,
`ZF_DELEGATED_SEND_RC_SMALL_HEADER` = 0x84, `ZF_DELEGATED_SEND_RC_SENDQ_BUSY` =
0x85, `ZF_DELEGATED_SEND_RC_NOARP` = 0x86 }
Return codes for functions in the Delegated Sends API.

Functions

- enum `zf_delegated_send_rc` `zf_delegated_send_prepare` (struct `zft` *ts, int max_delegated_wnd, int cong_wnd_override, unsigned flags, struct `zf_ds` *ds)
Delegate sends to the application.
- static void `zf_delegated_send_tcp_update` (struct `zf_ds` *ds, int bytes, int push)
Update packet headers with correct data length and PUSH flag.
- static void `zf_delegated_send_tcp_advance` (struct `zf_ds` *ds, int bytes)
Update packet headers to reflect that a packet has been sent.
- int `zf_delegated_send_complete` (struct `zft` *ts, const struct iovec *iov, int iovlen, int flags)
Notify TCPDirect that some data have been sent via delegated sends.
- int `zf_delegated_send_cancel` (struct `zft` *ts)
Notify TCPDirect that a reserved set of bytes are no longer required.

12.7.1 Detailed Description

TCPDirect Delegated Sends API.

12.7.2 Enumeration Type Documentation

12.7.2.1 `zf_delegated_send_rc`

```
enum zf_delegated_send_rc
```

Return codes for functions in the Delegated Sends API.

Enumerator

ZF_DELEGATED_SEND_RC_OK	Success
ZF_DELEGATED_SEND_RC_NOCWIN	Insufficient congestion window
ZF_DELEGATED_SEND_RC_NOWIN	Insufficient send window
ZF_DELEGATED_SEND_RC_BAD_SOCKET	Zocket not in a state to send
ZF_DELEGATED_SEND_RC_SMALL_HEADER	headers_size too small (headers_len gives size required)
ZF_DELEGATED_SEND_RC_SENDQ_BUSY	Zocket has data in send queue
ZF_DELEGATED_SEND_RC_NOARP	Remote MAC for peer not known

Definition at line 54 of file `zf_ds.h`.

12.7.3 Function Documentation

12.7.3.1 `zf_delegated_send_cancel()`

```
int zf_delegated_send_cancel (
    struct zft * ts )
```

Notify TCPDirect that a reserved set of bytes are no longer required.

Parameters

<code>ts</code>	TCP zocket
-----------------	------------

Returns

0 on success, or negative error on failure

Notify TCPDirect that a previously reserved set of bytes (obtained using [zf_delegated_send_prepare\(\)](#)) are no longer required.

This must be used if the caller has not called [zf_delegated_send_complete\(\)](#) for all the bytes reserved. After successful return, the caller can use other TCPDirect send API calls, or start another delegated send operation with [zf_delegated_send_prepare\(\)](#).

12.7.3.2 zf_delegated_send_complete()

```
int zf_delegated_send_complete (
    struct zft * ts,
    const struct iovec * iov,
    int iovlen,
    int flags )
```

Notify TCPDirect that some data have been sent via delegated sends.

Parameters

<i>ts</i>	TCP zocket
<i>iov</i>	Start of the iovec array describing the packet buffers
<i>iovlen</i>	Length of the iovec array
<i>flags</i>	Reserved for future use

Returns

Number of bytes completed on success (which may be less than the requested number of bytes for partial success)

Negative error on failure:

- EMSGSIZE: attempt to "complete" more bytes than were "prepared"
- EAGAIN: no space on send queue (prepare should already have failed)

Notify TCPDirect that some data have been sent via delegated sends. If successful, TCPDirect will handle all further aspects of the TCP protocol (e.g. acknowledgements, retransmissions) for those bytes.

12.7.3.3 zf_delegated_send_prepare()

```
enum zf_delegated_send_rc zf_delegated_send_prepare (
    struct zft * ts,
    int max_delegated_wnd,
    int cong_wnd_override,
    unsigned flags,
    struct zf_ds * ds )
```

Delegate sends to the application.

Parameters

<i>ts</i>	TCP zocket
<i>max_delegated_wnd</i>	Bytes to reserve for future delegated sends
<i>cong_wnd_override</i>	Minimum congestion window, or zero
<i>flags</i>	Reserved for future use
<i>ds</i>	Structure used for delegated sends

Returns

ZF_DELEGATED_SEND_RC_OK: Success
ZF_DELEGATED_SEND_RC_NOCWIN: Insufficient congestion window
ZF_DELEGATED_SEND_RC_NOWIN: Insufficient send window
ZF_DELEGATED_SEND_RC_BAD_SOCKET: Zocket not in a state to send
ZF_DELEGATED_SEND_RC_SMALL_HEADER: headers_size too small (headers_len gives size required)
ZF_DELEGATED_SEND_RC_SENDQ_BUSY: Zocket has data in send queue
ZF_DELEGATED_SEND_RC_NOARP: Remote MAC for peer not known

This function delegates sends to the application. It reserves up to `max_delegated_wnd` bytes for future delegated sends, and returns the Ethernet-IP-TCP headers. The maximum amount of data that the application can send is then returned in `ds->delegated_wnd`. Both `max_delegated_wnd` and `ds->delegated_wnd` are relative to sends already completed.

If `cong_wnd_override` is non-zero, it specifies a minimum congestion window. This call behaves as if the congestion window is the larger of `cong_wnd_override` and the zocket's actual congestion window.

Once a send has been completed, call [zf_delegated_send_complete\(\)](#) to indicate how many bytes were used, and make further calls to this function to extend the window for delegated sends.

