

# L4 Nucleus Version X Reference Manual

x86

Version X.0

Jochen Liedtke  
Universität Karlsruhe  
liedtke@ira.uka.de

September 3, 1999  
Under Construction



## How To Read This Manual

This reference manual consists of two parts, (1) a processor-independent description of the principles and mechanisms of L4 and (2) a more detailed processor-specific description. Part 2 refers to the Intel processors, Pentium<sup>®</sup><sup>1</sup> and Pentium<sup>®</sup>Pro.

## Credits

Helpful contributions for improving this reference manual and the L4 interface came from many persons, in particular from Bryan Ford, Hermann Härtig, Michael Hohmuth, Sebastian Schönberg, Jean Wolter, Kevin Elphinstone, Trent Jaeger, Yoonho Park, Volkmar Uhlig, and Gernot Heiser.

---

<sup>1</sup>Pentium<sup>®</sup> is a registered trademark of Intel Corp.



# Contents

<b>1</b>	<b>L4 in General</b>	<b>7</b>
<b>2</b>	<b>L4-X/x86</b>	<b>9</b>
2.1	Notational conventions . . . . .	9
2.2	Data Types . . . . .	10
2.2.1	Unique Ids . . . . .	10
2.2.2	Fpages . . . . .	10
2.2.3	Messages . . . . .	11
2.2.4	Timeouts . . . . .	13
2.3	L4 Calls . . . . .	15
	ipc . . . . .	16
	id_nearest . . . . .	23
	fpage_unmap . . . . .	24
	thread_switch . . . . .	25
	thread_schedule . . . . .	26
	lthread_ex_regs . . . . .	28
	task_new . . . . .	30
2.4	Processor Mirroring . . . . .	32
2.4.1	Segments . . . . .	32
2.4.2	Exception Handling . . . . .	32
2.4.3	Debug Registers . . . . .	32
2.5	The Kernel-Info Page . . . . .	33
2.6	Page-Fault and Preemption RPC . . . . .	34
2.7	$\sigma_0$ RPC protocol . . . . .	35
2.8	Starting L4 . . . . .	37
<b>A</b>	<b>Booting</b>	<b>39</b>



# Chapter 1

## L4 in General

Chapter omitted in this release.





# Chapter 2

## L4-X/x86

The L4/x86 Nucleus runs on x86-type processor's that offer at least the functionality of an Intel 486 processor, e.g. 486, Pentium, PentiumPro, Pentium II, and Pentium III.

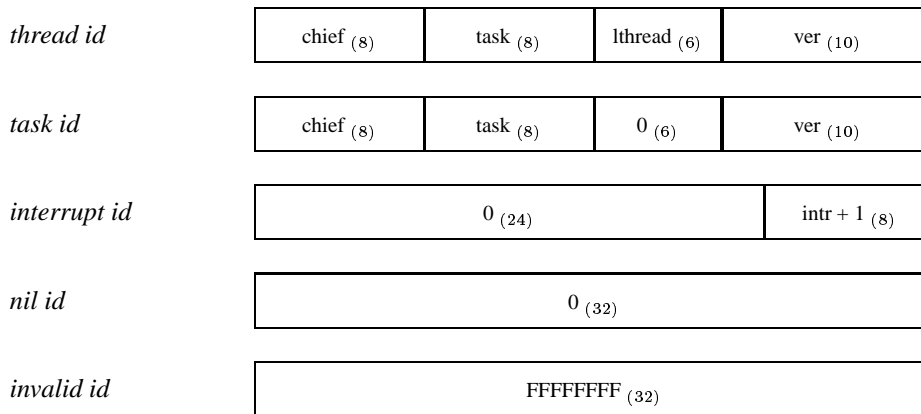
### 2.1 Notational conventions

- ~ If this refers to an input parameter, its value is meaningless. If it refers to an output parameter, its value is undefined.
- EAX,ECX... denote the processor's general registers.
- $\langle SP+n \rangle$  denotes the word on the user stack addressed by  $SP+n$ , where  $SP$  represents the user-level stack pointer.

## 2.2 Data Types

### 2.2.1 Unique Ids

Unique ids identify tasks, threads and hardware interrupts. Each unique id is a 32-bit value which is unique in time. An unique id in x86 format consists of one 32-bit word:

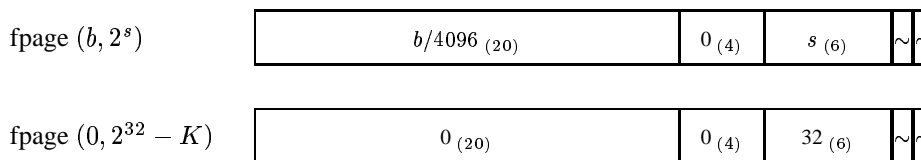


available tasks:	255	1–255	task 0 does not exist
available threads per task:	64	0–63	
usable version numbers:	1024	0–1023	

### 2.2.2 Fpages

Fpages (Flexpages) are regions of the virtual address space. An fpage consists of all pages actually mapped in this region. The minimal fpage size is 4 K, the minimal hardware-page size.

An fpage of size  $2^s$  has a  $2^s$ -aligned base address  $b$ , i.e.  $b \bmod 2^s = 0$ . On the x86 processors, the smallest possible value for  $s$  is 12, since hardware pages are at least 4K. The complete user address space (base address 0, size  $2^{32} - K$ , where  $K$  is the size of the kernel area) is denoted by  $b = 0, s = 32$ . An fpage with base address  $b$  and size  $2^s$  is denoted by a 32-bit word:



### IO-Ports

x86 IO-ports form a separate address space besides the conventional memory address space. Its size is 64 K and its granularity is 16 bytes. However, IO-ports can only be mapped idempotently, i.e. physical port  $x$  is either mapped at the address  $x$  in the task's IO address space or it is not mapped.

L4 handles IO-ports like memory, i.e. as fpages. IO-fpages can be mapped, granted and unmapped like memory fpages. However, since IO-ports can only mapped idempotent, always the complete IO space (64 K) should be specified as receive fpage.

An IO-fpage of size  $2^s$  ( $4 \leq s \leq 16$ ) has a  $2^s$ -aligned base address  $p$ , i.e.  $p \bmod 2^s = 0$ . An fpage with base port address  $p$  and size  $2^s$  is denoted by a 32-bit word:

IO-fpage ( $p, 2^s$ )	F0 <sub>(8)</sub>	$p$ <sub>(12)</sub>	0 <sub>(4)</sub>	$s$ <sub>(6)</sub>	~	~
-----------------------	-------------------	---------------------	------------------	--------------------	---	---

IO-fpage (0, $2^{16}$ )	F0 <sub>(8)</sub>	0 <sub>(12)</sub>	0 <sub>(4)</sub>	16 <sub>(6)</sub>	~	~
-------------------------	-------------------	-------------------	------------------	-------------------	---	---

### 2.2.3 Messages

#### Register Messages/Buffers

Register messages consist of up to 3 words of 32 bits. Upon sending, the message is located in the registers EDX, EBX, and EDI. Upon receiving, the same registers serve as a buffer, i.e. the registers EDX, EBX, and EDI contain the received message (where send EDX is received EDX, etc.).

dword 2 <sub>(32)</sub>	EDI
dword 1 <sub>(32)</sub>	EBX
dword 0 <sub>(32)</sub>	EDX

#### Simple Memory Messages/Buffers

If messages are longer than 3 dwords, memory messages have to be used. Such messages consist of a message descriptor (dope) that, e.g., specifies the current message length in dwords (msg snd dope) for messages to be sent.

A message receive buffer is structured similarly but the size of the receive buffer, i.e. the maximum message length than can be received in this buffer, is specified by the msg size dope, also in dwords. After a message was received, the *msg snd dope* contains the current length of the received message. *msg size dope* is unchanged.

