

# Übung zu Betriebssystemtechnik

## Nachrichtenaustausch und Copy-on-Write

---

27. Juni 2019

Bernhard Heinloth, Andreas Ziegler,  
Christian Eichler & Harald Böhm

Lehrstuhl für Informatik 4  
Friedrich-Alexander-Universität Erlangen-Nürnberg



Lehrstuhl für Verteilte Systeme  
und Betriebssysteme



FRIEDRICH-ALEXANDER  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG

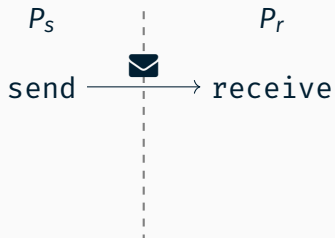
TECHNISCHE FAKULTÄT

## Ziel: Getrennte Adressräume überwinden

⇒ Nachrichtenaustausch

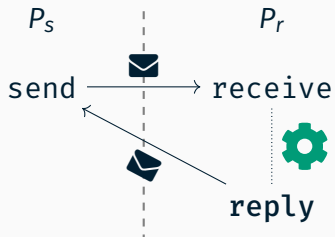
## Ziel: Getrennte Adressräume überwinden

⇒ Nachrichtenaustausch:  
send, receive



# Ziel: Getrennte Adressräume überwinden

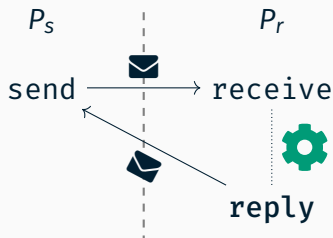
⇒ Nachrichtenaustausch:  
send, receive, reply



Warum **reply**?

## Ziel: Getrennte Adressräume überwinden

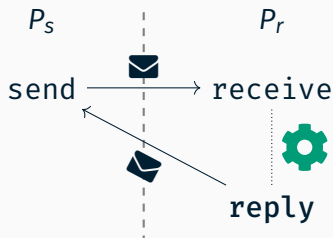
⇒ Nachrichtenaustausch:  
send, receive, **reply**



Warum **reply**? ⇒ Keine Timer nötig um DoS auf Server-Seite zu verhindern!

# Ziel: Getrennte Adressräume überwinden

⇒ Nachrichtenaustausch:  
send, receive, **reply**



Warum **reply**? ⇒ Keine Timer nötig um DoS auf Server-Seite zu verhindern!

## Begriffsklärung: (a-)synchron / (nicht-)blockierend

<b>synchron</b>	Prozess schläft bis Senden fertig
<b>asynchron</b>	Prozess kehrt während Senden zurück
<b>(nicht-)blockierend</b>	Synchronisation der Operationen

# Umsetzung

---

## Ablauf send

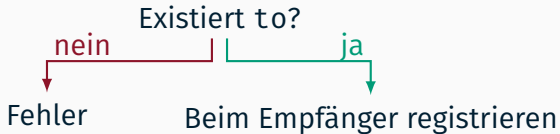
```
send(to, data, len, replybuf, replylen)
```

Existiert to?



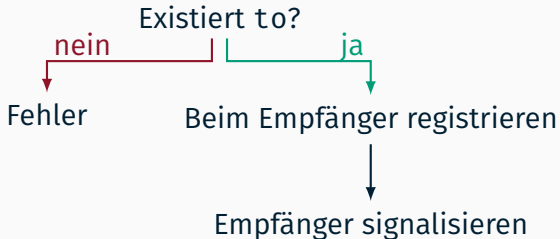
## Ablauf send

`send(to, data, len, replybuf, replylen)`



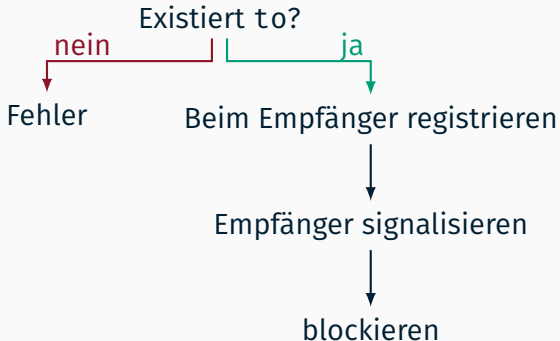
## Ablauf send

`send(to, data, len, replybuf, replylen)`



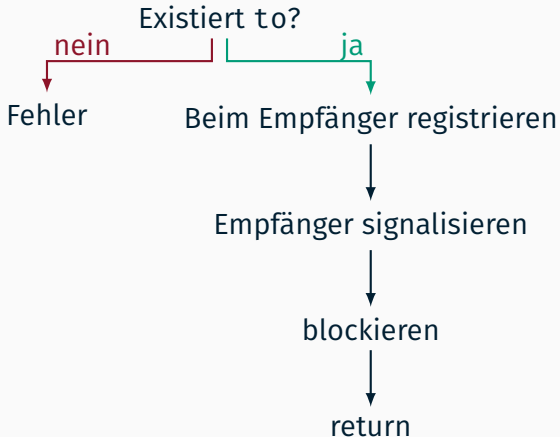
## Ablauf send

`send(to, data, len, replybuf, replylen)`



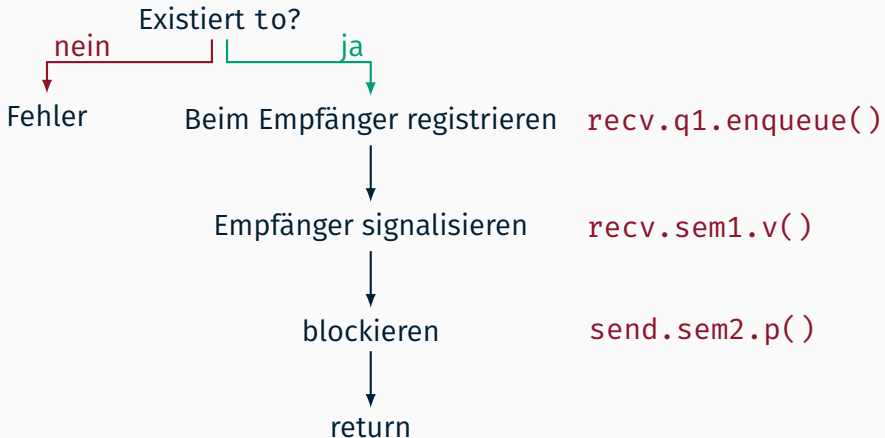
## Ablauf send

`send(to, data, len, replybuf, replylen)`



## Ablauf send

send(to, data, len, replybuf, replylen)



