

2026年3月24日 - IJ Lab Seminar

カーネルバイパスによる 通信高速化の基本

IJ 技術研究所 安形

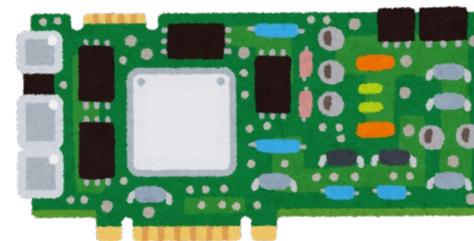
カーネルバイパスとは？

カーネルバイパスとは？

カーネル

カーネルバイパスとは？

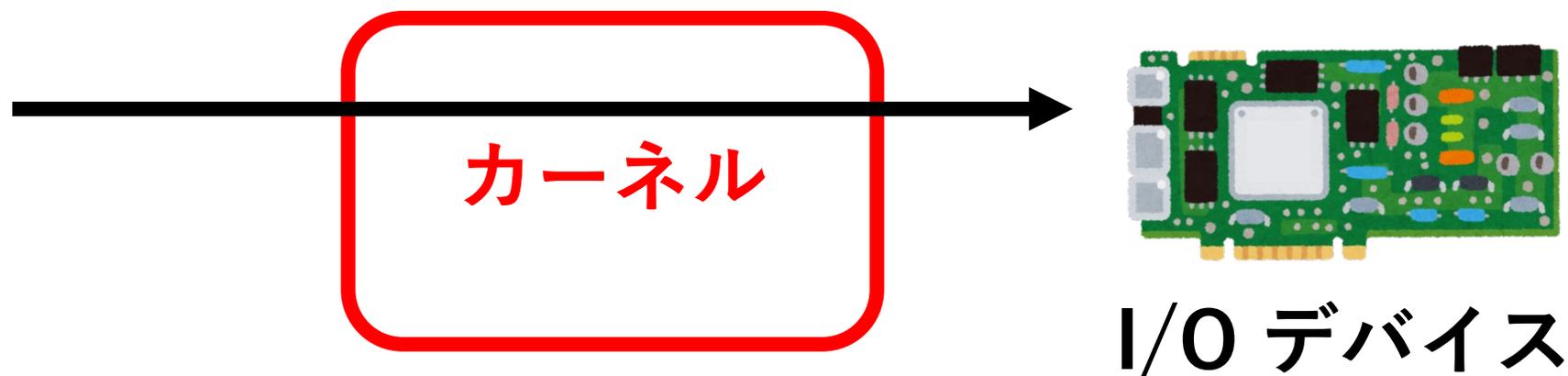
カーネル



I/O デバイス

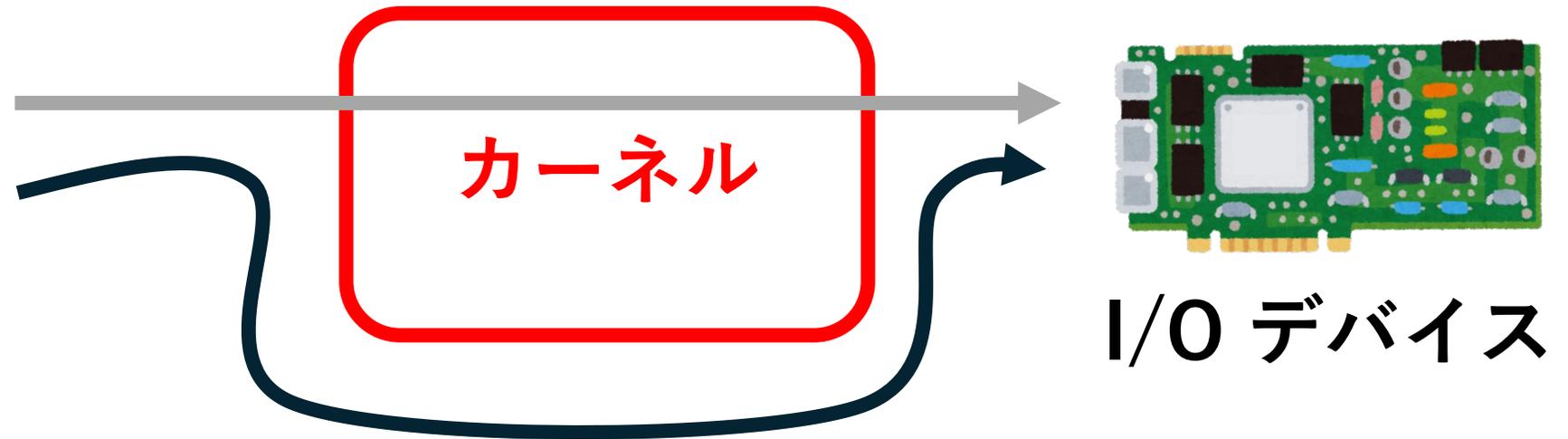
カーネルバイパスとは？

- 基本的にカーネルを経由していたところ

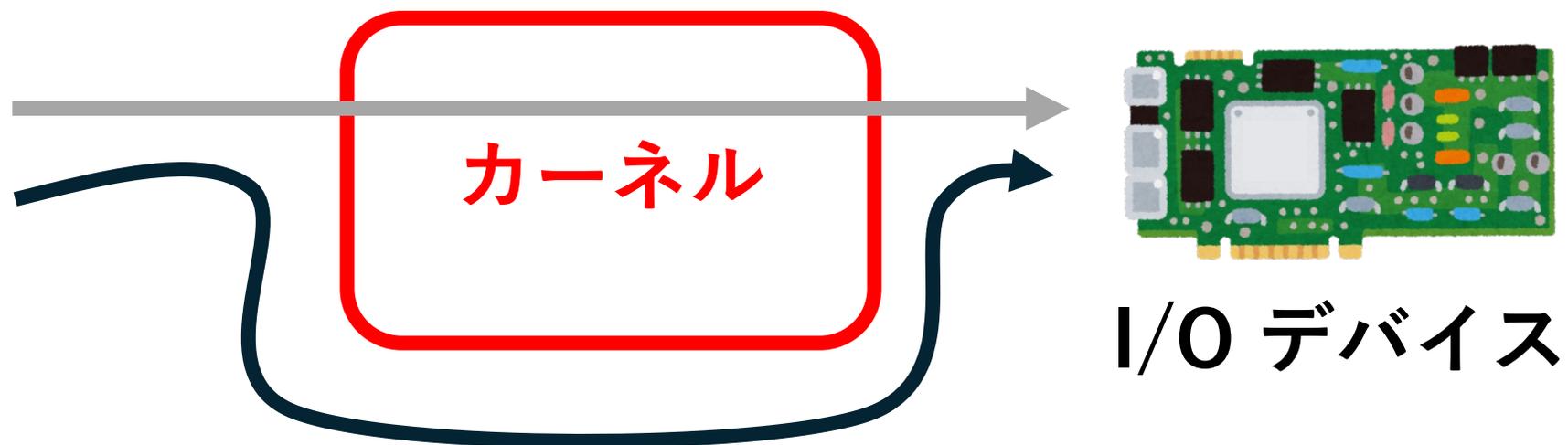