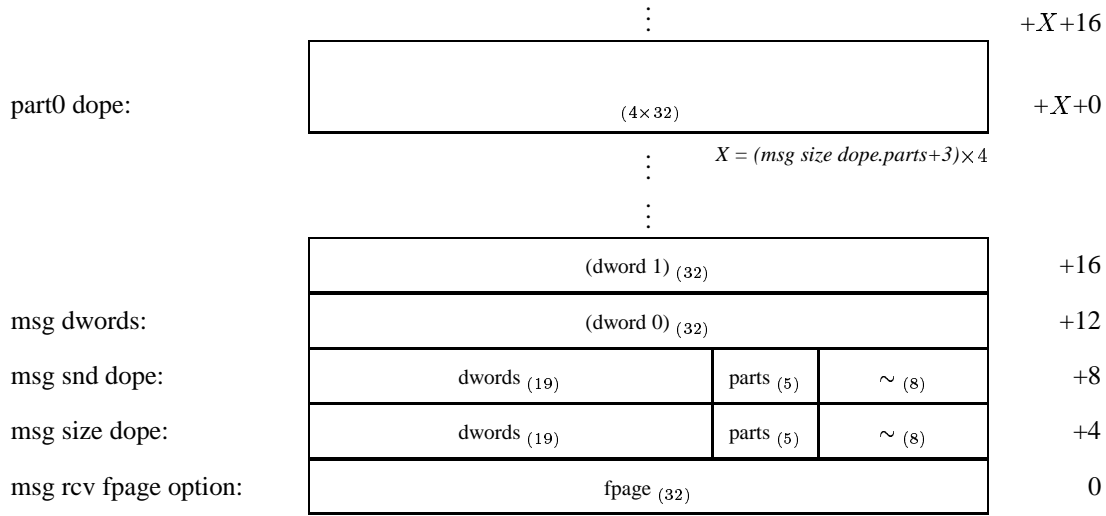
	⋮			
	dword 3 <sub>(32)</sub>		+32	
	(dword 2 — not transferred) <sub>(32)</sub>		+24	
	(dword 1 — not transferred) <sub>(32)</sub>		+16	
msg dwords:	(dword 0 — not transferred) <sub>(32)</sub>		+12	
msg snd dope:	dwords <sub>(19)</sub>	0 <sub>(5)</sub>	~ <sub>(8)</sub>	+8
msg size dope:	dwords <sub>(19)</sub>	0 <sub>(5)</sub>	~ <sub>(8)</sub>	+4
msg rcv fpage option:	fpage <sub>(32)</sub>			0

Using different send and size dopes permits to specify not only pure send messages and pure receive message buffers. It is as well possible to send a message and receive the reply using the same data structure. And, you can also receive a message in a message buffer and then forward this buffer as a message without changing the data structure.

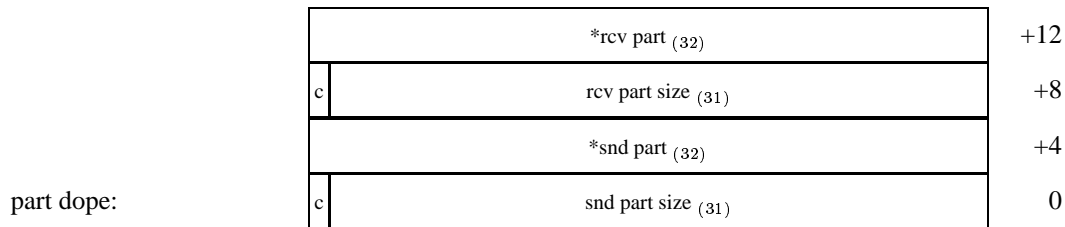
A further simplification for user-level programming is that *the first three dwords of a message are always transferred via register*. This permits to handle long (memory) and short (register) messages basically in the same way (the first three dwords are always in registers). Note that loading/storing those registers from/to the message/buffer data structure is *not* handled by the Nucleus. It can be done at user level.

## Memory Messages/Buffers With Indirect Parts

Complex messages or buffers can contain up to 31 indirect parts. The number of those parts is specified in the *parts* fields. For buffers, the *msg size dope* defines the maximum number of accepted indirect parts for receiving. For messages, *msg snd dope* defines the current indirect parts for sending, and the *dwords* field of *msg size dope* defines the position *X* where the first indirect part is located. Therefore, the *dwords* field of the *msg size dope* must be set correctly, even if it is a pure send message.



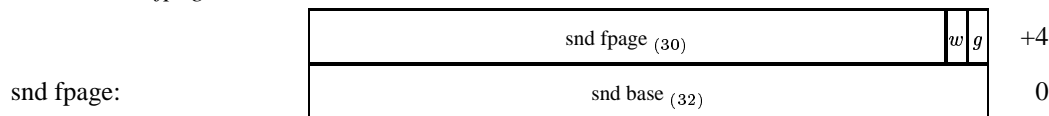
Every indirect part dope can specify an indirect send part and an indirect receive buffer. For a send message, the receive buffer part is ignored. For a receive buffer, the receive buffer specifies the buffer address and maximum length; after a message was received, the send part address is set to the beginning of the receive buffer and the send part length specifies the current length. Thus, received parts can be forwarded by a successive send operation without any change.



The *c* bits enable *scatter/gather* functionality. *c* = 0 specifies the begin of a logical string; *c* = 1 specifies that this is a continuation of the last logical string. On the sender side, *logical string* means a sequence of one or more parts that are transferred as if they were one contiguous string (gather). On the receiver side, *logical string* means a sequence of one or more buffers that are treated as one logical buffer; the corresponding received string is scattered among them. Continuations can be arbitrarily combined on sender and receiver side. Note that length and size fields are always per part.

## Map Messages/Buffers

Map messages are differentiated from copy messages in the *ipc* system call. The format of a map message is like a register message or memory message where the direct message part (not the optional indirect parts) consists of *snd fpages*:



$w=0$	The fpage will be mapped/granted read only.
$w=1$	The fpage will be mapped/granted with the full access right the mapper possesses.
$g=0$	The fpage will be mapped.
$g=1$	The fpage will be granted.

## 2.2.4 Timeouts

Timeouts are used to control ipc operations. The *send timeout* determines how long ipc should try to send a message. If the specified period is exhausted without that message transfer could start, ipc fails. The *receive timeout* specifies how long ipc should wait for an incoming message. Both timeouts specify the maximum period of time *before message transfer starts*. Once started, message transfer is no longer influenced by send or receive timeout.

Pagefaults occurring during ipc are controlled by *send* and *receive pagefault timeout*. A pagefault is translated to an RPC by the kernel. In the case of a pagefault in the receiver's address space, the corresponding RPC to the pager uses *send pagefault timeout* (specified by the sender) for both send and receive timeout. In the case of a pagefault in the sender's address space, *receive pagefault timeout* specified by the receiver is taken.

Besides the special timeouts 0 (do not wait at all) and  $\infty$  (wait forever), periods from 1  $\mu$ s up to approximately 19 hours can be specified. The complete quadruple is packed into one 32-bit word:

$m_r$ (8)	$m_s$ (8)	$p_s$ (4)	$p_r$ (4)	$e_s$ (4)	$e_r$ (4)
-----------	-----------	-----------	-----------	-----------	-----------

Note that for efficiency reasons the highest bit of any mantissa  $m$  must be 1, except for  $m=0$ .

$$\text{snd timeout} = \begin{cases} \infty & \text{if } e_s=0 \\ 4^{15-e_s} m_s \mu s & \text{if } e_s > 0 \\ 0 & \text{if } m_s=0, e_s \neq 0 \end{cases}$$

$$\text{rcv timeout} = \begin{cases} \infty & \text{if } e_r=0 \\ 4^{15-e_r} m_r \mu s & \text{if } e_r > 0 \\ 0 & \text{if } m_r=0, e_r \neq 0 \end{cases}$$

$$\text{snd pagefault timeout} = \begin{cases} \infty & \text{if } p_s=0 \\ 4^{15-p_s} \mu s & \text{if } 0 < p_s < 15 \\ 0 & \text{if } p_s=15 \end{cases}$$

$$\text{rcv pagefault timeout} = \begin{cases} \infty & \text{if } p_r=0 \\ 4^{16-p_r} \mu s & \text{if } 0 < p_r < 15 \\ 0 & \text{if } p_r=15 \end{cases}$$

approximate timeout ranges		
$e_s, e_r, p_s, p_r$	snd/rcv timeout	pf timeout
0	$\infty$	$\infty$
1	256 s ... 19 h	256 s
2	64 s ... 55 h	64 s
3	16 s ... 71 m	16 s
4	4 s ... 17 m	4 s
5	1 s ... 4 m	1 s
6	262 ms ... 67 s	256 ms
7	65 ms ... 17 s	64 ms
8	16 ms ... 4 s	16 ms
9	4 ms ... 1 s	4 ms
10	1 ms ... 261 ms	1 ms
11	256 $\mu$ s ... 65 ms	256 $\mu$ s
12	64 $\mu$ s ... 16 ms	64 $\mu$ s
13	16 $\mu$ s ... 4 ms	16 $\mu$ s
14	4 $\mu$ s ... 1 ms	4 $\mu$ s
15	1 $\mu$ s ... 255 $\mu$ s	0
$m=0, e>0$	0	—