If not all the bytes reserved with this call are used then [zf_delegated_send_cancel\(\)](#) must be called before further normal sends. A subsequent call to [zf_delegated_send_prepare\(\)](#) is safe without calling [zf_delegated_send_cancel\(\)](#).

When the return code is RC_OK, RC_NOWIN or RC_NOCWIN, the headers and other fields in `ds` are initialised. This set is indicated by `!(rc & ZF_DELEGATED_SEND_RC_FATAL)`. When RC_SMALL_HEADER is returned, `ds->headers_len` is initialised. In other cases `ds` is not filled in.

Note that the delegated window is never reduced by this call, so `ds->delegated_wnd` may be non-zero even if RC_NOWIN or RC_NOCWIN is returned.

12.7.3.4 zf_delegated_send_tcp_advance()

```
static void zf_delegated_send_tcp_advance (
    struct zf_ds * ds,
    int bytes ) [inline], [static]
```

Update packet headers to reflect that a packet has been sent.

Parameters

<i>ds</i>	Structure used for delegated sends
<i>bytes</i>	Bytes sent

Update packet headers created by [zf_delegated_send_prepare\(\)](#) to reflect that a packet of length `bytes` has been sent.

[zf_delegated_send_prepare\(\)](#) reserves a potentially long area for delegated sends. If these bytes are sent in multiple packets, this function must be used in between each delegated send to update the TCP headers appropriately.

Definition at line 179 of file `zf_ds.h`.

12.7.3.5 zf_delegated_send_tcp_update()

```
static void zf_delegated_send_tcp_update (
    struct zf_ds * ds,
    int bytes,
    int push ) [inline], [static]
```

Update packet headers with correct data length and PUSH flag.

Parameters

<i>ds</i>	Structure used for delegated sends
<i>bytes</i>	Correct data length
<i>push</i>	Zero to clear PUSH flag, non-zero to set PUSH flag

Update packet headers created by [zf_delegated_send_prepare\(\)](#) with correct data length and PUSH flag details.

[zf_delegated_send_prepare\(\)](#) assumes that the delegated send will be the maximum segment size, and that no PUSH flag will be set in the TCP header. If this assumption is correct there is no need to call [zf_delegated_send_tcp_update\(\)](#).

Definition at line 142 of file `zf_ds.h`.

12.8 zf_platform.h File Reference

TCPDirect platform API.

12.8.1 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.

12.9 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_reactor_perform_attr](#) (struct [zf_stack](#) *st, const struct [zf_attr](#) *attr)
Process events on a stack, with overridden attributes.
- int [zf_stack_has_pending_work](#) (const struct [zf_stack](#) *st)
Determine whether a stack has work pending.
- int [zf_stack_has_pending_events](#) (const struct [zf_stack](#) *st)
Determine whether a stack has events pending, but don't check TCP-specific non-event-based work.

12.9.1 Detailed Description

TCPDirect reactor API for processing stack events.

12.9.2 Function Documentation

12.9.2.1 `zf_reactor_perform()`

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

Process events on a stack.

Parameters

<code>st</code>	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 `zf_reactor_perform()`, `zf_reactor_perform_attr()` 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()`. To override `reactor_spin_count` for a single call, `zf_reactor_perform_attr()` can be used instead.

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_reactor_perform_attr\(\)](#) [zf_muxer_wait\(\)](#)

12.9.2.2 `zf_reactor_perform_attr()`

```
int zf_reactor_perform_attr (
    struct zf_stack * st,
    const struct zf_attr * attr )
```

Process events on a stack, with overridden attributes.

Parameters

<i>st</i>	Stack for which to process events.
<i>attr</i>	Overridden properties for event processing. Only reactor_spin_count is currently supported.

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 [zf_reactor_perform\(\)](#), [zf_reactor_perform_attr\(\)](#) or [zf_muxer_wait\(\)](#) frequently for each stack that is in use. Please see [Stack polling](#) in the User Guide for further information.

This function differs from [zf_reactor_perform\(\)](#) in that the reactor_spin_count stack attribute will be overridden using the provided attributes. In all other respects it is identical to [zf_reactor_perform\(\)](#).

reactor_spin_count is the only supported override at this time. Other attributes may be added in future versions so callers need to take care with the setting of other attributes to avoid unintended side effects when run against future versions.

This function polls repeatedly for events. The amount of time spent polling is controlled by the attribute reactor_spin_count. Setting reactor_spin_count to 1 disables polling and minimises the cost of [zf_reactor_perform_attr\(\)](#).

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_reactor_perform\(\)](#) [zf_muxer_wait\(\)](#)

12.9.2.3 zf_stack_has_pending_events()

```
int zf_stack_has_pending_events (
    const struct zf_stack * st )
```

Determine whether a stack has events pending, but don't check TCP-specific non-event-based work.

Parameters

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

This function is a cut-down version of [zf_stack_has_pending_work\(\)](#). It returns non-zero if the stack has events pending, and therefore the application should call [zf_reactor_perform\(\)](#), [zf_reactor_perform_attr\(\)](#) or [zf_muxer_wait\(\)](#).

This differs from [zf_stack_has_pending_work\(\)](#) in that it never tries to check whether there is non-event-based work (such as processing TCP timers) pending. If the calling application knows there is no TCP work (e.g. it is using only UDP sockets) this function may be a few cycles cheaper.