カーネルバイパスとは？

- 基本的にカーネルを経由していたところ、迂回するようにする

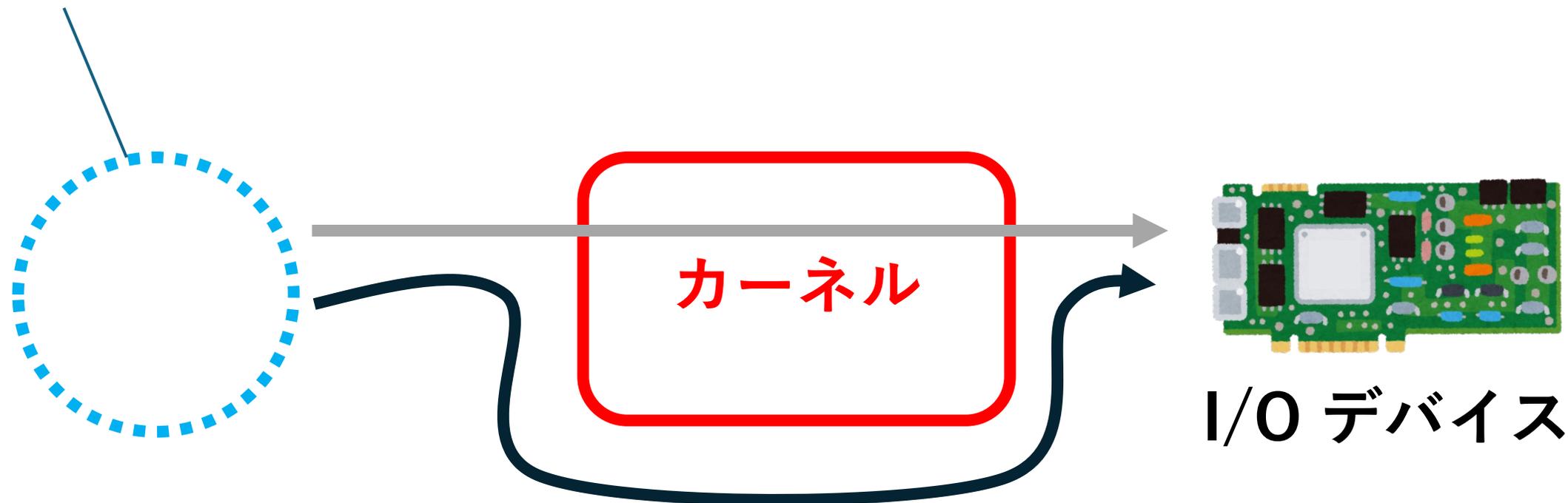


カーネルバイパスとは？



カーネルバイパスとは？

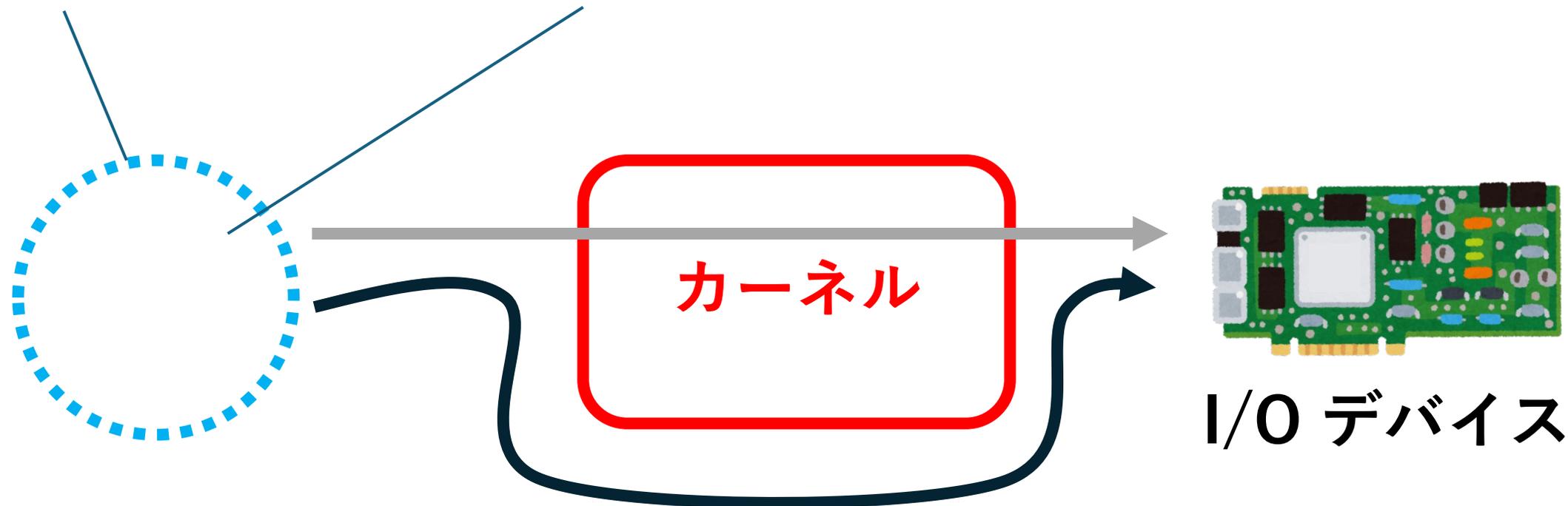
Q. これは何？



カーネルバイパスとは？

Q. これは何？

Q. 何がカーネルを中継しなくなる？

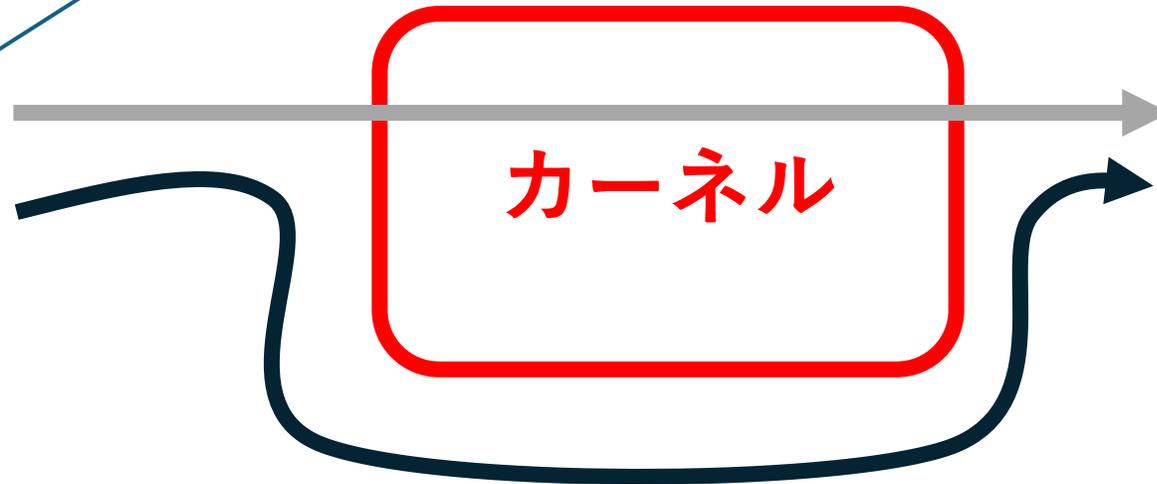
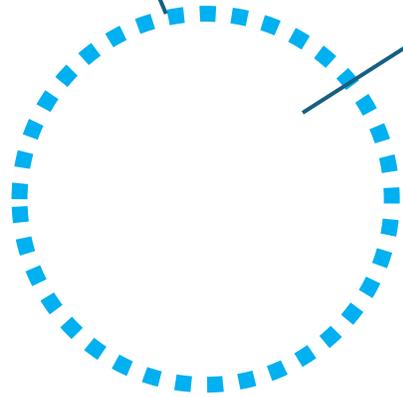