## Ablauf receive

```
recv(data, len)
```

Auf Signal eines Senders warten

## Ablauf receive

```
recv(data, len)
```

Auf Signal eines Senders warten



Infos des Senders entgegennehmen

## Ablauf receive

`recv(data, len)`

Auf Signal eines Senders warten



Infos des Senders entgegennehmen



Daten übertragen (Sender → Empfänger)



## Ablauf receive

`recv(data, len)`

Auf Signal eines Senders warten



Infos des Senders entgegennehmen



Daten übertragen (Sender → Empfänger)



Nachricht als empfangen markieren

## Ablauf receive

`recv(data, len)`

Auf Signal eines Senders warten



Infos des Senders entgegennehmen



Daten übertragen (Sender → Empfänger)



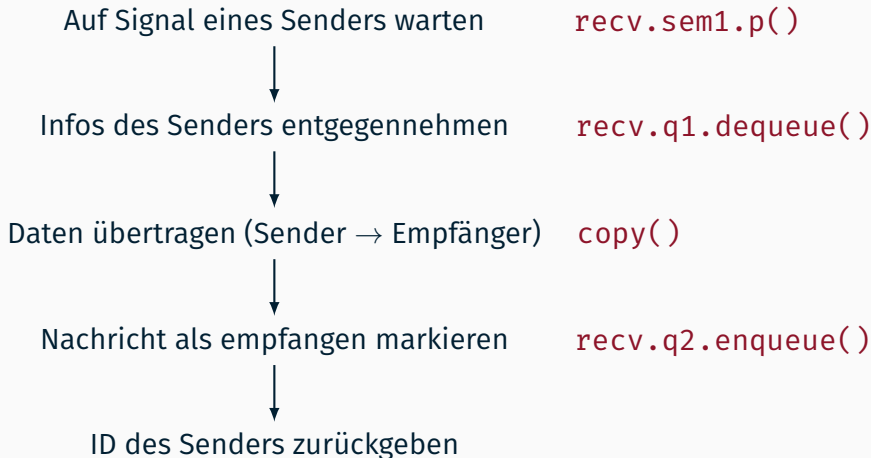
Nachricht als empfangen markieren



ID des Senders zurückgeben

## Ablauf receive

`recv(data, len)`



## Ablauf reply

```
reply(to, data, len)
```

Gab es ein send von to?

## Ablauf reply

reply(to, data, len)

Gab es ein send von to?

Fehler

Daten übertragen  
(Empfänger → Sender)

## Ablauf reply

reply(to, data, len)

Gab es ein send von to?

Fehler

Daten übertragen  
(Empfänger → Sender)

Sender deblockieren

## Ablauf reply

reply(to, data, len)

Gab es ein send von to?

Fehler

Daten übertragen  
(Empfänger → Sender)

Sender deblockieren

return

# Ablauf reply

reply(to, data, len)

Gab es ein send von to?

Fehler

Daten übertragen  
(Empfänger → Sender)

Sender deblockieren

return

Nachricht mit ID in Queue 2?

copy()

send.sem2.v()



# Kopieren zwischen Adressräumen

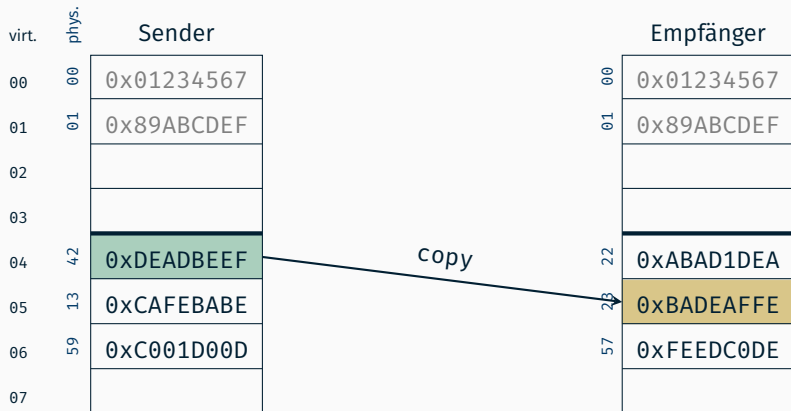
---

# Kopieren zwischen Adressräumen

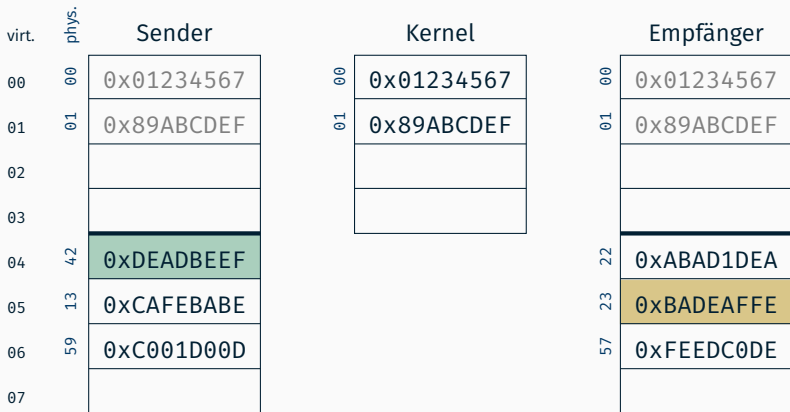
virt.	phys.	Sender
00	00	0x01234567
01	01	0x89ABCDEF
02		
03		
04	42	0xDEADBEEF
05	13	0xCAFEBAFE
06	59	0xC001D00D
07		