12.9.2.4 zf_stack_has_pending_work()

```
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\(\)](#), [zf_reactor_perform_attr\(\)](#) 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_reactor_perform_attr\(\)](#) [zf_muxer_wait\(\)](#)

12.10 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.
- `const char * zf_version (void)`
Returns library name and version.
- `void zf_print_version (void)`
Prints library name and version to stderr, then exits.

12.10.1 Detailed Description

TCPDirect stack API.

12.10.2 Macro Definition Documentation

12.10.2.1 EPOLLSTACKHUP

```
#define EPOLLSTACKHUP EPOLLRDHUP
```

Event indicating stack quiescence.

See also

[zf_stack_to_waitable\(\)](#)

Definition at line 90 of file `zf_stack.h`.

12.10.3 Function Documentation

12.10.3.1 zf_deinit()

```
int zf_deinit (
    void )
```

Deinitialize zf library.

Returns

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

12.10.3.2 zf_init()

```
int zf_init (
    void )
```

Initialize zf library.

Should be called exactly once per process, and before any other API calls are made.

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.

12.10.3.3 zf_stack_alloc()

```
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)`.

12.10.3.4 zf_stack_free()

```
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.

12.10.3.5 `zf_stack_is_quiescent()`

```
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.

12.10.3.6 `zf_stack_to_waitable()`

```
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.

12.11 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 zocket.
- struct [zft](#)
Opaque structure describing a TCP zocket 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 zocket.
- 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 zocket.
- int [zftl_free](#) (struct [zftl](#) *ts)
Release resources associated with a TCP listening zocket.
- 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 zocket.
- int [zft_handle_free](#) (struct [zft_handle](#) *handle)
Release a handle to a TCP zocket.
- 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 zocket.
- int [zft_shutdown_tx](#) (struct [zft](#) *ts)
Shut down outgoing TCP connection.
- int [zft_free](#) (struct [zft](#) *ts)
Release resources associated with a TCP zocket.
- 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_recv` (struct `zft` *ts, struct `zft_msg` *msg, int flags)

Zero-copy read of available packets.

 - int `zft_zc_recv_done` (struct `zft` *ts, struct `zft_msg` *msg)

Concludes pending `zc_recv` operation as done.

 - 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.

 - int `zft_recv` (struct `zft` *ts, const struct `iovec` *iov, int iovcnt, int flags)

Copy-based receive.

 - int `zft_pkt_get_timestamp` (struct `zft` *ts, const struct `zft_msg` *msg, struct `timespec` *ts_out, int pktind, unsigned *flags)

Retrieve the UTC timestamp associated with a received packet, and the clock sync status flags.

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

Send data specified in `iovec` array.

 - `ZF_NOCLONE` `ssize_t` `zft_send_single` (struct `zft` *ts, const void *buf, `size_t` buflen, int flags)

Send data given in single buffer.

 - `ssize_t` `zft_send_single_warm` (struct `zft` *ts, const void *buf, `size_t` buflen)

Warms code path used by `zft_send_single()` without sending data.

 - 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.

 - int `zft_get_tx_timestamps` (struct `zft` *ts, struct `zft_pkt_report` *reports, int *count_in_out)

Retrieve timestamp reports from previously sent data.

12.11.1 Detailed Description

TCPDirect TCP API.

12.11.2 Function Documentation

12.11.2.1 `zft_addr_bind()`

```
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.

12.11.2.2 zft_alloc()

```
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 zf_stack .
<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 Attributes 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\(\)](#)

12.11.2.3 zft_connect()

```
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.
- EINVAL Zocket in unexpected TCP state, or no raddr supplied

See also

[zft_addr_bind\(\)](#)

12.11.2.4 zft_error()

```
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.

12.11.2.5 zft_free()

```
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.

12.11.2.6 zft_get_header_size()

```
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.

12.11.2.7 zft_get_mss()

```
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.

12.11.2.8 zft_get_tx_timestamps()

```
int zft_get_tx_timestamps (
    struct zft * ts,
    struct zf_pkt_report * reports,
    int * count_in_out )
```

Retrieve timestamp reports from previously sent data.

Parameters

<i>ts</i>	TCP zocket.
<i>reports</i>	Array to fill with timestamp reports
<i>count_in_out</i>	IN: size of array, OUT: number of reports

Returns

- 0 on success, or negative error code.

If transmit timestamps are enabled, then one report will be generated for each segment. The segment can be identified by the "start" field of the report, which begins at 0 and increments for each byte sent on this zocket. Retransmission will cause multiple reports for that segment, and is indicated by the ZF_PKT_REPORT_TCP_RETRANS flag as well as discontinuities in the reported location. Timestamps are also reported for the packets sent to open and close the stream, indicated by the ZF_PKT_REPORT_TCP_SYN and ZF_PKT_REPORT_TCP_FIN flags.

12.11.2.9 zft_getname()

```
void zft_getname (
    struct zft * ts,
    struct sockaddr * laddr_out,
    socklen_t * laddr_len,
    struct sockaddr * raddr_out,
    socklen_t * raddr_len )
```

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.

12.11.2.10 zft_handle_free()

```
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.

12.11.2.11 zft_handle_getname()

```
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.

12.11.2.12 zft_pkt_get_timestamp()

```
int zft_pkt_get_timestamp (
    struct zft * ts,
    const struct zft_msg * msg,
    struct timespec * ts_out,
    int pktind,
    unsigned * flags )
```

Retrieve the UTC timestamp associated with a received packet, and the clock sync status flags.