カーネルバイパスとは？

Q. これは何？

Q. 何がカーネルを中継しなくなる？

Q. デバイスへのアクセスとは具体的に何？



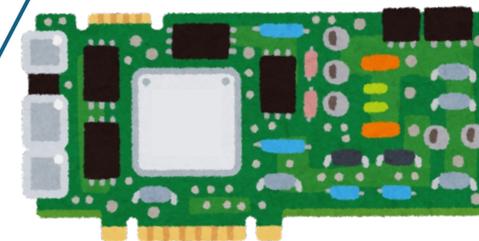
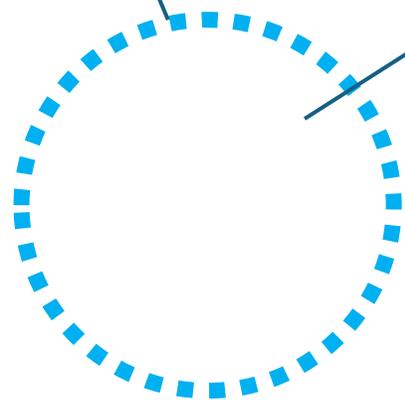
I/O デバイス

カーネルバイパスとは？

Q. これは何？

Q. 何がカーネルを中継しなくなる？

Q. デバイスへのアクセスとは具体的に何？



I/O デバイス

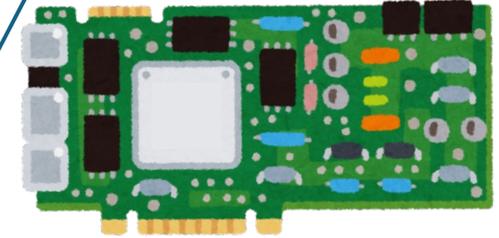
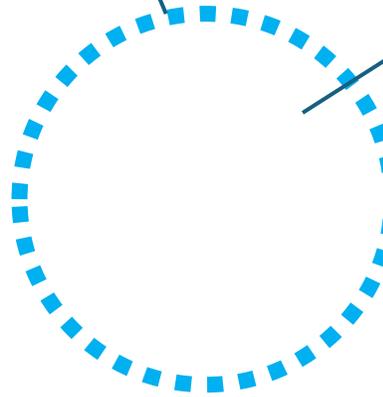
Q. 何故通常はカーネルを経由する？

カーネルバイパスとは？

Q. これは何？

Q. 何がカーネルを中継しなくなる？

Q. デバイスへのアクセスとは具体的に何？



I/O デバイス

Q. 何故通常はカーネルを経由する？

Q. どうやってバイパスする？

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのため

Q. デバイスへの

アクセスとは

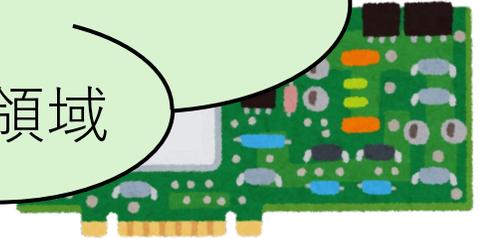
具体的に何？

本能的には

メモリの読み書き

Q. ユーザー空間とは何？

A. CPU が非特権モードで動作している間にアクセス可能なメモリ領域



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセス

Q. デバイスへの

アクセスとは

具体的に何？

本能的には

メモリの読み書き

Q. ユーザー空間とは何？

A. CPU が 非特権モード で動作している間にアクセス可能なメモリ領域



CPU の非特権モードとは？

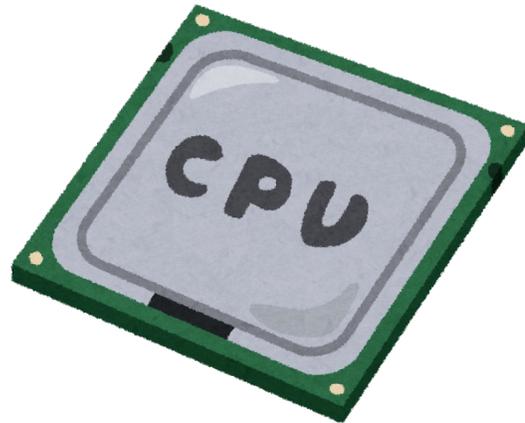
Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

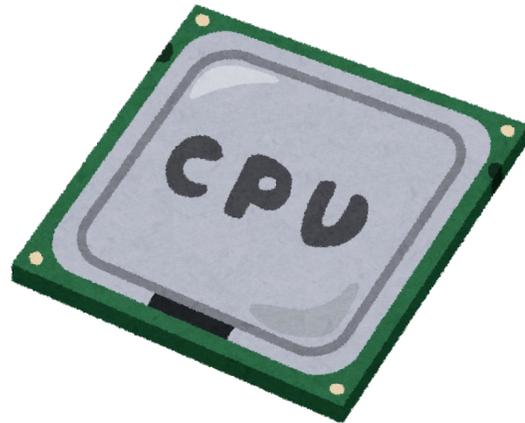
CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）
- 特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）
- 特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されていない場合



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類

1. 特権モード：大体の CPU 機能を使える（システムの破壊が容易）

2. 非特権モード：できることが制限されている（破壊操作が難しい）

特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります

- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されていない場合

このアプリ良さそう！
インストールして
使ってみよう！



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されていない場合

インストールしたアプリが
バグでクラッシュしたせいで
パソコンがフリーズした！



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されていない場合

アプリのクラッシュが
パソコン全体に影響を
与えている

インストールしたアプリが
バグでクラッシュしたせいで
パソコンがフリーズした！



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されている場合

このアプリ良さそう！
インストールして
使ってみよう！



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 - 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 - 2. **非特権モード**：できることが制限されている（破壊操作が難しい）
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります

この方針が適用されている場合

インストールしたアプリが
バグでクラッシュした！けど
パソコンはフリーズしない



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されている場合

アプリのクラッシュの
影響はパソコン全体には
波及しない

インストールしたアプリが
バグでクラッシュした！けど
パソコンはフリーズしない



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）

特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります

- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**
 - ほとんどのプログラムを**非特権モード**で実行（破壊操作が難しい）
 - 限られたプログラムのみを**特権モード**で実行

CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類

1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）

2. **非特権モード**：できることが制限されている（破壊操作が難しい）

特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります

- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

- ほとんどのプログラムを**非特権モード**で実行（破壊操作が難しい）

- 限られたプログラムのみを**特権モード**で実行



- 雑な「カーネル」の定義：**特権モード**で実行されるプログラム

OS の設計によっては必ずしもこの限りではないと思われます

CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて2つ

1. **特権モード**：大体の CPU 機能を使える

2. **非特権モード**：できることが制限される

特権モードはカーネルモード、非特権モードはユーザーモード

非特権モードで動作するアプリがクラッシュしても

パソコンはフリーズしないが

特権モードで動作するカーネルが

完全にクラッシュすると

パソコンはフリーズする

- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

- ほとんどのプログラムを**非特権モード**で実行（破壊操作が難しい）

- 限られたプログラムのみを**特権モード**で実行



- 雑な「カーネル」の定義：**特権モード**で実行されるプログラム

OS の設計によっては必ずしもこの限りではないと思われます

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスを直接行う

Q. デバイスへの

アクセスとは

具体的に何？

基本的には

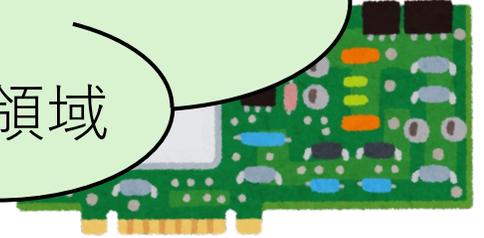
デバイスの読み書き

Q. ユーザー空間とは何？

A. CPU が **非特権モード** で

動作している間に

アクセス可能なメモリ領域



CPU の非特権モードとは？

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセス

Q. デバイスへの

アクセスとは

具体的に何？

本能的には

メモリの読み書き

Q. ユーザー空間とは何？

A. CPU が非特権モードで

動作している間に

アクセス可能なメモリ領域

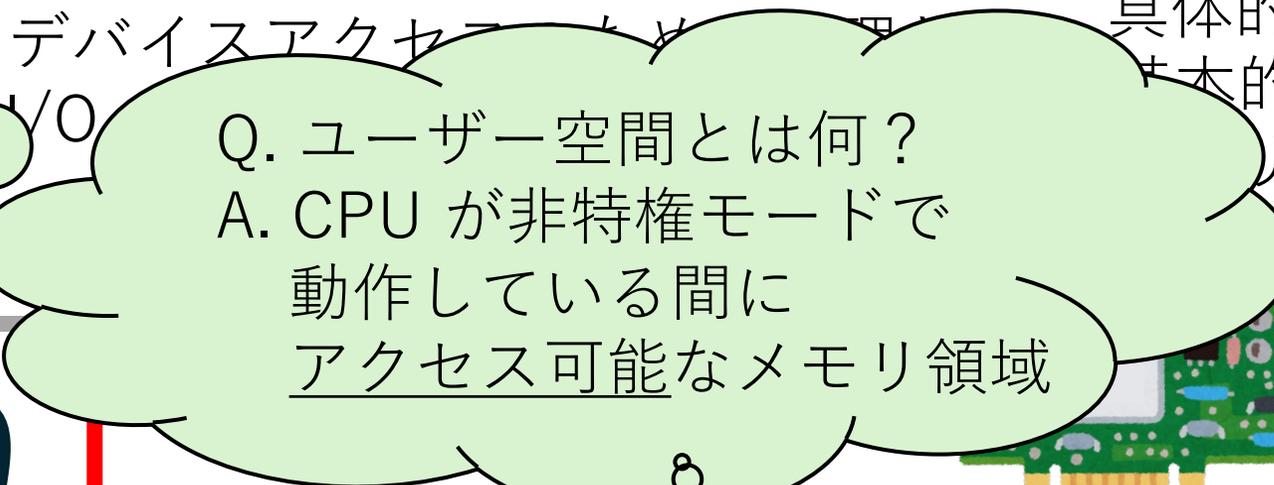
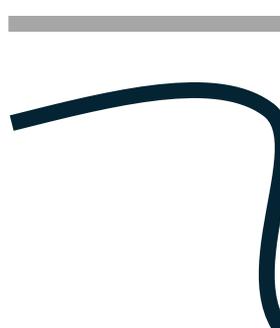
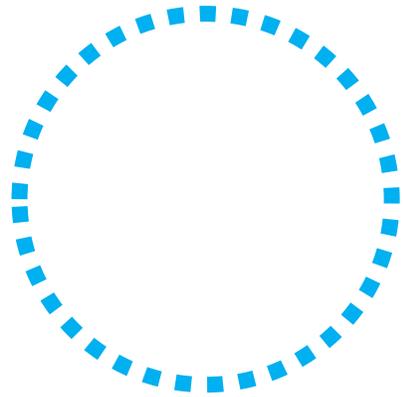
そもそもメモリアクセスとは？

Q. 何故通常はカーネルを経由する？

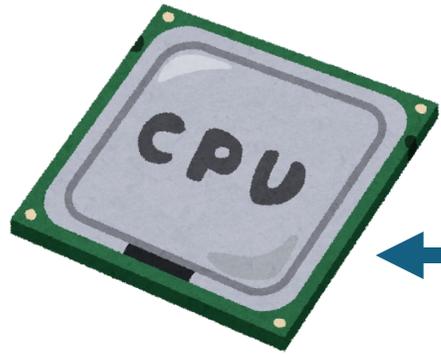
A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

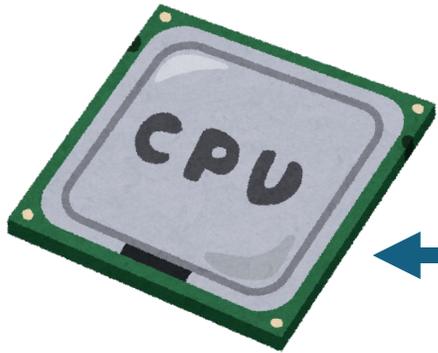


プログラムによるメモリアクセス



プログラムによるメモリアクセス

メモリ書き込みをするには
どんなプログラムを書けば良いか？



プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
    *((int *) 0x100000000) = 0x12345678;
}
```

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

0000000000001129 <main>:

1129:	f3 0f 1e fa	endbr64
112d:	55	push %rbp
112e:	48 89 e5	mov %rsp,%rbp
1131:	48 b8 00 00 00 00 01	movabs \$0x100000000,%rax
1138:	00 00 00	
113b:	c7 00 78 56 34 12	movl \$0x12345678,(%rax)
1141:	b8 00 00 00 00	mov \$0x0,%eax
1146:	5d	pop %rbp
1147:	c3	retq
1148:	0f 1f 84 00 00 00 00	nopl 0x0(%rax,%rax,1)
114f:	00	

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

0000000000001129 <main>:

```
1129: f3 0f 1e fa
112d: 55
112e: 48 89 e5
1131: 48 b8 00 00 00 00 01
1138: 00 00 00
113b: c7 00 00 12
1141: b8 00 00 00 00
1146: 5d
1147: c3
1148: 0f 1f 84 00 00 00 00
114f: 00
```

機械語

```
endbr64
push   %rbp
mov    %rsp,%rbp
movabs $0x100000000,%rax
movl  $0x12345678,%rax
mov    $0x0,%eax
pop    %rbp
retq
nopl  0x0(%rax,%rax,1)
```

アセンブリ言語

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

0000000000001129 <main>:

```
1129:    f3 0f 1e fa    endbr64
112d:    55            push   %rbp
112e:    48 89 e5      mov    %rsp,%rbp
1131:    48 b8 00 00 00 01  movabs $0x100000000,%rax
1138:    00 00 00
113b:    c7 00 78 56 34 12  movl   $0x12345678,(%rax)
1141:    b8 00 00 00 00  mov    $0x0,%eax
1146:    5d            pop    %rbp
1147:    c3            retq
1148:    0f 1f 84 00 00 00 00  nopl   0x0(%rax,%rax,1)
114f:    00
```

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

00000000000001129 <main>:

rax レジスタに 0x100000000 を設定する

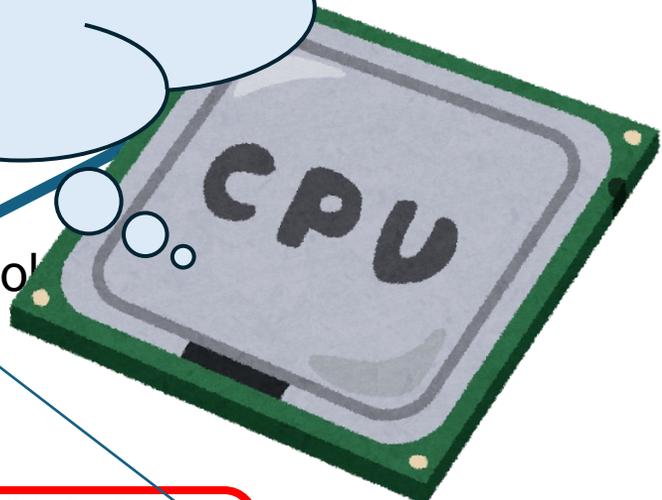
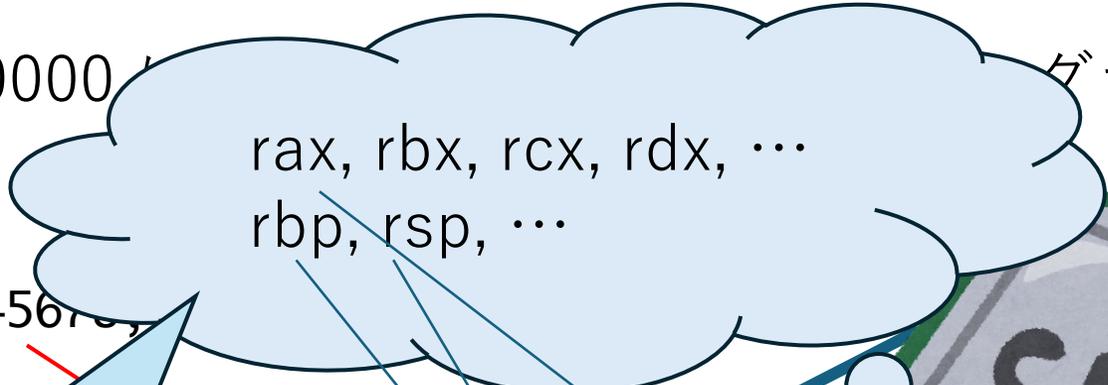
```
1131: 48 b8 00 00 00 00 01
1138: 00 00 00
113b: c7 00 78 56 34 12
1141: b8 00 00 00 00
1146: 5d
1147: c3
1148: 0f 1f 84 00 00 00 00
114f: 00
```

```
endbr64
push   %rbp
mov    %rsp,%rbp
movabs $0x100000000,%rax
movl   $0x12345678,(%rax)
mov    $0x0,%eax
pop    %rbp
retq
nopl   0x0(%rax,%rax,1)
```


プログラムによるメモリアクセス

例：メモリアドレス 0x1000000000

プログラム



Q. レジスタとは？
A. CPU に付属する記憶領域

レジスタそれぞれに名前がついており
操作のための専用の CPU 命令が
CPU に実装されている

rax, rbx, rcx, rdx, ...
rbp, rsp, ...

00000000000000001129 <main>:

rax

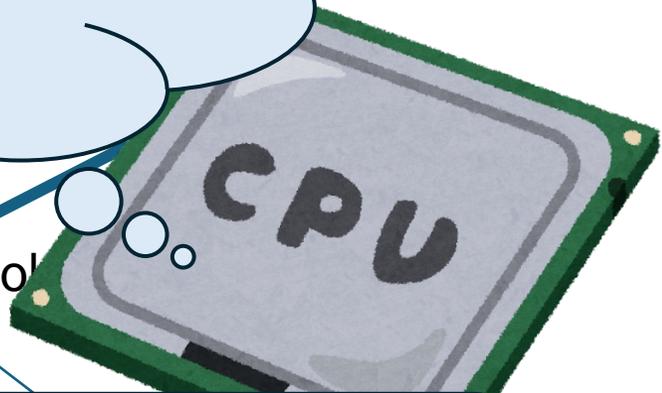
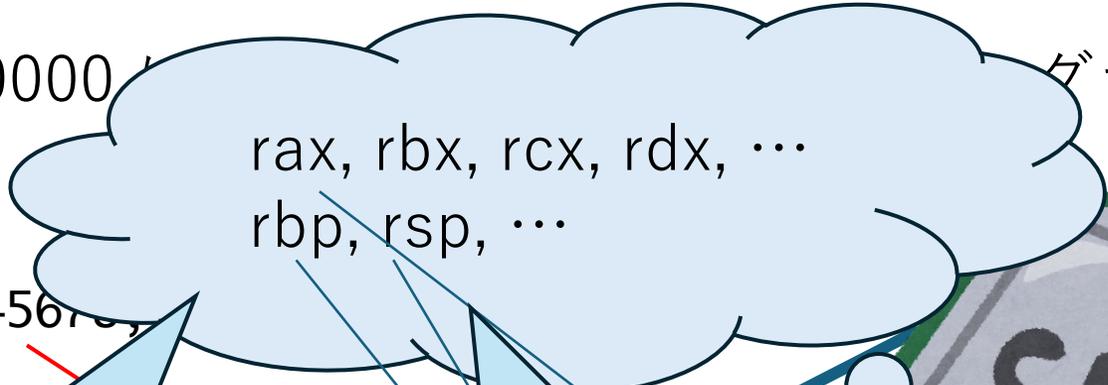
```
endbr64  
push    %rbp  
mov     %rsp,%rbp  
movabs  $0x1000000000,%rax  
movl   $0x12345678,(%rax)  
mov     $0x0,%eax  
pop     %rbp  
retq  
nopl   0x0(%rax,%rax,1)
```

```
1146: 5d  
1147: c3  
1148: 0f 1f 84 00 00 00 00  
114f: 00
```

プログラムによるメモリアクセス

例：メモリアドレス 0x1000000000

プログラム



Q. レジスタとは？
A. CPU に付属する記憶領域

レジスタそれぞれに名前がついており
操作のための専用の CPU 命令が
CPU に実装されている

ポイント

- CPU 付属の高速な記憶領域
- CPU 命令から直接操作可能

i
{
}
}

0000000000001129 <main>:

rax

endbr6

```
1146: 5d          pop     %rbp
1147: c3          retq
1148: 0f 1f 84 00 00 00 00  nopl   0x0(%rax,%rax,1)
114f: 00
```

```
pop     %rbp
retq
nopl   0x0(%rax,%rax,1)
```

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

00000000000001129 <main>:

rax レジスタに 0x100000000 を設定する

```
1131: 48 b8 00 00 00 00 01
1138: 00 00 00
113b: c7 00 78 56 34 12
```

rax レジスタの値により参照される
メモリアドレス (0x100000000) へ
0x12345678 を書き込む

```
endbr64
push  %rbp
mov   %rsp,%rbp
movabs $0x100000000,%rax
movl  $0x12345678,(%rax)
mov   $0x0,%eax
pop   %rbp
retq
nopl  0x0(%rax,%rax,1)
```

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

00000000000001129 <main>:

rax レジスタに 0x100000000 を設定する

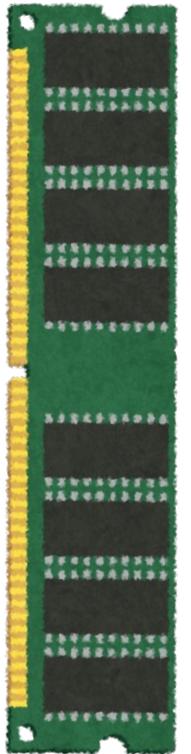
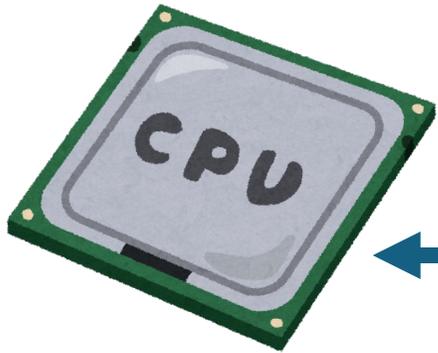
```
1131: 48 b8 00 00 00 00 01
1138: 00 00 00
113b: c7 00 78 56 34 12
```

rax レジスタの値により参照される
メモリアドレス (0x100000000) へ
0x12345678 を書き込む