## 2.3 L4 Calls

This section describes the 7 system calls of L4:

- ipc int 30
- id\_nearest int 31
- fpage\_unmap int 32
- thread\_switch int 33
- thread\_schedule int 34
- lthread\_ex\_regs int 35
- task\_new int 36

# ipc

<i>snd descriptor</i>	EAX	— INT 0x30 →	EAX	<i>msg.dope</i> + <i>cc</i> / <i>cc</i>
<i>timeouts</i>	ECX		ECX	~ / ~
<i>msg.w0</i>	EDX		EDX	<i>msg.w0</i> / ~
<i>msg.w1</i>	EBX		EBX	<i>msg.w1</i> / ~
<i>rcv descriptor</i>	EBP		EBP	~ / ~
<i>dest</i>	ESI		ESI	<i>source</i> / ~
<i>msg.w2</i>	EDI		EDI	<i>msg.w2</i> / ~
<i>propagatee</i> / ~	<SP>			

This is the basic system call for inter-process communication and synchronization. It may be used for intra- as inter-address-space communication. All communication is synchronous and unbuffered: a message is transferred from the sender to the recipient if and only if the recipient has invoked a corresponding ipc operation. The sender blocks until this happens or a period specified by the sender elapsed without that the destination became ready to receive.

Ipc can be used to copy data as well as to *map* or *grant* fpages from the sender to the recipient. For the description of messages see section 2.2.3.

12-byte messages (plus 32-bit sender id) can be transferred solely via the registers and are thus specially optimized. If possible, short messages should therefore be reduced to 12-byte messages.

A single ipc call combines an optional send operation with an optional receive operation. Whether it includes a send respectively a receive is determined by the actual parameters. If the send or receive address is specified as *nil* (0xFFFFFFFF), the corresponding operation is skipped.

No time is required for the transition between send and receive phase of one ipc operation.

## Parameters

<i>snd descriptor</i>	“ <i>nil</i> ”	<div style="border: 1px solid black; padding: 5px; text-align: center;">0xFFFFFFFF (32)</div>	Ipc does not include a send operation.
	“ <i>mem</i> ”	<div style="border: 1px solid black; padding: 5px; text-align: center;">*snd msg/4 (30) <span style="border: 1px solid black; padding: 0 2px;"><i>m</i></span> <span style="border: 1px solid black; padding: 0 2px;"><i>p</i></span></div>	Ipc includes sending a message to the destination specified by <i>dest id</i> . *snd msg must point to a valid message. The first two 32-bit words of the message ( <i>msg.w0</i> and <i>msg.w1</i> ) are <i>not</i> taken from the message data structure but must be contained in registers EDX and EBX.
	“ <i>reg</i> ”	<div style="border: 1px solid black; padding: 5px; text-align: center;">0 (30) <span style="border: 1px solid black; padding: 0 2px;">0</span> <span style="border: 1px solid black; padding: 0 2px;"><i>p</i></span></div>	Ipc includes sending a message to the destination specified by <i>dest id</i> . The message consists solely of the two 32-bit words <i>msg.w0</i> and <i>msg.w1</i> in registers EDX and EBX.
<i>m=0</i>			Value-copying send operation; the dwords or the message are simply copied to the recipient.



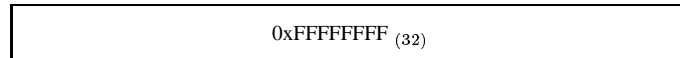
*snd descriptor*  $m=1$  Fpage-mapping send operation. The dwords of the message to be sent are treated as 'send fpages'. The described fpages are mapped (respectively granted) into the recipient's address space. Mapping/granting stops when either the end of the dwords is reached or when an invalid fpage denoter is found, in particular 0. The send fpage descriptors and all potentially following words are also transferred by simple copy to the recipient. Thus a message may contain some fpages and additional value parameters. The recipient can use the received fpage descriptors to determine what has been mapped or granted into its address space, including location and access rights.

$p=0$  Normal (unpropagated) send operation. The recipient gets the original sender's id.

$p=1$  Propagating send operation. The additional *propagatee* parameter on the stack specifies the id of the propagatee thread. If propagatee and current sender or propagatee and receiver belong to the same task, propagation is always permitted. Otherwise, the current sender is supposed to be a chief that uses deceiving and it is checked whether it is direction preserving. If not,  $p$  and the *propagatee* parameter are ignored. If propagation (or deceiving) is permitted, the receiver receives the propagatee's id instead of the current sender's id. If the propagatee is waiting (closed) for a reply from the current sender, the propagatee's status is additionally modified such that the propagatee now waits for a the new receiver instead of the current sender. The *propagatee* parameter on the user stack is only required if  $p=1$ .

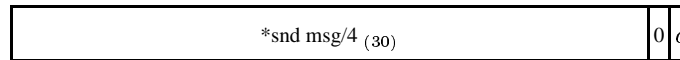
*rcv descriptor*

"nil"



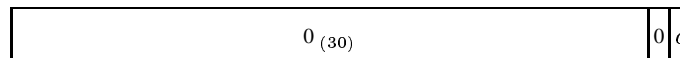
Ipc does not include a receive operation.

"mem"



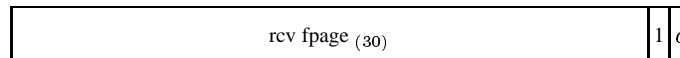
Ipc includes receiving a message respectively waiting to receive a message. \*rcv msg must point to a valid message. The first two 32-bit words of the received message (*msg.w0* and *msg.w1*) are *not* stored in the message data structure but are returned in registers EDX and EBX.

"reg"



Ipc includes receiving a message respectively waiting to receive a message. However, only messages up to two 32-bit words *msg.w0* and *msg.w1* are accepted. The received message is returned in registers EDX and EBX.

"rmap"



Ipc includes receiving a message respectively waiting to receive a message. However, only send-fpage messages or up to two 32-bit words *msg.w0* and *msg.w1* are accepted. The received message is returned in registers EDX and EBX. If a map message is received, "rcv fpage" describes the receive fpage (instead of "rcv fpage option" in a memory message buffer). Thus fpages can also be received without a message buffer in memory.

*rcv\_descriptor* *o*=0 Only messages from the thread specified as *dest id* are accepted (“closed wait”). Any send operation from a different thread (or hardware interrupt) will be handled exactly as if the actual thread would be busy.

*o*=1 Messages from any thread will be accepted (“open wait”). If the actual thread is associated with a hardware interrupt, also messages from this hardware interrupt can arrive.

*dest id* *≠nil* Sending is directed to the specified thread, if it resides in the sender’s clan. If the destination is outside the sender’s clan, the message is sent to the sender’s chief. If the destination is in an inner clan (a clan whose chief resides in the sender’s clan), it is redirected to that chief. (See also ‘nchief’ operation.) If no send part was specified (*snd\_descriptor = nil*), *dest id* specifies the source from which messages can be received. (However recall that the receive restriction is only effective if *o* = 0.)

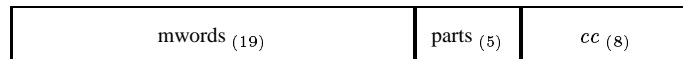
*=nil* (*nil*=0) Although specifying *nil* as the destination for a send operation is illegal (error: ‘destination not existent’), it can be legally specified for a receive-only operation. In this case, ipc will not receive any message but will wait the specified *rcv timeout* and then terminate with error code ‘receive timeout’. (However recall that the receive restriction is only effective if *o* = 0.)

*source id* If a message was received this is the id of its sender. (If a hardware interrupt was received this is the interrupt id.) The parameter is undefined if no message was received.

*msg.w0, w1, w2* “*snd*” First three 32-bit words of message to be sent. These message words are taken directly from registers EDX, EBX, and EDI. *They are not read from the message data structure.*

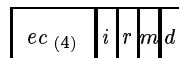
“*rcv*” First three 32-bit words of received message, undefined if no message was received. *These message words are available only in registers EDX, EBX, and EDI.* The Nucleus does not store it in the receive message buffer. The user program may store it or use it directly in the registers.