	Empfänger
00	0x01234567
01	0x89ABCDEF
22	0xABAD1DEA
23	0xBADEAFFE
57	0xFEEDC0DE

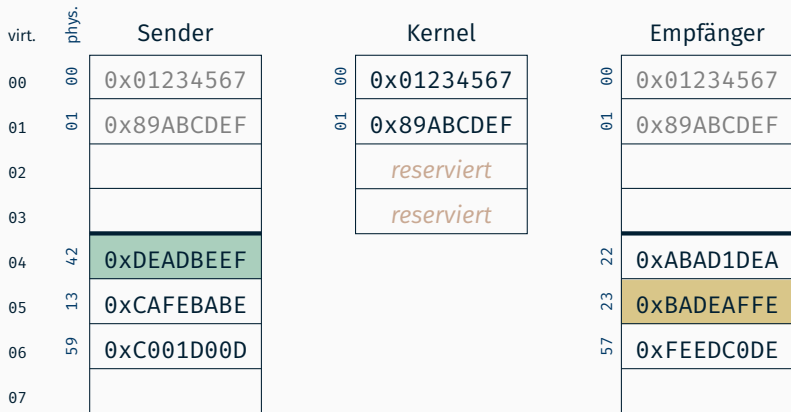
# Kopieren zwischen Adressräumen



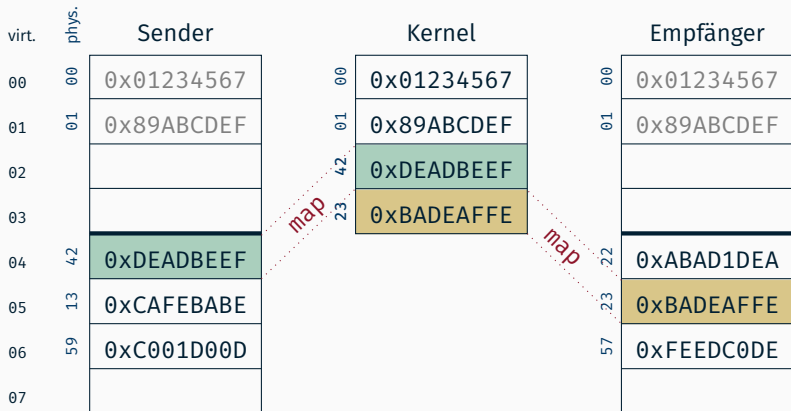
# Kopieren zwischen Adressräumen



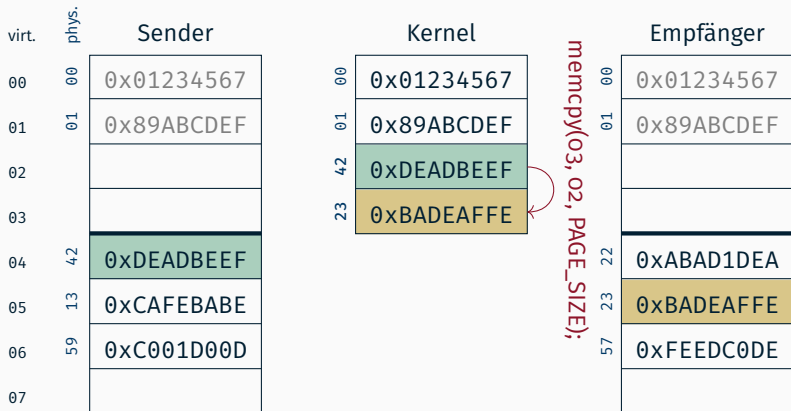
# Kopieren zwischen Adressräumen



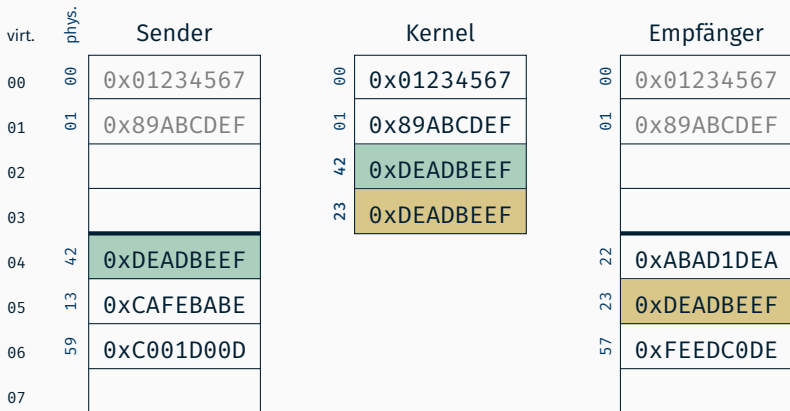
# Kopieren zwischen Adressräumen



# Kopieren zwischen Adressräumen

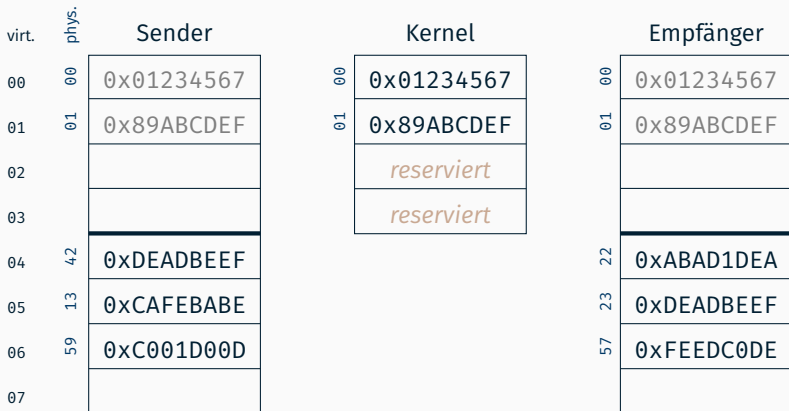


# Kopieren zwischen Adressräumen

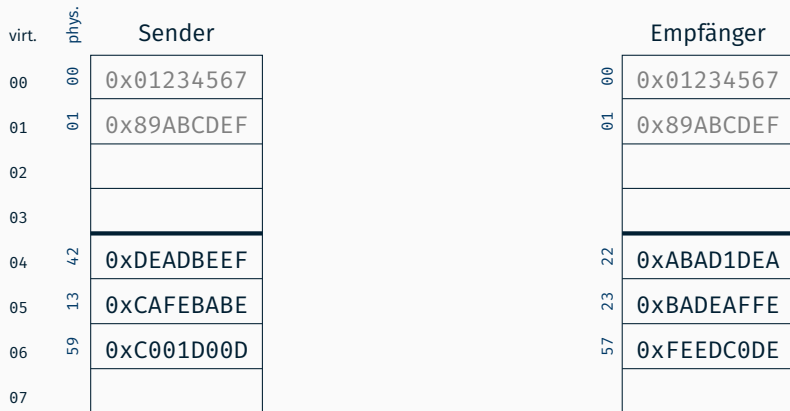




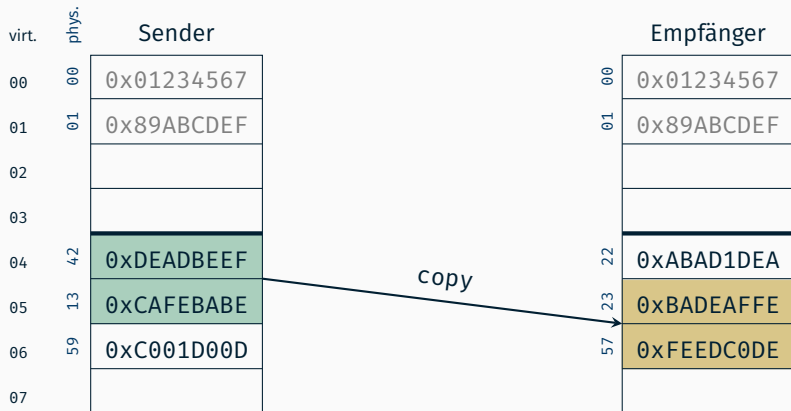