```
endbr64
push   %rbp
mov    %rsp,%rbp
movabs $0x100000000,%rax
movl   $0x12345678,(%rax)
mov    $0x0,%eax
pop    %rbp
retq
nopl   0x0(%rax,%rax,1)
```

プログラムによるメモリアクセス

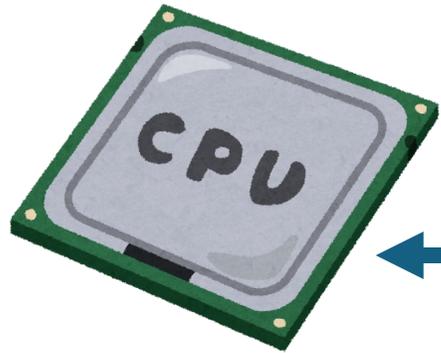
メモリ書き込みをするには
どんなプログラムを書けば良いか？



プログラムによるメモリアクセス

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



この命令の実行時に CPU から
メモリアクセスが試みられる

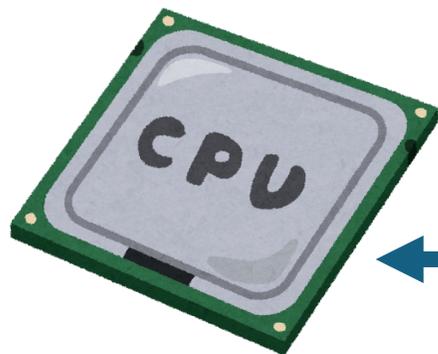
rax:0x100000000



プログラムによるメモリアクセス

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000

この命令の実行時に CPU から
メモリアクセスが試みられる

ポイント

- CPU は以下のような **CPU 命令** を実装している
- CPU レジスタの値を読み書きする
 - メモリのアドレスを指定して読み書きを行う



カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセス

Q. デバイスへの

アクセスとは

具体的に何？

本能的には

メモリの読み書き

Q. ユーザー空間とは何？

A. CPU が非特権モードで

動作している間に

アクセス可能なメモリ領域

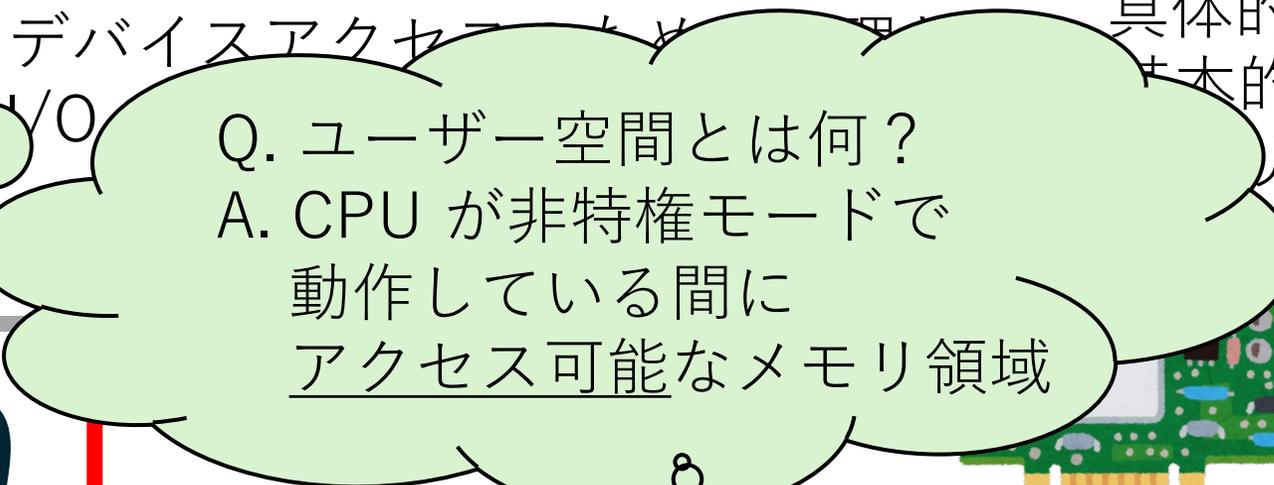
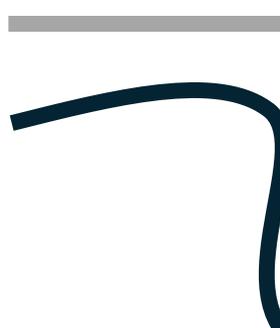
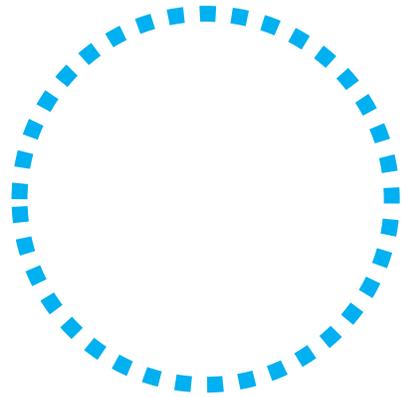
そもそもメモリアクセスとは？

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する



カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスを直接行う

Q. デバイスへの

アクセスとは

具体的に何？

基本的には

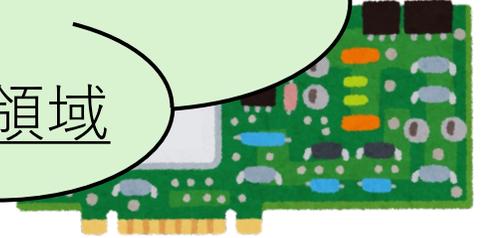
デバイスの読み書き

Q. ユーザー空間とは何？

A. CPU が非特権モードで

動作している間に

アクセス可能なメモリ領域



デバイス

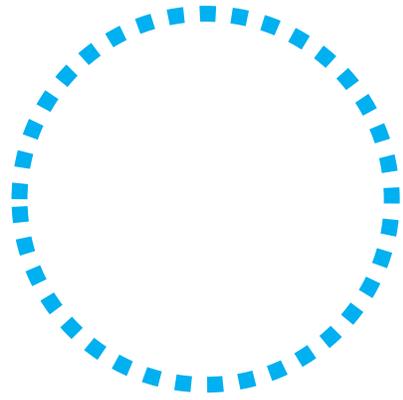
アクセスの可否？

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってカーネルバイパスする？

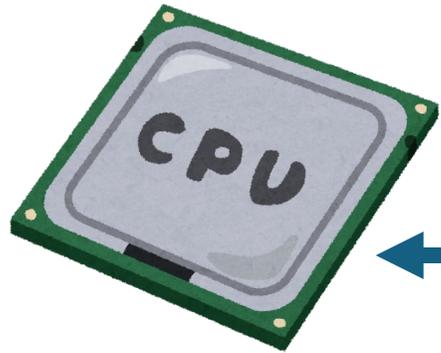
A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する



プログラムによるメモリアクセス

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000

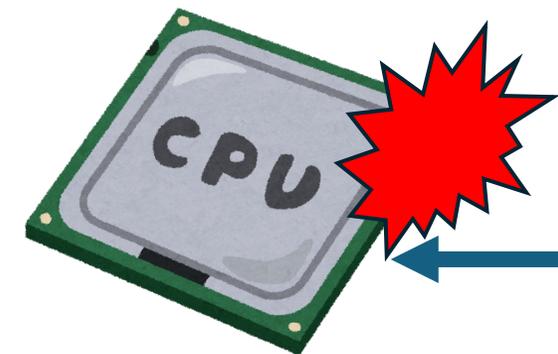
この命令の実行時に CPU から
メモリアクセスが試みられる



プログラムによるメモリアクセス

0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000

実行するとおそらく
Segmentation fault という
エラーで停止して
メモリアクセスができない

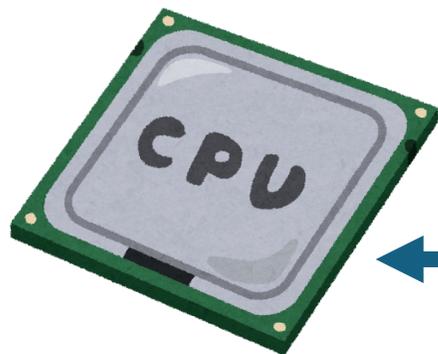
この命令の実行時に CPU から
メモリアクセスが試みられる



プログラムによるメモリアクセス