*msg.dope + cc*



Message dope describing received message. If no message was received, only *cc* is delivered. *The dope word of the received message is available only in register EAX.* The Nucleus does not store it in the receive message buffer. The user program may store it or use it directly in the register. (Note that the lowermost 8 bits of msg dope and size dope in the message data structure are undefined. So it is legal to store EAX in the msg-dope field, even if *cc*≠0.)

*cc*



*d*=0 The received message is transferred directly (“undeceived”) from *source id*.

*d*=1 The received message is “deceived” by a chief. *source id* is the virtual source id which was specified by the sending chief.

*m*=0 The received message did not contain fpages.

*m*=1 The sender mapped or granted fpages. The sender’s fpage descriptors were also (besides mapping/granting) transferred as mwords.

*r*=0 The received message was directed to the actual recipient, not redirected to a chief. I.e. sender and receiver a part of the same clan. The *i*-bit has no meaning in this case and is zero.

*r*=1 The received message was redirected to the chief which was next on the path to the true destination. Sender and addressed recipient belong to different clans.

*i*=0 If *r*=1: the received message comes from outside the own clan.

*i*=1 If *r*=1: the received message comes from an inner clan.

*ec*

- = 0 *ok*: the optional send operation was successful, and if a receive operation was also specified (*rcv\_descriptor ≠ nil*) a message was also received correctly.
- ≠ 0 If ipc fails the completion code is in the range 0x10...0xF0. If the send operation already failed, ipc is terminated without the potentially specified receive operation. *s* specifies whether the error occurred during the receive (*s* = 0) operation or during the send (*s* = 1) operation:
  - 1 *Non-existing* destination or source.
  - 2 + *s* *Timeout*.
  - 4 + *s* *Cancelled* by another thread (system call *lthread\_ex\_regs*).
  - 6 + *s* *Map failed* due to a shortage of page tables.
  - 8 + *s* *Send pagefault timeout*.
  - A + *s* *Receive pagefault timeout*.
  - C + *s* *Aborted* by another thread (system call *lthread\_ex\_regs*).
  - E + *s* *Cut* message. Potential reasons are (a) the recipient's mword buffer is too small; (b) the recipient does not accept enough parts; (c) at least one of the recipient's part buffers is too small.
- 1...5 The according operation was terminated before a real message transfer started. No partner was directly involved.
- 6...F The according operation was terminated while a message transfer was running. The message transfer was aborted. The current partner (sender or receiver) was involved and got the corresponding error code. It is not defined which parts of the message are already transferred and which parts are not yet transferred.

*timeouts*

This 32-bit word specifies all 4 timeouts, the quadruple (*snd*, *rcv*, *snd pf*, *rcv pf*). For A detailed description see section 2.2.4. Frequently used values are

	snd	rcv	snd pf	rcv pf
0x00000000	∞	∞	∞	∞
0x00000001	0	∞	∞	∞
0x00000011	0	0	∞	∞

- “*snd*” If the required send operation cannot start transfer data within the specified time, ipc is terminated and fails with completion code ‘send timeout’ (0x18). If ipc does not include a send operation, this parameter is meaningless.
- “*rcv*” If ipc includes a receive operation and no message transfer starts within the specified time, ipc is terminated and fails with completion code ‘receive timeout’ (0xA0). If ipc does not include a receive operation, this parameter is meaningless.
- “*spf*” If during sending data a pagefault *in the receiver's address space* occurs, *snd pf* specified by the sender is used as send and receive timeout for the pagefault RPC.
- “*rcpf*” If during receiving data a pagefault *in the sender's address space* occurs, *rcv pf* specified by the receiver is used as send and receive timeout for the pagefault RPC.

## Basic Ipc Types

CALL	<i>*snd msg / 0</i>	EAX		- INT 0x30 →	EAX	<i>msg.dope + cc</i>
	<i>timeouts</i>	ECX			ECX	~
	<i>msg.w0</i>	EDX			EDX	<i>msg.w0</i>
	<i>msg.w1</i>	EBX			EBX	<i>msg.w1</i>
	<i>*rcv msg / 0</i>	EBP			EBP	~
	<i>dest id</i>	ESI			ESI	<i>unchanged</i>
	<i>msg.w2</i>	EDI			EDI	<i>msg.w2</i>

This is the usual blocking RPC. *snd msg* is sent to *dest id* and the invoker waits for a reply from *dest id*. Messages from other sources are not accepted. Note that since the send/receive transition needs no time, the destination can reply with *snd timeout = 0*.

This operation can also be used for a server with one dedicated client. It sends the reply to the client and waits for the client's next order.

SEND/RECEIVE	<i>*snd msg / 0</i>	EAX		- INT 0x30 →	EAX	<i>msg.dope + cc</i>
	<i>timeouts</i>	ECX			ECX	~
	<i>msg.w0</i>	EDX			EDX	<i>msg.w0</i>
	<i>msg.w1</i>	EBX			EBX	<i>msg.w1</i>
	<i>*rcv msg + 1 / 0 + 1</i>	EBP			EBP	~
	<i>dest id</i>	ESI			ESI	<i>source id</i>
	<i>msg.w2</i>	EDI			EDI	<i>msg.w2</i>

*snd msg* is sent to *dest id* and the invoker waits for a reply from any source. This is the standard server operation: it sends a reply to the actual client and waits for the next order which may come from a different client.

SEND	<i>*snd msg / 0</i>	EAX		- INT 0x30 →	EAX	<i>cc</i>
	<i>timeouts</i>	ECX			ECX	~
	<i>msg.w0</i>	EDX			EDX	~
	<i>msg.w1</i>	EBX			EBX	~
	<i>0xFFFFFFFF</i>	EBP			EBP	~
	<i>dest id</i>	ESI			ESI	~
	<i>msg.w2</i>	EDI			EDI	~

*snd msg* is sent to *dest id*. There is no receive phase included. The invoker continues working after sending the message.

RECEIVE FROM	<i>0xFFFFFFFF</i>	EAX		- INT 0x30 →	EAX	<i>msg.dope + cc</i>
	<i>timeouts</i>	ECX			ECX	~
	~	EDX			EDX	<i>msg.w0</i>
	~	EBX			EBX	<i>msg.w1</i>
	<i>*rcv msg / 0</i>	EBP			EBP	~
	<i>dest id</i>	ESI			ESI	<i>unchanged</i>
	~	EDI			EDI	<i>msg.w2</i>

This operation includes no send phase. The invoker waits for a message from *source id*. Messages from other

sources are not accepted. Note that also a hardware interrupt might be specified as source.

RECEIVE	<i>0xFFFFFFFF</i>	EAX		- INT 0x30 →		EAX	<i>msg.dope + cc</i>
	<i>timeouts</i>	ECX				ECX	~
	~	EDX				EDX	<i>msg.w0</i>
	~	EBX				EBX	<i>msg.w1</i>
	<i>*rcv msg+1 / 0+1</i>	EBP				EBP	~
	~	ESI				ESI	<i>source id</i>
	~	EDI				EDI	<i>msg.w2</i>

This operation includes no send phase. The invoker waits for a message from any source (including a hardware interrupt).

RECEIVE INTR	<i>0xFFFFFFFF</i>	EAX		- INT 0x30 →		EAX	<i>msg.dope + cc</i>
	<i>timeouts</i>	ECX				ECX	~
	~	EDX				EDX	~
	~	EBX				EBX	~
	<i>*rcv msg / 0</i>	EBP				EBP	~
	<i>intr + 1</i>	ESI				ESI	<i>unchanged</i>
	~	EDI				EDI	~

This operation includes no send phase. The invoker waits for an interrupt message coming from interrupt source *intr*. Note that interrupt messages come *only* from the interrupt which is currently associated with this thread.

The *intr* parameter is only evaluated if *rcv timeout = 0* is specified, see ‘associate intr’.

ASSOCIATE INTR	<i>0xFFFFFFFF</i>	EAX		- INT 0x30 →		EAX	<i>msg.dope + cc</i>
	<i>rcv timeout = 0</i>	ECX				ECX	~
	~	EDX				EDX	~
	~	EBX				EBX	~
	<i>*rcv msg / 0</i>	EBP				EBP	~
	<i>intr + 1</i>	ESI				ESI	<i>unchanged</i>
	~	EDI				EDI	~

The *intr* parameter is evaluated if *rcv timeout = 0* is specified. If no (currently associated) interrupt is pending, the current thread is (1) detached from its currently associated interrupt (if any) and (2) associated with the specified interrupt provided that this one is free, i.e. not associated with another thread. If the association succeeds, the completion code is *receive timeout* (0x20) and no interrupt is received.