Parameters

<i>ts</i>	TCP zocket.
<i>msg</i>	Pointer to the received message for which the RX timestamp will be retrieved.
<i>ts_out</i>	Pointer to a timespec that is updated on return with the UTC timestamp for the packet.
<i>pktind</i>	Index of packet within <code>msg->iov</code> .
<i>flags</i>	Pointer to an unsigned that is updated on return with the sync flags for the packet.

Returns

- 0 Success.
- ENOMSG Synchronisation with adapter has not yet been achieved. This only happens with old firmware.
- ENODATA Packet does not have a timestamp. On current Solarflare adapters, packets that are switched from TX to RX do not get timestamped.
- EL2NSYNC Synchronisation with adapter has been lost. This should never happen!

Note

This function must be called after `zf_reactor_perform()` returns a value greater than zero, and before `zf_reactor_perform()` is called again.
If RX timestamps were not enabled during stack initialisation, the behaviour of this function is undefined.

On success the `ts_out` and `flags_out` fields are updated, and a value of zero is returned. The `flags_out` field contains the following flags:

- `EF_VI_SYNC_FLAG_CLOCK_SET` is set if the adapter clock has ever been set (in sync with system)
- `EF_VI_SYNC_FLAG_CLOCK_IN_SYNC` is set if the adapter clock is in sync with the external clock (PTP).

12.11.2.13 zft_recv()

```
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.

12.11.2.14 zft_send()

```
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 MSG_MORE.

This function adds the supplied data (as indicated by the *iov* argument) to the zocket's send queue and if possible will send it (or part of it) on the wire. To prevent a small packet being sent the MSG_MORE flag can be used: it will prevent a packet that is not filled up to MSS from being sent.

There is no guarantee that separate calls to this function, or separate entries in the iovec array, will result in separate packets. To achieve control over packet boundaries the delegated sends API can be used instead.

Provided buffers may be re-used on 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 to send any portion of the data.
- 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`.

This function will send only part of the data provided if there is insufficient space in the send queue to send all of it (and there are no error conditions). Use `zft_send_space()` immediately before this call to determine in advance whether only part of the data would be sent.

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.

12.11.2.15 zft_send_single()

```
ZF_NOCLONE 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 zocket 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 <code>MSG_MORE</code> .

This function adds the supplied data (as indicated by the `buf` argument) to the zocket's send queue and if possible will send it (or part of it) on the wire. To prevent a small packet being sent the `MSG_MORE` flag can be used: it will prevent a packet that is not filled up to MSS from being sent.

There is no guarantee that separate calls to this function will result in separate packets. To achieve control over packet boundaries the delegated sends API can be used instead. The "single" in the name of the function refers to it taking a single buffer rather than an `iovec` of buffers.

Provided buffer may be re-used on 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 to send any portion of the data.
- 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.

This function will send only part of the data provided if there is insufficient space in the send queue to send all of it (and there are no error conditions). Use [zft_send_space\(\)](#) immediately before this call to determine in advance whether only part of the data would be sent.

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.

12.11.2.16 zft_send_single_warm()

```
ssize_t zft_send_single_warm (  
    struct zft * ts,  
    const void * buf,  
    size_t buflen )
```

Warms code path used by [zft_send_single\(\)](#) without sending data.

Parameters

<i>ts</i>	The TCP zocket to send on.
<i>buf</i>	The buffer of data to send.
<i>buflen</i>	The length of buffer.

Returns

Number of bytes warmed on success.

- EAGAIN Events need to be processed before warming. Call [zf_reactor_perform\(\)](#)
- EMSGSIZE Data buffer too long.
- ENOTCONN Zocket is not in a valid TCP state for sending.
- ENOMEM Not enough packet buffers available.

This function can be called repeatedly while the application waits for an input that will trigger a call to `zft_send_single()`. Doing so warms the code path to avoid cache and TLB misses when actually sending data in the subsequent `zft_send_single()` call. `buf` need not contain the exact data that will eventually be sent.

This function only supports warming the code path where the send queue is empty and a PIO or CTPIO send would be performed. If `buflen` is too large for PIO then `-EMSGSIZE` will be returned. If previous sends may still be in progress `-EAGAIN` will be returned. In this case, the application can call `zf_reactor_perform()` and then try again.

See also

[zft_send_single\(\)](#)

12.11.2.17 zft_send_space()

```
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`.

12.11.2.18 zft_shutdown_tx()

```
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.
- EBUSY Delegated send in progress

12.11.2.19 zft_state()

```
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).

12.11.2.20 zft_to_waitable()

```
struct zf_waitable* zft_to_waitable (  
    struct zft * ts )
```

Returns a [zf_waitable](#) representing the given [zft](#).

Parameters

<i>ts</i>	The zft to return as a zf_waitable
-----------	--

Returns

The [zf_waitable](#)

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

12.11.2.21 zft_zc_rcv()

```
void zft_zc_rcv (  
    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 and its referenced `iovec` array 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 the reason of stream termination check the 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()`.

12.11.2.22 zft_zc_recv_done()

```
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()`.

12.11.2.23 zft_zc_recv_done_some()

```
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_rcv_done\(\)](#).

Can be called after each successful [zft_zc_rcv\(\)](#) operation as an alternative to [zft_zc_rcv_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_rcv_done\(\)](#). The function releases resources and enables subsequent calls to [zft_zc_rcv\(\)](#) or [zft_rcv\(\)](#). [zft_zc_rcv\(\)](#) or [zft_rcv\(\)](#) 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_rcv\(\)](#). *len* must not be greater than total payload returned by [zft_zc_rcv\(\)](#).

12.11.2.24 zftl_accept()

```
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.