# Kopieren zwischen Adressräumen



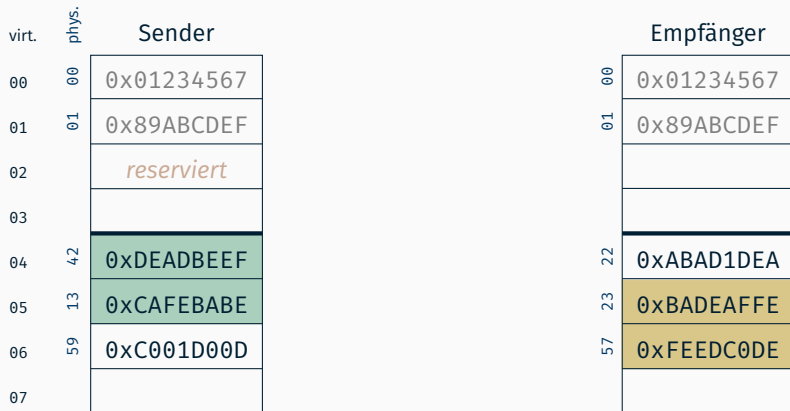
# Kopieren zwischen Adressräumen (optimiert)



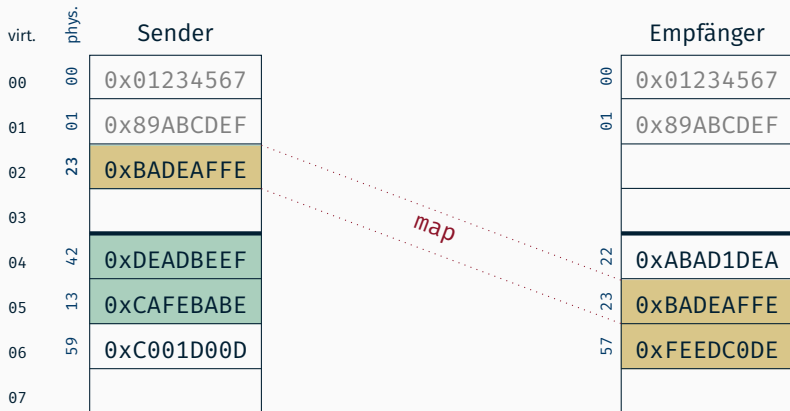
# Kopieren zwischen Adressräumen (optimiert)



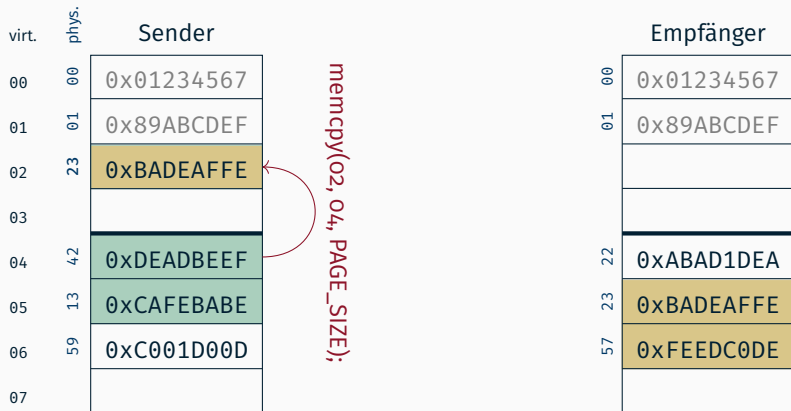
# Kopieren zwischen Adressräumen (optimiert)



# Kopieren zwischen Adressräumen (optimiert)



# Kopieren zwischen Adressräumen (optimiert)

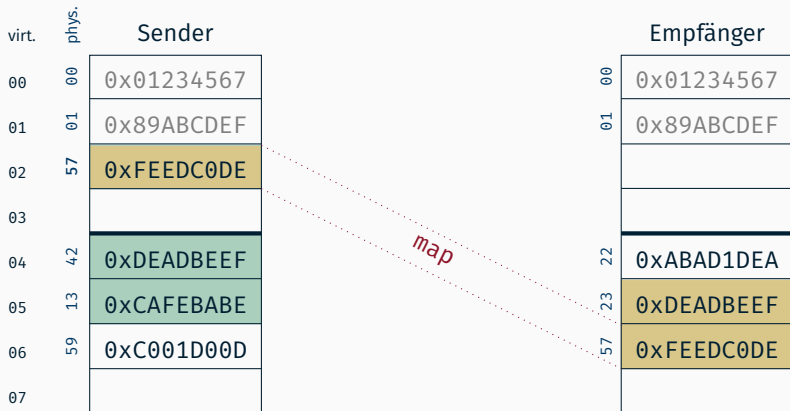


# Kopieren zwischen Adressräumen (optimiert)

virt.	phys.	Sender
00	00	0x01234567
01	01	0x89ABCDEF
02	23	0xDEADBEEF
03		
04	42	0xDEADBEEF
05	13	0xCAFEBAFE
06	59	0xC001D00D
07		