If an interrupt from the currently associated interrupt was pending, this one is delivered together with completion code *ok* (0x00); the interrupt association is *not* modified. If the requested new interrupt is already associated to another thread or is not existing, completion code *non existing* (0x10) is delivered and the interrupt association is not modified.

Getting rid of an associated interrupt without associating a new one is done by issuing a receive from *nilthread* (0) with *rcv timeout = 0*.

SLEEP	<i>0xFFFFFFFF</i>	EAX		- INT 0x30 →		EAX	<i>cc = 0xA0</i>
	<i>timeouts</i>	ECX				ECX	~
	~	EDX				EDX	~
	~	EBX				EBX	~
	0	EBP				EBP	~
	0	ESI				ESI	~
	~	EDI				EDI	~

This operation includes no send phase. Since *nil* (0) is specified as source, no message can arrive and the ipc will be terminated with 'receive timeout' after the time specified by the *rcv-timeout* parameter is elapsed.

## id\_nearest

~	EAX	— INT 0x31 →	EAX	<i>type</i>
~	ECX		ECX	~
~	EDX		EDX	~
~	EBX		EBX	~
~	EBP		EBP	~
<i>dest id</i>	ESI		ESI	<i>nearest id</i>
~	EDI		EDI	~

If *nil* is specified as destination, the system call delivers the uid of the current thread. Otherwise, it delivers the nearest partner which would be engaged when sending a message to the specified destination. If the destination does not belong to the invoker's clan, this call delivers the chief that is nearest to the invoker on the path from the invoker to the destination.

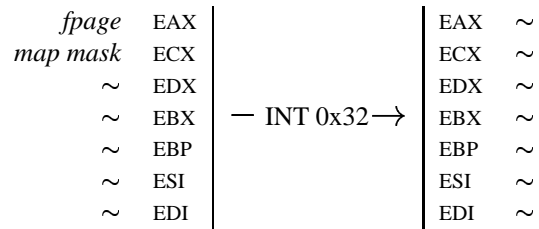
- If the destination resides outside the invoker's clan, it delivers the invoker's own chief.
- If the destination is inside a clan or a clan nesting whose chief *C* is direct member of the invoker's clan, the call delivers *C*.
- If the destination is a direct member of the invoker's clan, the call delivers the destination itself.
- If the destination is *nil*, the call delivers the current thread's id.

Concluding: *nchief* (*dest id*  $\neq$  *nil*) delivers exactly that partner to which the kernel would physically send a message which is targeted to *dest id*. On the other hand, a message from *dest id* would physically come from exactly this partner.

### Parameters

<i>dest id</i>		Id of the destination.
<i>type</i>		Note that the <i>type</i> values correspond exactly to the completion codes of ipc.
	=0	Destination resides in the same clan. <i>dest id</i> is delivered as <i>nearest id</i> .
	= <i>C</i>	Destination is in an inner clan. The chief of this clan or clan nesting is delivered as <i>nearest id</i> .
	=4	Destination is outside the invoker's clan. The invoker's chief is delivered as <i>nearest id</i> .
<i>nearest id</i>		Either the current thread's id or the id of the nearest partner towards <i>dest id</i> .

# fpage\_unmap



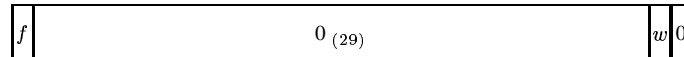
The specified *fpage* is unmapped in all address spaces into which the invoker mapped it directly or indirectly.

## Parameters

*fpage*

Fpage to be unmapped.

*map mask*



- w*=0 Fpage will partially unmapped. Already read/write mapped parts will be set to read only. Read only mapped parts are not affected.
- w*=1 Fpage will be completely unmapped.
- f*=0 Unmapping happens in all address spaces into which pages of the specified fpage have been mapped directly or indirectly. The *original* pages in the own task remain mapped.
- f*=1 Additionally, also the original pages in the own task are unmapped (flushing).





## thread\_schedule

<i>param word</i>	EAX	— INT 0x34 —>	EAX	<i>old param word</i>
~	ECX		ECX	<i>time.low</i>
~	EDX		EDX	<i>time.high</i>
<i>ext preempter</i>	EBX		EBX	<i>old preempter</i>
~	EBP		EBP	~
<i>dest id</i>	ESI		ESI	<i>partner</i>
~	EDI		EDI	~

The system call can be used by schedulers to define the *priority*, *timeslice length* and *external preempter* of other threads. Furthermore, it delivers thread states. Note that due to security reasons thread state information must be retrieved through the appropriate scheduler.

The system call is only effective, if the current priority of the specified destination is less or equal than the current task's *maximum controlled priority (mcp)*.

## Parameters

*dest id* Destination thread id. The destination thread must currently exist and run on a priority level less than or equal to the current thread's *mcp*. Otherwise, the destination thread is not affected by this system call and all result parameters except *old param word* are undefined.

<i>param word</i>	valid	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="border: 1px solid black; padding: 2px 10px;"><i>m<sub>t</sub></i> (8)</td> <td style="border: 1px solid black; padding: 2px 10px;"><i>e<sub>t</sub></i> (4)</td> <td style="border: 1px solid black; padding: 2px 10px;">0 (4)</td> <td style="border: 1px solid black; padding: 2px 10px;">small (8)</td> <td style="border: 1px solid black; padding: 2px 10px;">prio (8)</td> </tr> </table>	<i>m<sub>t</sub></i> (8)	<i>e<sub>t</sub></i> (4)	0 (4)	small (8)	prio (8)
<i>m<sub>t</sub></i> (8)	<i>e<sub>t</sub></i> (4)	0 (4)	small (8)	prio (8)			
	<i>prio</i>	New priority for destination thread. Must be less than or equal to current thread's <i>mcp</i> .					
	<i>small</i>	(Only effective for Pentium.) Sets the <i>small address space number</i> for the addressed <i>task</i> . On Pentium, small address spaces from 1 to 127 currently available. A value of 0 or 255 in this field does not change the current setting for the task. This field is currently ignored for 486 and PPro.					
	<i>m<sub>t</sub>, e<sub>t</sub></i>	New timeslice length for the destination thread. The timeslice quantum is encoded like a timeout: $4^{15-e_t} m_t \mu s$ . The kernel rounds this value up towards the nearest possible value. Thus the timeslice granularity can be determined by trying to set the timeslice to 1 $\mu s$ . However note that the timeslice granularity may depend on the priority. Timeslice length 0 ( $m_t = 0, e_t \neq 0$ ) is always a possible value. It means that the thread will get no ordinary timeslice, i.e. is blocked. However, even a blocked thread may execute in a timeslice donated to it by ipc.					
	"inv"	(0xFFFFFFFF) The current priority and timeslice length of the thread is not modified.					
	<i>ext preempter</i>	valid Defines the external preempter for the destination thread. (Nilthread is a valid id.)					
	"inv"	(0xFFFFFFFF,~) The current external preempter of the thread is not changed.					

*old param word*

valid

$m_t$ (8)	$e_t$ (4)	$ts$ (4)	$\sim$ (8)	prio (8)
-----------	-----------	----------	------------	----------

*prio*

Old priority of destination thread.

$m_t, e_t$

Old timeslice length of the destination thread:  $4^{15-e_t} m_t \mu s$ .

$ts =$

Thread state:

$0 + k$

*Running*. The thread is ready to execute at user-level.

$4 + k$

*Sending*. A user-invoked ipc send operation currently transfers an outgoing message.

$8 + k$

*Receiving*. A user-invoked ipc receive operation currently receives an incoming message.

C

*Waiting* for receive. A user-invoked ipc receive operation currently waits for an incoming message.

D

*Pending* send. A user-invoked ipc send operation currently waits for the destination (recipient) to become ready to receive.

E

Reserved.

F

*Dead*. The thread is unable to execute.

$k = 0$

*Kernel inactive*. The kernel does not execute an automatic RPC for the thread.

1

*Pager*. The kernel executes a pagefault RPC to the thread's pager.