12.11.2.25 zftl_free()

```
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.

Note

The listening zocket is not removed until all accepted zockets have also been freed. If any connections to the listening zocket have been accepted, but the resulting zocket has not been freed by calling `zft_free()`, the listening zocket remains. It will not accept any new connections, and is shown in the output from `zf_stackdump`. Attempting to create an additional listening zocket on the same port results in an error.

Returns

0 Success.

12.11.2.26 zftl_getname()

```
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.

12.11.2.27 zftl_listen()

```
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 <code>zf_stack</code> 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 <code>INADDR_ANY</code>).
<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 Attributes 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.
- ENOBUFS No zockets of this type available.
- ENOMEM Out of memory.
- EOPNOTSUPP `laddr` is INADDR_ANY.

12.11.2.28 `zftl_to_waitable()`

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

Returns a `zf_waitable` representing `tl`.

Parameters

<code>tl</code>	The <code>zftl</code> to return as a <code>zf_waitable</code>
-----------------	---

Returns

The `zf_waitable`

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

12.12 `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 socket.
- `int zfur_free (struct zfur *us)`
Release UDP-receive socket 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 socket 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 socket from address.
- `void zfur_zc_recv (struct zfur *us, struct zfur_msg *msg, int flags)`
Zero-copy read of single datagram.
- `void zfur_zc_recv_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.
- `int zfur_pkt_get_timestamp (struct zfur *us, const struct zfur_msg *msg, struct timespec *ts_out, int pktind, unsigned *flags)`
Retrieve the UTC timestamp associated with a received packet, and the clock sync status flags.
- `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 socket.
- `int zfut_free (struct zfut *us)`
Free UDP-transmit socket.
- `int zfut_get_mss (struct zfut *us)`
Get the maximum segment size which can be transmitted.
- `ZF_NOCLONE 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_single_warm (struct zfut *us, const void *buf, size_t buflen)`
Warms code path used by `zfut_send_single()` without sending data.
- `int zfut_send (struct zfut *us, const struct iovec *iov, int iov_cnt, int flags)`
Copy-based send of single UDP packet (possibly fragmented).
- `int zfut_get_tx_timestamps (struct zfut *us, struct zf_pkt_report *reports_out, int *count_in_out)`
Retrieve timestamp reports from previously sent data.
- `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 socket.

12.12.1 Detailed Description

TCPDirect UDP API.

12.12.2 Function Documentation

12.12.2.1 zfur_addr_bind()

```
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 laddr.
<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 raddr.
<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.

12.12.2.2 zfur_addr_unbind()

```
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\(\)](#).

12.12.2.3 zfur_alloc()

```
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 zf_stack 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 Attributes for details.

Returns

- 0 Success.
- ENOBUFS 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.

12.12.2.4 zfur_free()

```
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.

12.12.2.5 zfur_pkt_get_header()

```
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->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`.

12.12.2.6 zfur_pkt_get_timestamp()

```
int zfur_pkt_get_timestamp (
    struct zfur * us,
    const struct zfur_msg * msg,
    struct timespec * ts_out,
    int pktind,
    unsigned * flags )
```

Retrieve the UTC timestamp associated with a received packet, and the clock sync status flags.

Parameters

<i>us</i>	UDP zocket.
<i>msg</i>	Pointer to the received message for which the RX timestamp will be retrieved.
<i>ts_out</i>	Pointer to a <code>timespec</code> that is updated on return with the UTC timestamp for the packet.
<i>pktind</i>	Index of packet within <code>msg->iov</code> .
<i>flags</i>	Pointer to an unsigned that is updated on return with the sync flags for the packet.

Returns

- 0 Success.
- ENOMSG Synchronisation with adapter has not yet been achieved. This only happens with old firmware.
- ENODATA Packet does not have a timestamp. On current Solarflare adapters, packets that are switched from TX to RX do not get timestamped.
- EL2NSYNC Synchronisation with adapter has been lost. This should never happen!

Note

This function must be called after `zf_reactor_perform()` returns a value greater than zero, and before `zf_reactor_perform()` is called again.
If RX timestamps were not enabled during stack initialisation, the behaviour of this function is undefined.

On success the `ts_out` and `flags_out` fields are updated, and a value of zero is returned. The `flags_out` field contains the following flags:

- `EF_VI_SYNC_FLAG_CLOCK_SET` is set if the adapter clock has ever been set (in sync with system)
- `EF_VI_SYNC_FLAG_CLOCK_IN_SYNC` is set if the adapter clock is in sync with the external clock (PTP).

12.12.2.7 `zfur_to_waitable()`

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

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

Parameters

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

Returns

The `zf_waitable`

This is necessary for use with the multiplexer.

12.12.2.8 `zfur_zc_recv()`

```
void zfur_zc_recv (
    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()`.

12.12.2.9 zfur_zc_rcv_done()

```
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_rcv()` operation that returns at least one packet. It must not be called otherwise (in particular, when `zfur_zc_rcv()` returned no packets). The function releases resources and enables subsequent calls to `zfur_zc_rcv()`. `msg` must be passed unmodified from the call to `zfur_zc_rcv()`.

12.12.2.10 zfut_alloc()

```
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 Attributes 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.

12.12.2.11 `zfut_free()`

```
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.

12.12.2.12 `zfut_get_header_size()`

```
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.

12.12.2.13 `zfut_get_mss()`

```
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.

12.12.2.14 `zfut_get_tx_timestamps()`

```
int zfut_get_tx_timestamps (  
    struct zfut * us,  
    struct zf_pkt_report * reports_out,  
    int * count_in_out )
```

Retrieve timestamp reports from previously sent data.