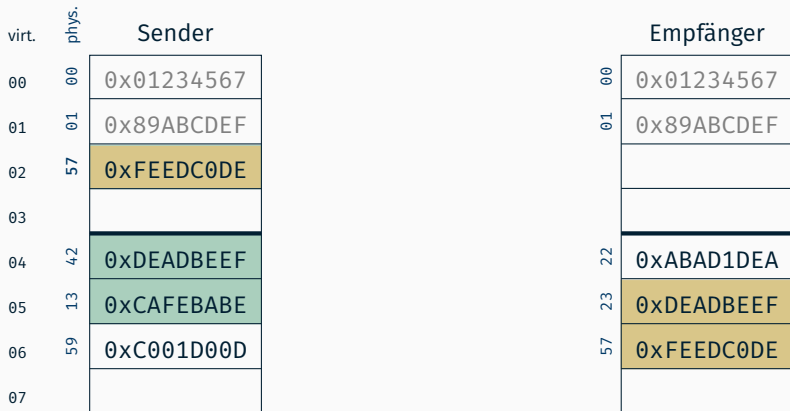
	Empfänger
00	0x01234567
01	0x89ABCDEF
22	0xABAD1DEA
23	0xDEADBEEF
57	0xFEEDC0DE

# Kopieren zwischen Adressräumen (optimiert)





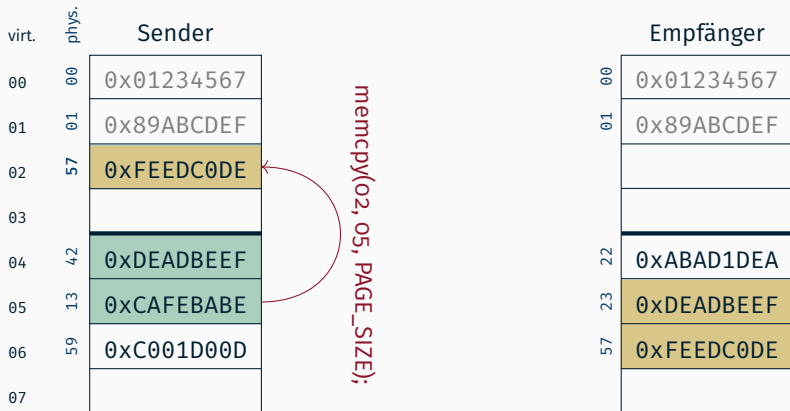
# Kopieren zwischen Adressräumen (optimiert)



**Achtung:** Nach Änderung des Mappings TLB invalidieren!

```
asm volatile("invlpg_{%0}": : "r"(address) : "memory");
```

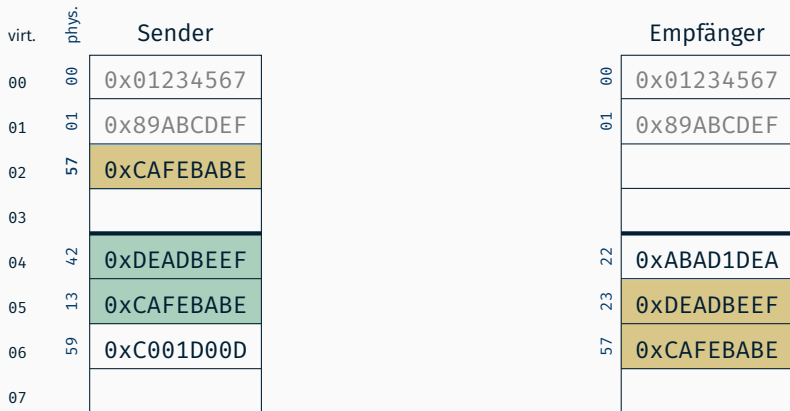
# Kopieren zwischen Adressräumen (optimiert)



**Achtung:** Nach Änderung des Mappings TLB invalidieren!

```
asm volatile("invlpg_{%0}": : "r"(address) : "memory");
```

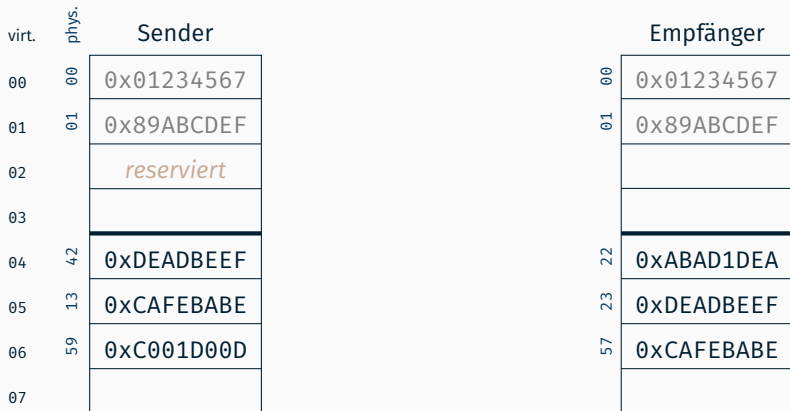
# Kopieren zwischen Adressräumen (optimiert)



**Achtung:** Nach Änderung des Mappings TLB invalidieren!

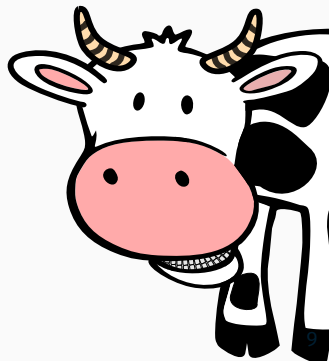
```
asm volatile("invlpg_{%0}": : "r"(address) : "memory");
```

# Kopieren zwischen Adressräumen (optimiert)



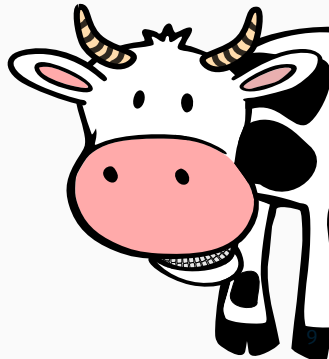
**Achtung:** Nach Änderung des Mappings TLB invalidieren!

```
asm volatile("invlpg_{%0}": : "r"(address) : "memory");
```