2

*Internal preempter*. The kernel executes a preemption RPC to the thread's internal preempter.

3

*External preempter*. The kernel executes a preemption RPC to the thread's external preempter.

"inv"

(0xFFFFFFFF) The addressed thread does either not exist or has a priority which exceeds the current thread's *mcp*. All other return parameters are undefined ( $\sim$ ).

*old ext preempter*

Old external preempter of the destination thread.

*partner*

Partner of an active user-invoked ipc operation. This parameter is only valid, if the thread's user state is *sending*, *receiving*, *pending* or *waiting* (4..D). An invalid thread id (0xFFFFFFFF,  $\sim$ ) is delivered if there is no specific partner, i.e. if the thread is in an open receive state.

*time*

$m_w$ (8)	$e_w$ (4)	$e_p$ (4)	$T_{high}$ (16)	EDX
$T_{low}$ (32)				ECX

$T$

Cpu time (48-bit value) in microseconds which has been consumed by the destination thread.

$m_w, e_w$

Current user-level wakeup of the destination thread, encoded like a timeout. The value denotes the still remaining timeout interval. Valid only if the user state is *waiting* (C) or *pending* (D).

$e_p$

Effective pagefault wakeup of the destination thread, encoded like a 4-bit pagefault timeout. The value denotes the still remaining timeout interval. Valid only if the kernel state is *pager* ( $k = 1$ ).

## lthread\_ex\_regs

<i>lthread no</i>	EAX	— INT 0x35 →	EAX	<i>old EFLAGS</i>
<i>ESP</i>	ECX		ECX	<i>old ESP</i>
<i>EIP</i>	EDX		EDX	<i>old EIP</i>
<i>int preempter</i>	EBX		EBX	<i>old preempter</i>
~	EBP		EBP	~
<i>pager</i>	ESI		ESI	<i>old pager</i>
~	EDI		EDI	~

This function reads and writes some register values of a thread in the current task.

It also creates threads. Conceptually, creating a task includes creating all of its threads. Except lthread 0, all these threads run an idle loop. Of course, the kernel does neither allocate control blocks nor time slices etc. to them. Setting stack and instruction pointer of such a thread to valid values then really generates the thread.

Note that this operation reads and writes the *user-level* registers (ESP, EIP and EFLAGS). Ongoing kernel activities are not affected. However an ipc operation is cancelled or aborted. If the ipc is either waiting to send a message or waiting to receive a message, i.e. a message transfer is not yet running, ipc is cancelled (completion code 0x40 or 0x50). If a message transfer is currently running, ipc is aborted (completion code 0xC0 or 0xD0).

## Parameters

*lthread no*

<i>v</i>	0 <sub>(24)</sub>	lthread <sub>(7)</sub>
----------	-------------------	------------------------

Number of addressed lthread (0..127) inside the current task.

<i>v = 0</i>	If ESP and EIP are valid, they specify 32-bit protected-mode values. The addressed thread will execute in 32-bit protected mode afterwards.
<i>v = 1</i>	If ESP and EIP are valid, they specify 16-bit V86-mode values. The addressed thread will execute in V86 mode afterwards.
<i>u = 0</i>	The thread's <i>auto-propagating</i> status is not updated.
<i>u = 1</i>	The thread's <i>auto-propagating</i> status is updated by <i>a</i> .
<i>a = 0</i>	If <i>u = 1</i> , the thread is set to <i>non auto-propagating</i> .
<i>a = 1</i>	If <i>u = 1</i> , the thread is set to <i>auto-propagating</i> .
<i>ESP</i>	valid      New stack pointer (ESP) for the thread. It must point into the user-accessible part of the address space.
	"inv"      (0xFFFFFFFF) The existing stack pointer is not modified.
<i>EIP</i>	valid      New instruction pointer (EIP) for the thread. It must point into the user-accessible part of the address space.
	"inv"      (0xFFFFFFFF) The existing instruction pointer is not modified.
<i>int preempter</i>	valid      Defines the internal preempter used by the thread. ( <i>Nil</i> is a valid id.)
	"inv"      (0xFFFFFFFF,~) The existing internal preempter id is not modified.
<i>pager</i>	valid      Defines the pager used by the thread.
	"inv"      (0xFFFFFFFF,~) The existing pager id is not modified.

<i>old EFLAGS</i>	Flags of the thread. The <i>VM</i> flag specifies whether the thread currently runs in 32-bit protected mode ( <i>VM</i> =0) or in V86 mode ( <i>VM</i> =1). Note that this flag determines the format of the delivered old ESP and EIP.
<i>old ESP</i>	Old stack pointer (ESP) of the thread.
<i>old EIP</i>	Old instruction pointer (EIP) of the thread.
<i>old int preempter</i>	Id of the thread's old internal preempter (may be nilthread).
<i>old pager</i>	Id of the thread's old pager.

## V86-mode Pointers

<i>ESP, old ESP</i>	SS <sub>(16)</sub>	SP <sub>(16)</sub>
<i>EIP, old EIP</i>	CS <sub>(16)</sub>	IP <sub>(16)</sub>

## Example

Signalling can be implemented as follows:

```

signal (lthread) :
    esp := receive signal stack ;
    eip := receive signal ;
    mem [esp --] := 0 ;
    lthread ex regs (lthread, esp, eip, eflags, -, -) ;
    mem [esp --] := eflags ;
    mem [esp --] := eip ;
    mem [idle stack - wordlength] := esp .

```

```

receive signal :
    push all regs ;
    while mem [esp + 8 × wordlength] = 0 do
        thread switch (nilthread)
    od ;
    pop all regs ;
    pop (esp) ;
    jmp (signal eip) .

```

## task\_new

<i>mcp / new chief</i>	EAX	— INT 0x36 →	EAX	~
<i>initial ESP</i>	ECX		ECX	~
<i>initial EIP</i>	EDX		EDX	~
<i>pager</i>	EBX		EBX	~
~	EBP		EBP	~
<i>dest task</i>	ESI		ESI	<i>new task</i>
~	EDI		EDI	~

This function deletes and/or creates a task. Deletion of a task means that the address space of the task and all threads of the task disappear. The cputime of all deleted threads is added to the cputime of the deleting thread. If the deleted task was chief of a clan, all tasks of the clan are deleted as well.

Tasks may be created as *active* or *inactive*. For an active task, a new address space is created together with 128 threads. Lthread 0 is started, the other ones wait for a “real” creation by lthread\_ex\_regs. An inactive task is empty. It occupies no resources, has no address space and no threads. Communication with inactive tasks is not possible. Loosely speaking, inactive tasks are not really existing but represent only the right to create an active task.

A newly created task gets the creator as its chief, i.e. it is created inside the creator’s clan. Symmetrically, a task can only be deleted either directly by its chief (its creator) or indirectly by a higher-level chief.

### Parameters

<i>dest task</i>		Task id of an <i>existing</i> task (active or inactive) whose chief is the current task. If one of these preconditions is not fulfilled, the system call has no effect. Simultaneously, a new task <i>with the same task number</i> is created. It may be active or inactive (see next parameter).
<i>pager</i>	$\neq nil$	The new task is created as <i>active</i> . The specified pager is associated with lthread 0.
	$= nil$	(0) The new task is created as <i>inactive</i> . No lthread is created.
<i>ESP</i>		Initial stack pointer for lthread 0 if the new task is created as an active one. Ignored otherwise.
<i>EIP</i>		Initial instruction pointer for lthread 0 if the new task is created as an active one. Ignored otherwise.
<i>mcp</i>		Maximum controlled priority (mcp) defines the highest priority which can be ruled by the new task acting as a scheduler. The new task’s effective mcp is the minimum of the creator’s mcp and the specified mcp. EAX contains this parameter, if the newly generated task is an <i>active</i> task, i.e. has a pager and at least lthread 0.
<i>new chief</i>		Specifies the chief of the new inactive task. This mechanism permits to transfer inactive (“empty”) tasks to other tasks. Transferring an inactive task to the specified chief means to transfer the related right to create a task. Note that the task number remains unchanged. EAX contains this parameter, if the newly generated task is an <i>inactive</i> task, i.e. has no pager and no threads. EAX contains only the lower 32 bits of the new chief’s task id. (The chief must reside in the same site.)