Parameters

<i>us</i>	UDP zocket.
<i>reports_out</i>	Array to fill with timestamp reports
<i>count_in_out</i>	IN: size of array, OUT: number of reports

Returns

0 on success
negative error code on failure. None are currently specified.

If transmit timestamps are enabled, then one report will be generated for each packet. The packet can be identified by the "start" field of the report, which begins at 0 and increments after each packet is sent on this zocket. If a packet is fragmented, then a single report will be generated with the timestamp for the final fragment.

12.12.2.15 `zfut_send()`

```
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.
- EAGAIN Hardware queue full. Call [zf_reactor_perform\(\)](#) until it returns non-zero and try again.
- EMSGSIZE Message too large.
- ENOBUFS 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\(\)](#)

12.12.2.16 zfut_send_single()

```
ZF_NOCLONE 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\(\)](#)

12.12.2.17 `zfut_send_single_warm()`

```
int zfut_send_single_warm (
    struct zfut * us,
    const void * buf,
    size_t buflen )
```

Warms code path used by `zfut_send_single()` without sending data.

Parameters

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

Returns

- Payload bytes that would have been sent (i.e. `buflen`) on success.
- EAGAIN Events need to be processed before warming. Call `zf_reactor_perform()`
- EMSGSIZE Message too large.
- ENOBUFS Out of packet buffers.

This function can be called repeatedly while the application waits for an input that will trigger a call to `zfut_send_single()`. Doing so warms the code path to avoid cache and TLB misses when actually sending data in the subsequent `zfut_send_single()` call. `buf` need not contain the data that will eventually be sent.

This function only supports warming the code path where a PIO or CTPIO send would be performed. If `buflen` is too large for PIO then -EMSGSIZE will be returned. If PIO is currently in use then -EAGAIN will be returned. In this case, the application can call `zf_reactor_perform()` and then try again.

See also

[zfut_send_single\(\)](#)

12.12.2.18 `zfut_to_waitable()`

```
struct zf_waitable* zfut_to_waitable (
    struct zfut * us )
```

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

Parameters

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

Returns

The `zf_waitable`.

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

Index

- attr.h, [75](#)
 - [zf_attr_alloc, 76](#)
 - [zf_attr_doc, 76](#)
 - [zf_attr_dup, 76](#)
 - [zf_attr_free, 77](#)
 - [zf_attr_get_int, 77](#)
 - [zf_attr_get_str, 77](#)
 - [zf_attr_reset, 79](#)
 - [zf_attr_set_from_fmt, 79](#)
 - [zf_attr_set_from_str, 79](#)
 - [zf_attr_set_int, 80](#)
 - [zf_attr_set_str, 80](#)
- bytes
 - [zf_pkt_report, 67](#)
- cong_wnd
 - [zf_ds, 64](#)
- delegated_wnd
 - [zf_ds, 64](#)
- dgrams_left
 - [zfur_msg, 73](#)
- EPOLLSTACKHUP
 - [zf_stack.h, 102](#)
- flags
 - [zf_pkt_report, 68](#)
 - [zft_msg, 71](#)
 - [zfur_msg, 73](#)
- headers
 - [zf_ds, 65](#)
- headers_len
 - [zf_ds, 65](#)
- headers_size
 - [zf_ds, 65](#)
- iov
 - [zft_msg, 71](#)
 - [zfur_msg, 73](#)
- iovcnt
 - [zft_msg, 71](#)
 - [zfur_msg, 74](#)
- ip_len_offset
 - [zf_ds, 65](#)
- ip_tcp_hdr_len
 - [zf_ds, 65](#)
- mss
 - [zf_ds, 66](#)
- muxer.h, [81](#)
 - [zf_muxer_add, 82](#)
 - [zf_muxer_alloc, 82](#)
 - [zf_muxer_del, 83](#)
 - [zf_muxer_free, 83](#)
 - [zf_muxer_mod, 84](#)
 - [zf_muxer_wait, 84](#)
 - [zf_waitable_event, 85](#)
 - [zf_waitable_fd_get, 85](#)
 - [zf_waitable_fd_prime, 86](#)
- pkts_left
 - [zft_msg, 71](#)
- reserved
 - [zf_ds, 66](#)
 - [zft_msg, 71](#)
 - [zfur_msg, 74](#)
- send_wnd
 - [zf_ds, 66](#)
- start
 - [zf_pkt_report, 68](#)
- tcp_seq_offset
 - [zf_ds, 66](#)
- timestamp
 - [zf_pkt_report, 68](#)
- types.h, [87](#)
 - [ZF_PKT_REPORT_CLOCK_SET, 87](#)
 - [ZF_PKT_REPORT_DROPPED, 87](#)
 - [ZF_PKT_REPORT_IN_SYNC, 87](#)
 - [ZF_PKT_REPORT_NO_TIMESTAMP, 88](#)
 - [ZF_PKT_REPORT_TCP_FIN, 88](#)
 - [ZF_PKT_REPORT_TCP_RETRANS, 88](#)
 - [ZF_PKT_REPORT_TCP_SYN, 88](#)
- x86.h, [89](#)
 - [ZF_PKT_REPORT_CLOCK_SET](#)
 - [types.h, 87](#)
 - [ZF_PKT_REPORT_DROPPED](#)
 - [types.h, 87](#)
 - [ZF_PKT_REPORT_IN_SYNC](#)
 - [types.h, 87](#)
 - [ZF_PKT_REPORT_NO_TIMESTAMP](#)
 - [types.h, 88](#)