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000

実行するとおそらく
Segmentation fault という
エラーで停止して
メモリアクセスができない

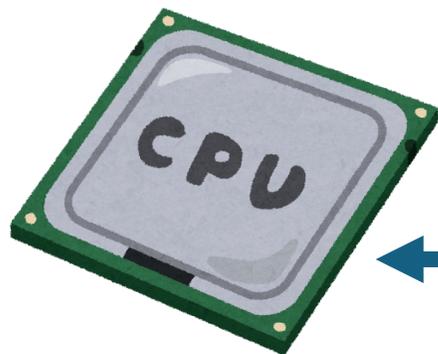
何故? : メモリアクセスが制限されているから



プログラムによるメモリアクセス

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000

実行するとおそらく
Segmentation fault という
エラーで停止して
メモリアクセスができない

何故? : メモリアクセスが制限されているから

何によって? : カーネル
どうやって? : **MMU** の機能を利用

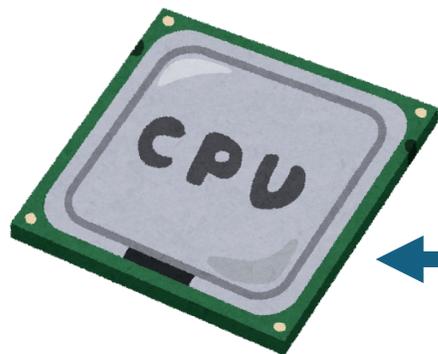


プログラムによるメモリアクセス

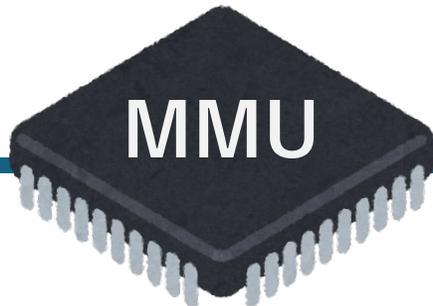
CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000



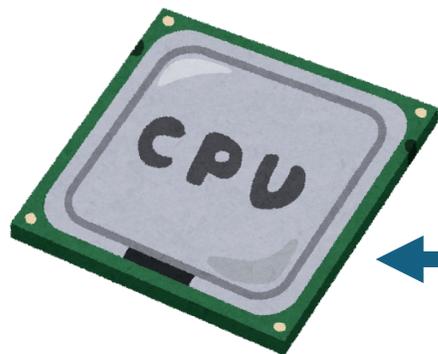
プログラムによるメモリアクセス

CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

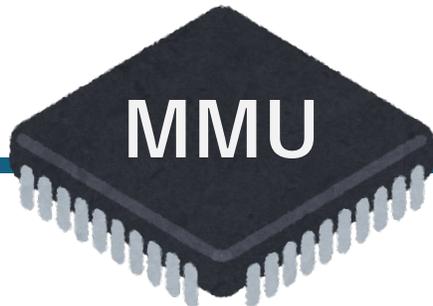
MMU は **仮想**アドレスから**物理**アドレスへの変換を担当

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000



プログラムによるメモリアクセス

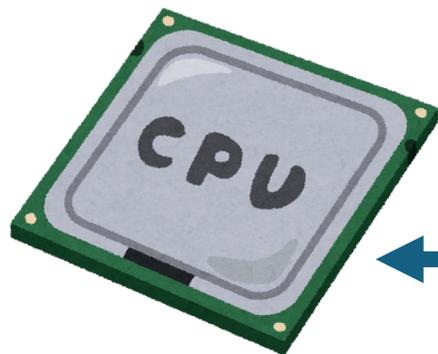
CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

MMU は **仮想**アドレスから**物理**アドレスへの変換を担当
ポイント

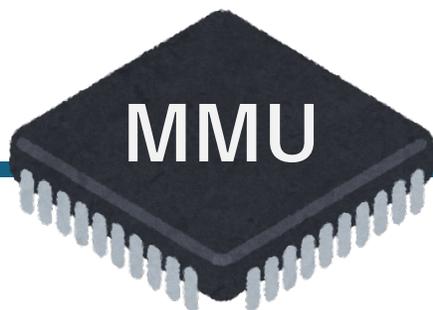
一般的なプログラムが参照する
メモリアドレスは**仮想**メモリアドレス

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000



プログラムによるメモリアクセス

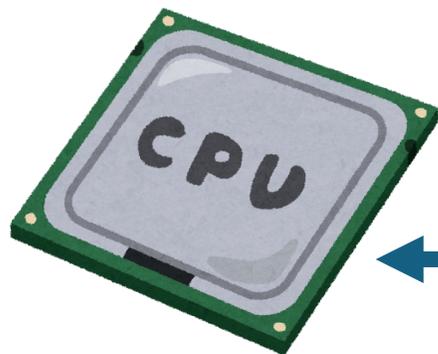
CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

MMU は **仮想**アドレスから**物理**アドレスへの変換を担当
ポイント

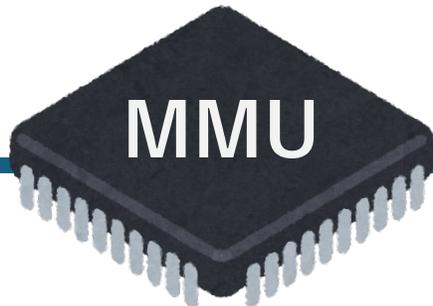
一般的なプログラムが参照する
メモリアドレスは**仮想**メモリアドレス

OS が起動時に初期化の早い段階で
仮想メモリアドレスを基本とする
CPU モードを有効にする

0x100000000 ~
0x12345678 を書き込み
`movl $0x12345678, (%rax)`



`rax:0x100000000`



ちなみに
カーネルも**仮想**メモリアドレスを
基本として動作しています

プログラムによるメモリアクセス

CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

MMU は **仮想**アドレスから**物理**アドレスへの変換を担当

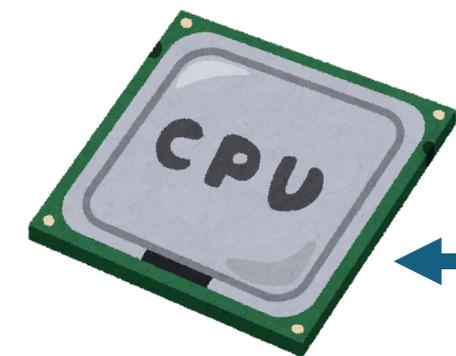
0x100000000 へ
0x12345678 を書き込み

ポイント

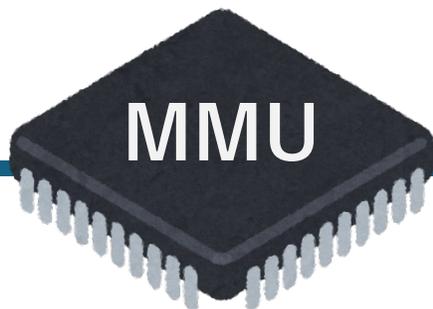
一般的なプログラムが参照する
メモリアドレスは**仮想**メモリアドレス