*new task id*    *≠nil*    Task creation succeeded. If the new task is active, the new task id will have a new version number so that it differs from all task ids used earlier. Chief and task number are the same as in *dest task*. If the new task is created inactive, the chief is taken from the *chief* parameter; the task number remains unchanged. The version is undefined so that the new task id might be identical with a formerly (but not currently and not in future) valid task id. This is safe since communication with inactive tasks is impossible.

*=nil*    (0) The task creation failed.

## 2.4 Processor Mirroring

### 2.4.1 Segments

L4 uses a flat (unsegmented) memory model. There are only two segments available: *user\_space*, a read/write segment, and *user\_space\_exec*, an executable segment. Both cover (at least) the complete user-level address space.

The values of both segment selectors are *undefined*. When a thread is created, its segment registers SS, DS, ES, FS and GS are initialized with *user\_space*, CS with *user\_space\_exec*. Whenever the kernel detects a general protection exception and the segment registers are not loaded properly, it reloads them with the above mentioned selectors. From the user's point of view, the segment registers cannot be modified.

However, the binary representation of *user\_space* and *user\_space\_exec* may change at any point during program execution. Never rely on this value.

Furthermore, the LSL (load segment limit) machine instruction may deliver wrong segment limits, even floating ones. The result of this instruction is always undefined.

### 2.4.2 Exception Handling

#PF (page fault), #MC (machine check exception) and some #GP (general protection) exceptions are handled by the kernel. The other exceptions are mirrored to the virtual processor of the thread which raised the exception.

The mirrored exception handling works as described in the processor manuals.

*LIDT [EAX]*

This machine instruction is emulated by the kernel and operates per thread. Any thread should install an IDT by this instruction. The length field of the IDT vector is not interpreted. The IDT has always a length of 256 bytes, i.e. covers the Intel-reserved exceptions 0 to 31.

The IDT must have the format described in the processor manuals. However, only trap gates can be used in a user-level IDT. The segment selectors in the IDT are ignored, since all segments describe the flat address space.

Invalid IDT addresses or invalid exception-handler addresses do not raise a double fault; instead, the current thread is shut down.

Note that this mechanism deals only with exceptions, not INT *n* instructions. Executing an INT *n* in 32-bit mode will always raise a #GP (general protection). The general-protection handler may interpret the error code ( $8n + 2$ , see processor manual) and emulate the INT *n* accordingly.

### 2.4.3 Debug Registers

User-level debug registers exist per thread. DR0...3, DR6 and DR7 can be accessed by the machine instructions MOV DR*x*,*n* and MOV *r*,DR*x*. However, only task-local breakpoints can be activated, i.e. bits L0...3 in DR7 cannot be set. Breakpoints operate per thread. Breakpoints are signalled as #DB exception (INT 1).

Note that user-level breakpoints are suspended when kernel breakpoints are set by the kernel debugger.



## 2.5 The Kernel-Info Page

The kernel-info page contains kernel-version data, memory descriptors *and the clock*. The remainder of the page is undefined. (In fact, it contains kernel code.) The kernel-info page is mapped *read-only* in the  $\sigma_0$ -address space.  $\sigma_0$  can use the memory descriptors for its memory management.  $\sigma_0$  can map the page read-only to other address spaces.

L4 version parts				$L \times 16$
~		bus frequency	processor frequency	0xB0
~		clock		0xA0
dedicated mem4.high	dedicated mem4.low	dedicated mem3.high	dedicated mem3.low	0x90
dedicated mem2.high	dedicated mem2.low	dedicated mem1.high	dedicated mem1.low	0x80
dedicated mem0.high	dedicated mem0.low	reserved mem1.high	reserved mem1.low	0x70
reserved mem0.high	reserved mem0.low	main mem.high	main mem.low	0x60
~				0x50
~				0x40
~				0x30
~				0x20
~				0x10
~	2	$L$	version word	“L4 $\mu$ K”
+0xC		+8	+4	+0

- mem.low*                      Physical address of first byte of region. Must be page aligned.
- mem.high*                    Physical address of first byte beyond the region. Must be page aligned. If *mem.high* = 0, the region is empty.
- main mem*                     Main memory region.
- reserved mem*                This region must not be used. It contains kernel code (reserved mem0) or data (reserved mem1) grabbed by the kernel before  $\sigma_0$  was initialized.
- dedicated mem*               This region contains dedicated memory which cannot be used as standard memory. For example, [640K, 1M] is a popular dedicated memory region.
- clock*                         System clock in  $\mu$ s.
- processor frequency*        Processor’s external clock rate in MHz.
- bus frequency*                Processor’s external clock rate in MHz.

## 2.6 Page-Fault and Preemption RPC

### Page Fault RPC

<b>kernel sends:</b>	w0 (EDX)	<table border="1"> <tr> <td colspan="2">fault address / 4 <sub>(30)</sub></td> <td><i>w</i></td> <td>~</td> </tr> <tr> <td colspan="4">faulting user-level EIP <sub>(32)</sub></td> </tr> </table>	fault address / 4 <sub>(30)</sub>		<i>w</i>	~	faulting user-level EIP <sub>(32)</sub>			
	fault address / 4 <sub>(30)</sub>		<i>w</i>	~						
faulting user-level EIP <sub>(32)</sub>										
	w1 (EBX)									

*w* = 0      Read page fault.  
*w* = 1      Write page fault.

**kernel receives:**      The receive fpage covers the complete user address space. The kernel accepts mappings or grants into this region as well as a simple 2-word copy message. The received message is ignored!

timeouts used for pagefault RPC	PF at user level	PF at ipc in receiver's space	PF at ipc in sender's space
snd	∞	sender's snd pf	receiver's rcv pf
rcv	∞	sender's snd pf	receiver's rcv pf
snd pf	∞	sender's snd pf	receiver's rcv pf
rcv pf	∞	sender's snd pf	receiver's rcv pf

### Preemption RPC

<b>kernel sends:</b>	w0 (EDX)	<table border="1"> <tr> <td>user-level ESP <sub>(32)</sub></td> </tr> <tr> <td>user-level EIP <sub>(32)</sub></td> </tr> </table>	user-level ESP <sub>(32)</sub>	user-level EIP <sub>(32)</sub>
	user-level ESP <sub>(32)</sub>			
user-level EIP <sub>(32)</sub>				
	w1 (EBX)			

**kernel receives:**      The kernel accepts only a simple 2-word reply. Its content is ignored!

timeouts used for preemption RPC	
snd	∞
rcv	∞
snd pf	∞
rcv pf	∞

## 2.7 $\sigma_0$ RPC protocol

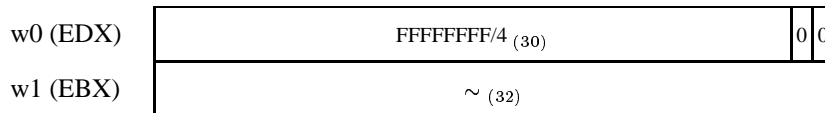
$\sigma_0$  is the initial address space. Although it is *not* part of the kernel, its basic protocol is defined with the Nucleus. Special  $\sigma_0$  implementations may extend this protocol.

The address space  $\sigma_0$  is idempotent, i.e. all virtual addresses in this address space are identical to the corresponding physical address. Note that pages requested from  $\sigma_0$  continue to be mapped idempotent if the receiver specifies its complete address space as receive fpage.

$\sigma_0$  gives pages to the kernel and to arbitrary tasks, but only once. The idea is that all pagers request the memory they need in the startup phase of the system so that afterwards  $\sigma_0$  has spent all its memory. Further requests will then automatically be denied.

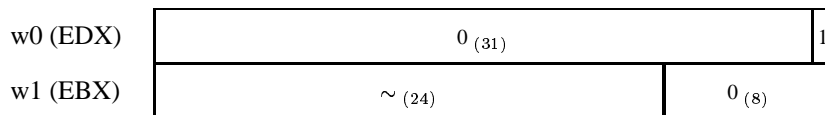
### Kernel Communication

$\sigma_0$  receives from a *kernel* thread:

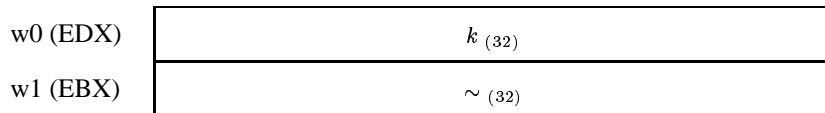


*Intended Action:* For reply,  $\sigma_0$  should *grant* a page (4K) to the kernel thread. The page might be located at an arbitrary position but must contain ordinary memory. If no more memory is available,  $\sigma_0$  should reply a 0-word instead of a page.