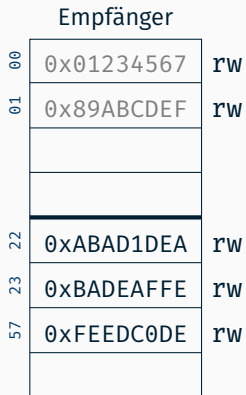


## Copy-on-Write (CoW)

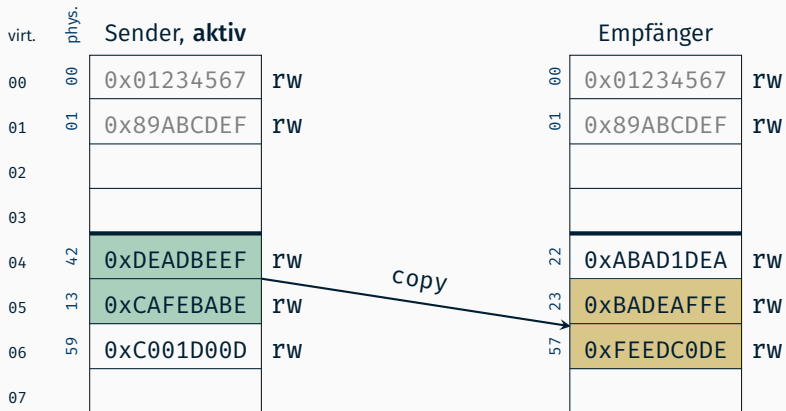
---



# Copy-on-Write

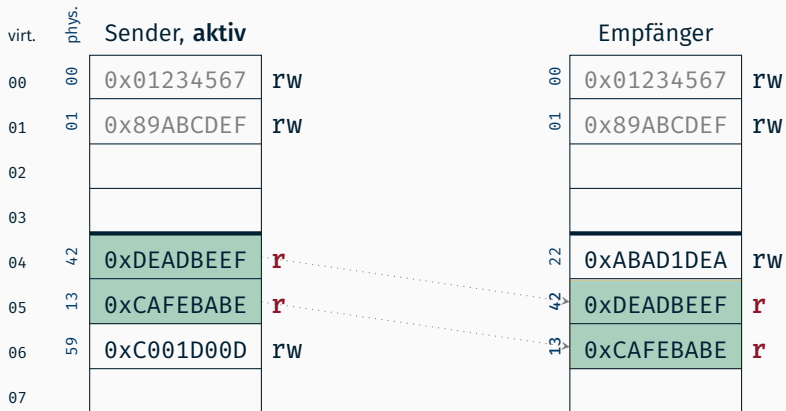


# Copy-on-Write

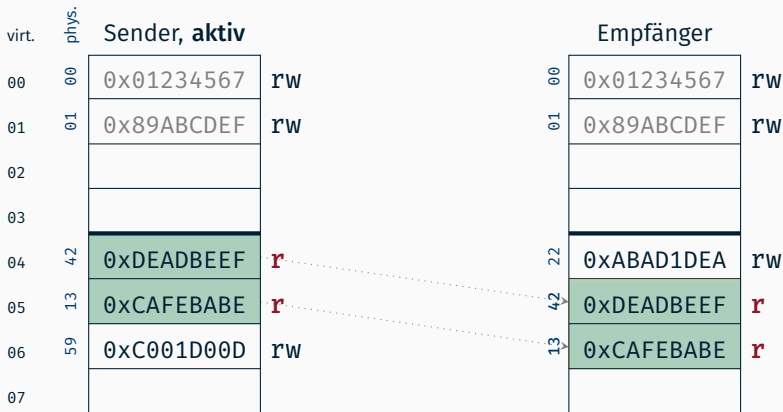




# Copy-on-Write



# Copy-on-Write



**Achtung:** Nach Änderung des Mappings TLB invalidieren!

```
asm volatile("invlpg_{%0}": : "r"(address) : "memory");
```

**Wann/wo werden Seiten tatsächlich kopiert?**

Wann/wo werden Seiten tatsächlich kopiert?

- Im Pagefault-Handler (Trap bei Schreibzugriffen)



**Wann/wo werden Seiten tatsächlich kopiert?**

- Im Pagefault-Handler (Trap bei Schreibzugriffen)

**Was passiert, wenn Puffer weitergeschickt wird?**

### **Wann/wo werden Seiten tatsächlich kopiert?**

- Im Pagefault-Handler (Trap bei Schreibzugriffen)

### **Was passiert, wenn Puffer weitergeschickt wird?**

- das Gleiche!

### Wann/wo werden Seiten tatsächlich kopiert?

- Im Pagefault-Handler (Trap bei Schreibzugriffen)

### Was passiert, wenn Puffer weitergeschickt wird?

- das Gleiche!
- Referenzen zählen (auch bei `fork()`)

### Implementierung des Referenzzählers?

## Wann/wo werden Seiten tatsächlich kopiert?

- Im Pagefault-Handler (Trap bei Schreibzugriffen)

## Was passiert, wenn Puffer weitergeschickt wird?

- das Gleiche!
- Referenzen zählen (auch bei `fork()`)

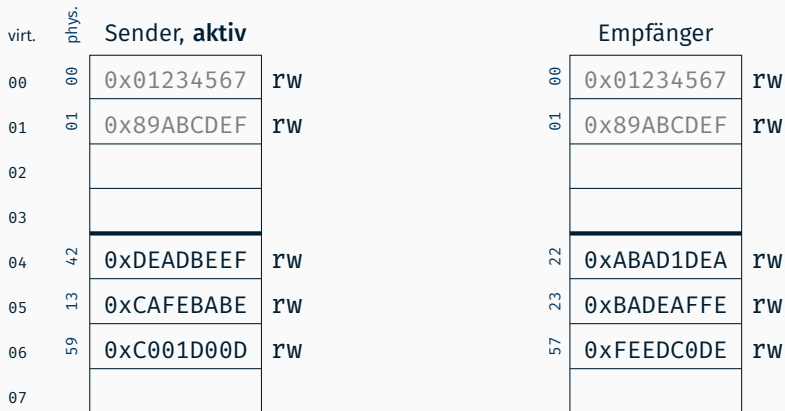
## Implementierung des Referenzzählers?