ZF_PKT_REPORT_TCP_FIN
types.h, 88

ZF_PKT_REPORT_TCP_RETRANS
types.h, 88

ZF_PKT_REPORT_TCP_SYN
types.h, 88

zf.h, 89

zf_alternatives_alloc
zf_alts.h, 90

zf_alternatives_cancel
zf_alts.h, 90

zf_alternatives_free_space
zf_alts.h, 91

zf_alternatives_query_overhead_tcp
zf_alts.h, 91

zf_alternatives_release
zf_alts.h, 92

zf_alternatives_send
zf_alts.h, 92

zf_alts.h, 89
zf_alternatives_alloc, 90
zf_alternatives_cancel, 90
zf_alternatives_free_space, 91
zf_alternatives_query_overhead_tcp, 91
zf_alternatives_release, 92
zf_alternatives_send, 92
zft_alternatives_queue, 93

zf_attr, 63

zf_attr_alloc
attr.h, 76

zf_attr_doc
attr.h, 76

zf_attr_dup
attr.h, 76

zf_attr_free
attr.h, 77

zf_attr_get_int
attr.h, 77

zf_attr_get_str
attr.h, 77

zf_attr_reset
attr.h, 79

zf_attr_set_from_fmt
attr.h, 79

zf_attr_set_from_str
attr.h, 79

zf_attr_set_int
attr.h, 80

zf_attr_set_str
attr.h, 80

zf_deinit
zf_stack.h, 102

zf_delegated_send_cancel
zf_ds.h, 95

zf_delegated_send_complete
zf_ds.h, 95

zf_delegated_send_prepare
zf_ds.h, 96

zf_delegated_send_rc
zf_ds.h, 95

zf_delegated_send_tcp_advance
zf_ds.h, 97

zf_delegated_send_tcp_update
zf_ds.h, 97

zf_ds, 64
cong_wnd, 64
delegated_wnd, 64
headers, 65
headers_len, 65
headers_size, 65
ip_len_offset, 65
ip_tcp_hdr_len, 65
mss, 66
reserved, 66
send_wnd, 66
tcp_seq_offset, 66

zf_ds.h, 94
zf_delegated_send_cancel, 95
zf_delegated_send_complete, 95
zf_delegated_send_prepare, 96
zf_delegated_send_rc, 95
zf_delegated_send_tcp_advance, 97
zf_delegated_send_tcp_update, 97

zf_init
zf_stack.h, 102

zf_muxer_add
muxer.h, 82

zf_muxer_alloc
muxer.h, 82

zf_muxer_del
muxer.h, 83

zf_muxer_free
muxer.h, 83

zf_muxer_mod
muxer.h, 84

zf_muxer_set, 67

zf_muxer_wait
muxer.h, 84

zf_pkt_report, 67
bytes, 67
flags, 68
start, 68
timestamp, 68

zf_platform.h, 98

zf_reactor.h, 98
zf_reactor_perform, 99
zf_reactor_perform_attr, 99
zf_stack_has_pending_events, 100
zf_stack_has_pending_work, 100

zf_reactor_perform
zf_reactor.h, 99

zf_reactor_perform_attr

- zf_reactor.h, 99
- zf_stack, 68
- zf_stack.h, 101
 - EPOLLSTACKHUP, 102
 - zf_deinit, 102
 - zf_init, 102
 - zf_stack_alloc, 102
 - zf_stack_free, 103
 - zf_stack_is_quiescent, 103
 - zf_stack_to_waitable, 104
- zf_stack_alloc
 - zf_stack.h, 102
- zf_stack_free
 - zf_stack.h, 103
- zf_stack_has_pending_events
 - zf_reactor.h, 100
- zf_stack_has_pending_work
 - zf_reactor.h, 100
- zf_stack_is_quiescent
 - zf_stack.h, 103
- zf_stack_to_waitable
 - zf_stack.h, 104
- zf_tcp.h, 105
 - zft_addr_bind, 106
 - zft_alloc, 107
 - zft_connect, 107
 - zft_error, 108
 - zft_free, 109
 - zft_get_header_size, 109
 - zft_get_mss, 109
 - zft_get_tx_timestamps, 110
 - zft_getname, 110
 - zft_handle_free, 111
 - zft_handle_getname, 111
 - zft_pkt_get_timestamp, 111
 - zft_rcv, 112
 - zft_send, 113
 - zft_send_single, 114
 - zft_send_single_warm, 115
 - zft_send_space, 116
 - zft_shutdown_tx, 116
 - zft_state, 117
 - zft_to_waitable, 117
 - zft_zc_rcv, 117
 - zft_zc_rcv_done, 118
 - zft_zc_rcv_done_some, 118
- zftl_accept, 119
- zftl_free, 119
- zftl_getname, 120
- zftl_listen, 120
- zftl_to_waitable, 121
- zf_udp.h, 121
 - zfur_addr_bind, 123
 - zfur_addr_unbind, 123
 - zfur_alloc, 124
 - zfur_free, 124
 - zfur_pkt_get_header, 125
 - zfur_pkt_get_timestamp, 125
 - zfur_to_waitable, 126
 - zfur_zc_rcv, 126
 - zfur_zc_rcv_done, 127
 - zfut_alloc, 127
 - zfut_free, 128
 - zfut_get_header_size, 128
 - zfut_get_mss, 129
 - zfut_get_tx_timestamps, 129
 - zfut_send, 129
 - zfut_send_single, 130
 - zfut_send_single_warm, 130
 - zfut_to_waitable, 131
- zf_waitable, 69
- zf_waitable_event
 - muxer.h, 85
- zf_waitable_fd_get
 - muxer.h, 85
- zf_waitable_fd_prime
 - muxer.h, 86
- ZF_DELEGATED_SEND_RC_BAD_SOCKET
 - zf_ds.h, 95
- ZF_DELEGATED_SEND_RC_NOARP
 - zf_ds.h, 95
- ZF_DELEGATED_SEND_RC_NOCWIN
 - zf_ds.h, 95
- ZF_DELEGATED_SEND_RC_NOWIN
 - zf_ds.h, 95
- ZF_DELEGATED_SEND_RC_OK
 - zf_ds.h, 95
- ZF_DELEGATED_SEND_RC_SENDQ_BUSY
 - zf_ds.h, 95
- ZF_DELEGATED_SEND_RC_SMALL_HEADER
 - zf_ds.h, 95
- zf_ds.h
 - ZF_DELEGATED_SEND_RC_BAD_SOCKET, 95
 - ZF_DELEGATED_SEND_RC_NOARP, 95
 - ZF_DELEGATED_SEND_RC_NOCWIN, 95
 - ZF_DELEGATED_SEND_RC_NOWIN, 95
 - ZF_DELEGATED_SEND_RC_OK, 95
 - ZF_DELEGATED_SEND_RC_SENDQ_BUSY, 95
 - ZF_DELEGATED_SEND_RC_SMALL_HEADER, 95
- zft, 69
- zft_addr_bind
 - zf_tcp.h, 106
- zft_alloc
 - zf_tcp.h, 107
- zft_alternatives_queue
 - zf_alts.h, 93
- zft_connect
 - zf_tcp.h, 107
- zft_error
 - zf_tcp.h, 108
- zft_free