$\sigma_0$  receives from a *kernel* thread:



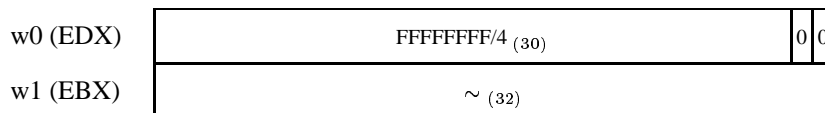
*Intended reply:*



$k$  is the number of pages recommended by  $\sigma_0$  for kernel use (pagetables and other kernel-internal data).

### General Memory Mapping

$\sigma_0$  receives from a *non-kernel* thread:



*Intended Action:* For reply,  $\sigma_0$  should *map* a page (4K) writable to the requester. The page might be located at an arbitrary position but must contain ordinary memory. If no more free page is available,  $\sigma_0$  should reply a 0-word instead of a page. The mapped page must be marked and must not further be mapped or granted by  $\sigma_0$  to another task. However, multiple mapping to the same requester should be supported.

$\sigma_0$  receives from any thread:

w0 (EDX)	address $\leq$ 40000000/4 <sub>(30)</sub>	0	0
w1 (EBX)	~ <sub>(32)</sub>		

*Intended Action:*

For reply,  $\sigma_0$  should *map* the specified physical page frame (4K) writable to the requester. If the page is already mapped or is not available,  $\sigma_0$  should reply a 0-word instead of a page.

The mapped page must be marked and must not further be mapped or granted by  $\sigma_0$  to another task. However, multiple mapping to the same requester should be supported.

$\sigma_0$  receives from any thread:

w0 (EDX)	40000000/4 < a $\leq$ C0000000/4 <sub>(30)</sub>	0	0
w1 (EBX)	~ <sub>(32)</sub>		

*Intended Action:*

For reply,  $\sigma_0$  should *map* the physical 4M superpage with address  $a - 40000000$  writable to the requester. If the page is already mapped or is not available,  $\sigma_0$  should reply a 0-word instead of a page.

The mapped superpage must be marked and must not further be mapped or granted by  $\sigma_0$  to another task. However, multiple mapping to the same requester should be supported.

$\sigma_0$  receives from any thread:

w0 (EDX)	0 <sub>(31)</sub>	1
w1 (EBX)	~ <sub>(24)</sub>	1 <sub>(8)</sub>

*Intended Action:*

For reply,  $\sigma_0$  should *map read only* the kernel info page to the requester. This page can be mapped multiply.

## 2.8 Starting L4

For booting L4, see appendix A.

After booting, L4 enters protected mode if started in real mode, enables paging and initializes itself. It generates the basic address space-servers  $\sigma_0$  and a task *root server* which is intended to boot the higher-level system.

$\sigma_0$  and the *root server* are user-level tasks and not part of the pure Nucleus. The predefined ones can be replaced by modifying the following table in the L4 image before starting L4. The kernel debugger *kdebug* is also not part of the Nucleus and can accordingly be replaced by modifying the table.

dedicated mem4.high	dedicated mem4.low	dedicated mem3.high	dedicated mem3.low	0x1090
dedicated mem2.high	dedicated mem2.low	dedicated mem1.high	dedicated mem1.low	0x1080
~				0x1070
~				0x1060
kdebug permissions	kdebug configuration	~	L4 configuration	0x1050
<i>root server</i> end+1	<i>root server</i> begin	<i>root server</i> start EIP	<i>root server</i> start ESP	0x1040
$\sigma_1$ end+1	$\sigma_1$ begin	$\sigma_1$ start EIP	$\sigma_1$ start ESP	0x1030
$\sigma_0$ end+1	$\sigma_0$ begin	$\sigma_0$ start EIP	$\sigma_0$ start ESP	0x1020
<i>kdebug</i> end+1	<i>kdebug</i> begin	<i>kdebug</i> exception	<i>kdebug</i> init	0x1010
	+0xC	+8	+4	+0

0x1010 ... are offsets relative to the load address. The EIP and ESP values however, are absolute 32-bit addresses. The appropriate code must be loaded at these addresses before L4 is started. Note that the predefined root server currently executes only a brief kernel test.

<i>mem.low</i>	Physical address of first byte of region. Must be page aligned.
<i>mem.high</i>	Physical address of first byte beyond the region. Must be page aligned. If <i>mem.high</i> = 0, the region is empty.
<i>dedicated mem</i>	This region contains dedicated memory which cannot be used as standard memory. For example, [640K, 1M] is a popular dedicated memory region.
<i>begin</i>	Physical address of the first byte of a server's code+data area.
<i>end</i>	Physical address of the last byte of a server's code+data area.
<i>start EIP</i>	Physical address of a server's initial instruction pointer (start).
<i>start ESP</i>	Physical address of a server's initial stack pointer (stack bottom).
<i>kdebug init</i>	Physical address of <i>kdebug</i> 's initialization routine.
<i>kdebug exception</i>	Physical address of <i>kdebug</i> 's debug-exception handler.

### L4 configuration

~ (16)	pnodes (8)	ptabs (8)
--------	------------	-----------

ptabs number of ptabs (4K each) per 4M of physical memory which the Nucleus should allocate for page tables. 0 indicates to use the default value. A value of 128, for example, specifies that 1/8 of memory should be reserved for page tables.

pnodes number of pnode entries (16 bytes each) the Nucleus should allocate per physical frame. A value of 0 indicates to use the default value.

### kdebug configuration

port (12)	rate (4)	~ (7)	s	pages (8)
-----------	----------	-------	---	-----------

pages The number of 4K pages that kdebug should allocate for its trace buffer. A value of 0 indicates no trace buffer.

s = 1 The Nucleus enters kdebug before starting the root server.

port Initially, kdebug should use the serial line base IO address *port* for output and input. A port address of 0 indicates to use the integrated console (keyboard and display) instead of a serial line.

rate determines the default baud rate for kdebug when using a serial line. If *port* ≠ 0, this is also the initial baud rate. Possible values:

- = 1 115.2 Kbd
- = 2 57.6 Kbd
- = 3 38.4 Kbd
- = 6 19.2 Kbd
- = 12 9.6 Kbd

### kdebug permissions

~ (18)	p	i	w	d	r	m	k (8)
--------	---	---	---	---	---	---	-------

k = 0 Any task can use kdebug from user-level.

k > 0 Only tasks 1 to *k* can use kdebug from user-level. Threads of other tasks will be shut down when invoking kdebug.

m = 1 Kdebug may display mapping.

r = 1 Kdebug may display user registers.

d = 1 Kdebug may display user memory.

w = 1 Kdebug may modify memory, register, mappings and tcbs.

i = 1 Kdebug may read from and write to IO ports.

p = 1 Kdebug may protocol page faults and IPC.

# Appendix A

## Booting

### PC-compatible Machines

L4 can be loaded at any 16-byte-aligned location beyond 0x1000 in physical memory. It can be started in real mode or in 32-bit protected mode at address 0x100 or 0x1000 relative to its load address. The protected-mode conditions are compliant to the Multiboot Standard Version 0.6.

Start Preconditions		
	Real Mode	32-bit Protected Mode
load base ( $L$ )	$L \geq 0x1000$ , 16-byte aligned	$L \geq 0x1000$
load offset ( $X$ )	$X = 0x100$ or $X = 0x1000$	$X = 0x100$ or $X = 0x1000$
Interrupts	disabled	disabled
Gate A20	~	open
EFLAGS	I=0	I=0, VM=0
CR0	PE=0	PE=1, PG=0
(E)IP	$X$	$L + X$
CS	$L/16$	0, 4GB, 32-bit exec
SS,DS,ES	~	0, 4GB, read/write
EAX	~	0x2BADB002
EBX	~	* $P$
$\langle P + 0 \rangle$		~ OR 1
$\langle P + 4 \rangle$	n/a	below 640 K mem in K
$\langle P + 8 \rangle$		beyond 1M mem in K
all remaining registers & flags (general, floating point, ESP, xDT, TR, CRx, DRx)	~	~

L4 relocates itself to 0x1000, enters protected mode if started in real mode, enables paging and initializes itself.