- (z.B.) Shadow-Pagetable

```
struct PTE_Counter {
    uint32_t present:1,
            counter:31;
};
```



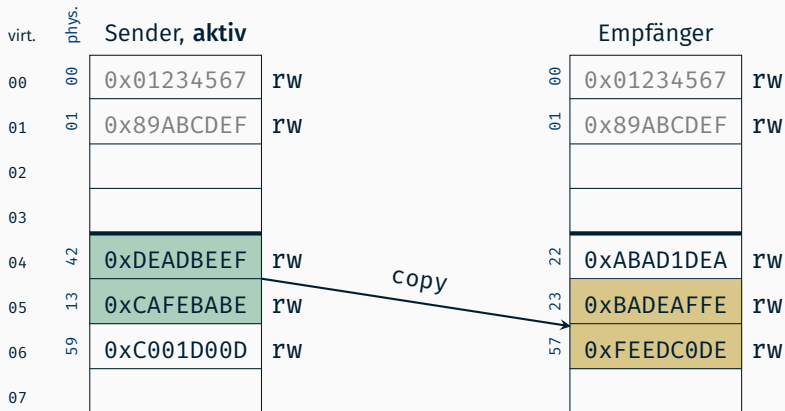
# Copy-on-Write



Referenzzähler:

phys.	...	12	13	14	...	21	22	23	...	42	...	57	58	59	...
Ref.	...	0	1	0	...	0	1	1	...	1	...	1	0	1	...

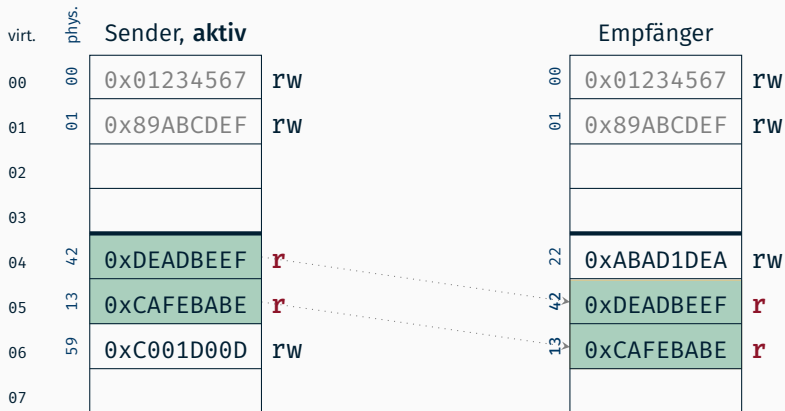
# Copy-on-Write



Referenzzähler:

<b>phys.</b>	...	12	13	14	...	21	22	23	...	42	...	57	58	59	...
<b>Ref.</b>	...	0	1	0	...	0	1	1	...	1	...	1	0	1	...

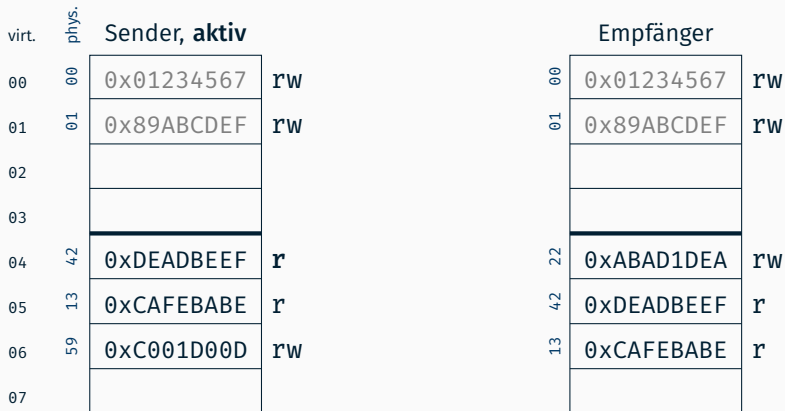
# Copy-on-Write



Referenzzähler:

phys.	...	12	13	14	...	21	22	23	...	42	...	57	58	59	...
Ref.	...	0	2	0	...	0	1	0	...	2	...	0	0	1	...

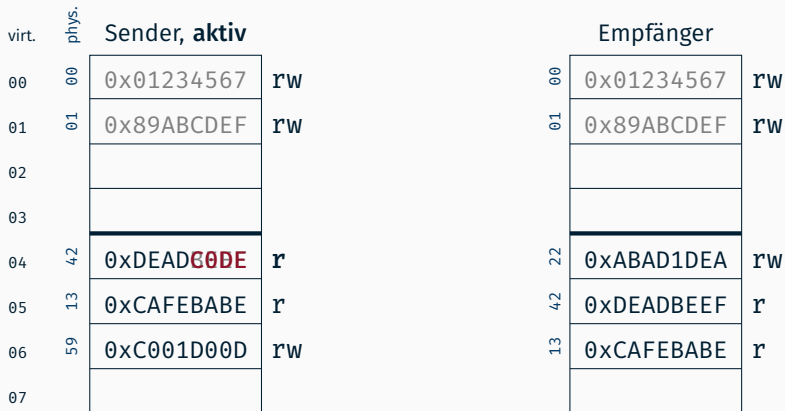
# Copy-on-Write



Referenzzähler:

phys.	...	12	13	14	...	21	22	23	...	42	...	57	58	59	...
Ref.	...	0	2	0	...	0	1	0	...	2	...	0	0	1	...