[zf_tcp.h](#), [109](#)
[zft_get_header_size](#)
 [zf_tcp.h](#), [109](#)
[zft_get_mss](#)
 [zf_tcp.h](#), [109](#)
[zft_get_tx_timestamps](#)
 [zf_tcp.h](#), [110](#)
[zft_getname](#)
 [zf_tcp.h](#), [110](#)
[zft_handle](#), [70](#)
[zft_handle_free](#)
 [zf_tcp.h](#), [111](#)
[zft_handle_getname](#)
 [zf_tcp.h](#), [111](#)
[zft_msg](#), [70](#)
 [flags](#), [71](#)
 [iov](#), [71](#)
 [iocnt](#), [71](#)
 [pkts_left](#), [71](#)
 [reserved](#), [71](#)
[zft_pkt_get_timestamp](#)
 [zf_tcp.h](#), [111](#)
[zft_rcv](#)
 [zf_tcp.h](#), [112](#)
[zft_send](#)
 [zf_tcp.h](#), [113](#)
[zft_send_single](#)
 [zf_tcp.h](#), [114](#)
[zft_send_single_warm](#)
 [zf_tcp.h](#), [115](#)
[zft_send_space](#)
 [zf_tcp.h](#), [116](#)
[zft_shutdown_tx](#)
 [zf_tcp.h](#), [116](#)
[zft_state](#)
 [zf_tcp.h](#), [117](#)
[zft_to_waitable](#)
 [zf_tcp.h](#), [117](#)
[zft_zc_rcv](#)
 [zf_tcp.h](#), [117](#)
[zft_zc_rcv_done](#)
 [zf_tcp.h](#), [118](#)
[zft_zc_rcv_done_some](#)
 [zf_tcp.h](#), [118](#)
[zftl](#), [72](#)
[zftl_accept](#)
 [zf_tcp.h](#), [119](#)
[zftl_free](#)
 [zf_tcp.h](#), [119](#)
[zftl_getname](#)
 [zf_tcp.h](#), [120](#)
[zftl_listen](#)
 [zf_tcp.h](#), [120](#)
[zftl_to_waitable](#)
 [zf_tcp.h](#), [121](#)
[zfur](#), [72](#)
[zfur_addr_bind](#)
 [zf_udp.h](#), [123](#)
[zfur_addr_unbind](#)
 [zf_udp.h](#), [123](#)
[zfur_alloc](#)
 [zf_udp.h](#), [124](#)
[zfur_free](#)
 [zf_udp.h](#), [124](#)
[zfur_msg](#), [73](#)
 [dgrams_left](#), [73](#)
 [flags](#), [73](#)
 [iov](#), [73](#)
 [iocnt](#), [74](#)
 [reserved](#), [74](#)
[zfur_pkt_get_header](#)
 [zf_udp.h](#), [125](#)
[zfur_pkt_get_timestamp](#)
 [zf_udp.h](#), [125](#)
[zfur_to_waitable](#)
 [zf_udp.h](#), [126](#)
[zfur_zc_rcv](#)
 [zf_udp.h](#), [126](#)
[zfur_zc_rcv_done](#)
 [zf_udp.h](#), [127](#)
[zfut](#), [74](#)
[zfut_alloc](#)
 [zf_udp.h](#), [127](#)
[zfut_free](#)
 [zf_udp.h](#), [128](#)
[zfut_get_header_size](#)
 [zf_udp.h](#), [128](#)
[zfut_get_mss](#)
 [zf_udp.h](#), [129](#)
[zfut_get_tx_timestamps](#)
 [zf_udp.h](#), [129](#)
[zfut_send](#)
 [zf_udp.h](#), [129](#)
[zfut_send_single](#)
 [zf_udp.h](#), [130](#)
[zfut_send_single_warm](#)
 [zf_udp.h](#), [130](#)
[zfut_to_waitable](#)
 [zf_udp.h](#), [131](#)