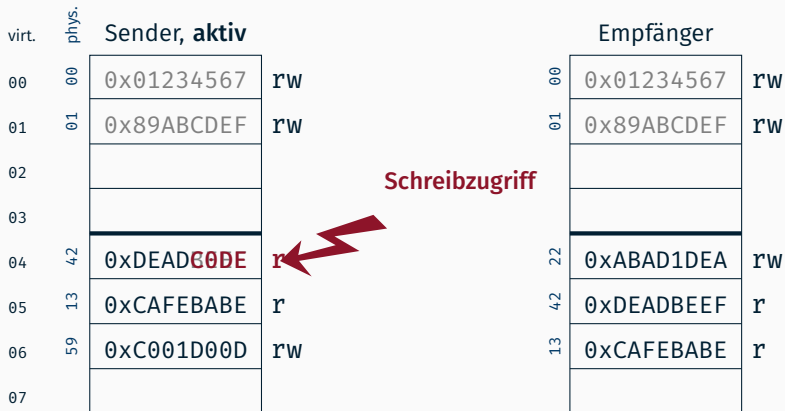
# Copy-on-Write



Referenzzähler:

phys.	...	12	13	14	...	21	22	23	...	42	...	57	58	59	...
Ref.	...	0	2	0	...	0	1	0	...	2	...	0	0	1	...

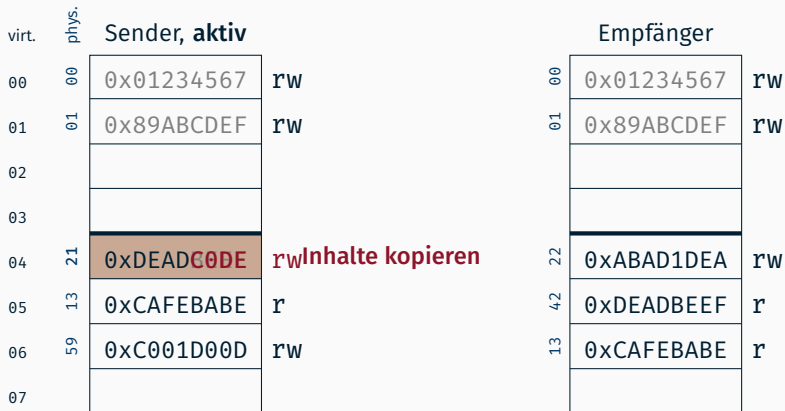
# Copy-on-Write



Referenzzähler:

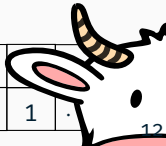
phys.	...	12	13	14	...	21	22	23	...	42	...	57	58	59	...
Ref.	...	0	2	0	...	0	1	0	...	2	...	0	0	1	...

# Copy-on-Write

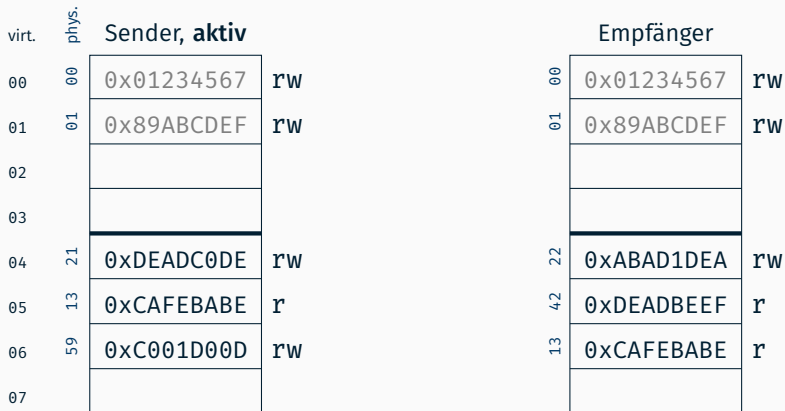


Referenzzähler:

phys.	...	12	13	14	...	21	22	23	...	42	...	57	58	...	
Ref.	...	0	2	0	...	1	1	0	...	1	...	0	0	1	...



# Copy-on-Write



Referenzzähler:

phys.	...	12	13	14	...	21	22	23	...	42	...	57	58	59	...
Ref.	...	0	2	0	...	1	1	0	...	1	...	0	0	1	...



## Umsetzung von Copy-on-Write mittels Pagefaults

- eigener Handler für int 14/0xe
- parallel zu Systemaufruf und Interrupts

Umsetzung von Copy-on-Write mittels Pagefaults

- eigener Handler für `int 14/0xe`
- parallel zu Systemaufruf und Interrupts










## Informationen über Seitenfehler

- Adresse, die Fehler verursacht hat, liegt in `cr2`
- Fehlercode auf dem Kernel-Stack

# Seitenfehler-Behandlung (cont.)

## Fehlercode auf dem Kernel-Stack

Von Hardware nach dem Kontext abgelegt (→ wegräumen!)

31		<b>Reserviert</b>
16		1 falls Software Guard Extensions
14		<b>Reserviert</b>
6		1 falls Memory Protection Key
5		Trap während Data Fetch (0) / Instruction Fetch (1)
4		1 falls reserviertes Bit in Pagingstruktur auf 1 gesetzt
3		Trap im Supervisormodus (0) / Usermodus (1)
2		Trap durch lesen (0) / schreiben (1)
1		Seite abwesend (0) / Schutzfehler (1)
0		

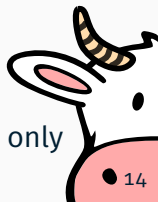
# Seitenfehler-Behandlung (cont.)

## Fehlercode auf dem Kernel-Stack

Von Hardware nach dem Kontext abgelegt (→ wegräumen!)

31		<b>Reserviert</b>
16	SGX	1 falls Software Guard Extensions
14		<b>Reserviert</b>
6		
5	PK	1 falls Memory Protection Key
4	I/D	Trap während Data Fetch (0) / Instruction Fetch (1)
3	RSVD	1 falls reserviertes Bit in Pagingstruktur auf 1 gesetzt
2	U/S	Trap im Supervisormodus (0) / <b>Usermodus (1)</b>
1	W/R	Trap durch lesen (0) / <b>schreiben (1)</b>
0	P	Seite abwesend (0) / <b>Schutzfehler (1)</b>

**Copy-on-Write:** P=1, W/R=1, U/S=1, PT-Eintrag read only



## Neue Prozesse erzeugen

---

- Soll **Copy-on-Write** für gesamten User-Bereich nutzen (auch Stack!)
- Rückgabewert: ID des jeweils anderen
- Neuer Prozess → Wechsel von Ring 0 nach Ring 3  
⇒ Rückgabewert, Instruktions-/Stapel-Zeiger

**Fragen?**