```
movl $0x12345678, (%rax)
```

OS が起動時に初期化の早い段階で
仮想メモリアドレスを基本とする
CPU モードを有効にする



rax:0x100000000



ちなみに

カーネルも**仮想**メモリアドレスを
基本として動作しています



プログラムによるメモリアクセス

CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

MMU は **仮想**アドレスから**物理**アドレスへの変換を担当

ポイント

一般的なプログラムが参照する
メモリアドレスは**仮想**メモリアドレス

OS が起動時に初期化の早い段階で
仮想メモリアドレスを基本とする
CPU モードを有効にする

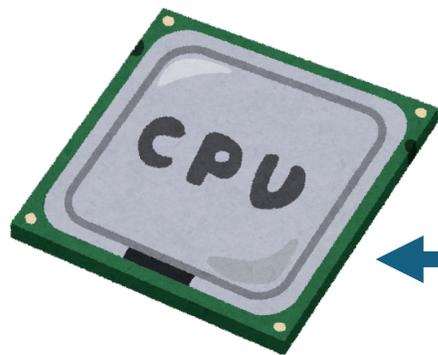
ちなみに

カーネルも**仮想**メモリアドレスを
基本として動作しています

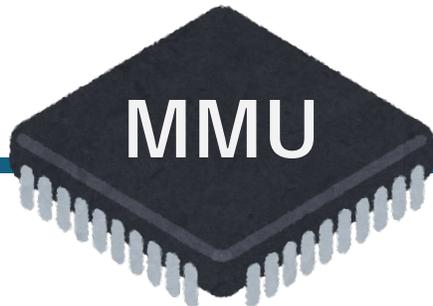


仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000



プログラムによるメモリアクセス

CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

MMU は **仮想**アドレスから**物理**アドレスへの変換を担当

仮想アドレス 0x1000000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



プログラムによるメモリアクセス

CPU からメモリへのアクセスには
MMU (Memory Management Unit) という部品が介在する

MMU は **仮想**アドレスから**物理**アドレスへの変換を担当

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
ページテーブルを参照する

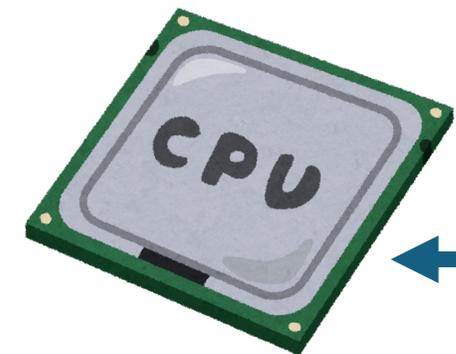
仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

プログラムによるメモリアクセス

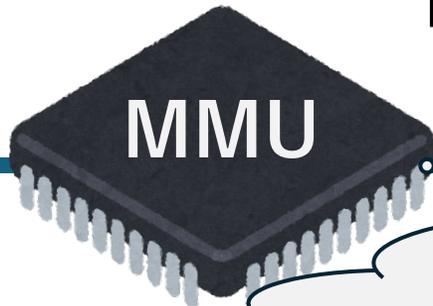
ページテーブルはソフトウェアによりメモリ上に用意される**仮想**アドレスと**物理**アドレスの対応を保持するテーブル
(テーブルのフォーマットはハードウェア依存)

仮想アドレス 0x1000000000 ~
0x12345678 を書き込み

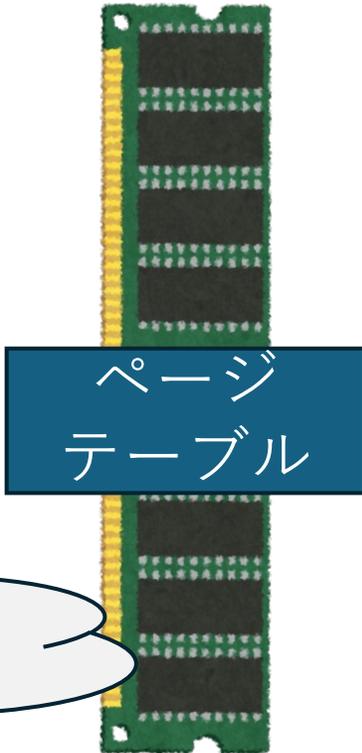
```
movl $0x12345678, (%rax)
```



rax:0x1000000000



MMU はアドレス変換に際して
ページテーブルを参照する



仮想アドレス 0x1000000000 は
物理アドレスのどこだろう？

プログラムによるメモリアクセス

ページテーブルはソフトウェアによりメモリ上に用意される**仮想**アドレスと**物理**アドレスの対応を保持するテーブル
(テーブルのフォーマットはハードウェア依存)

仮想アドレス 0x1000000000 ~ 0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

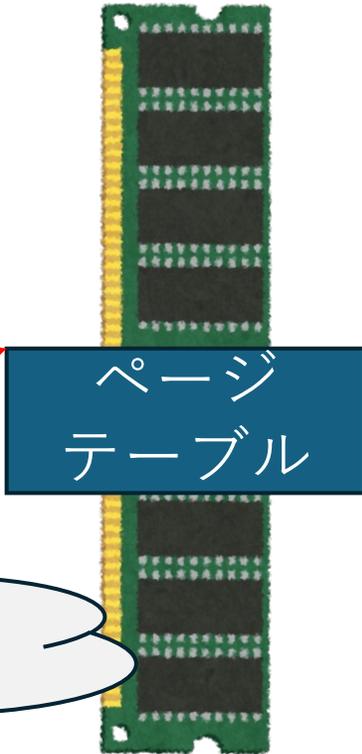
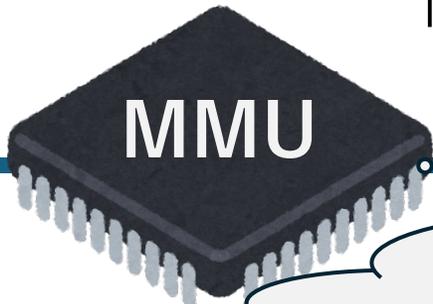
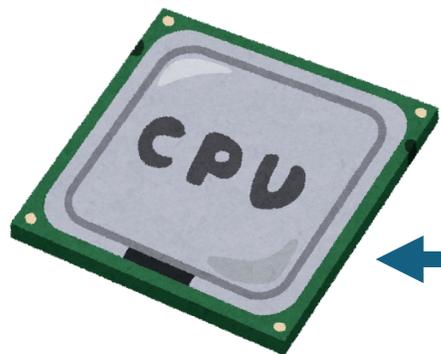
MMU が参照すべきページテーブルは cr3 レジスタを通してソフトウェアにより設定される

MMU はアドレス変換に際してページテーブルを参照する

仮想アドレス 0x1000000000 は物理アドレスのどこだろう？

rax:0x1000000000

cr3: ページテーブルの物理アドレス



プログラムによるメモリアクセス

ページテーブルはソフトウェアによりメモリ上に用意される**仮想**アドレスと**物理**アドレスの対応を保持するテーブル
(テーブルのフォーマットはハードウェア依存)

仮想アドレス 0x1000000000 ~ 0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

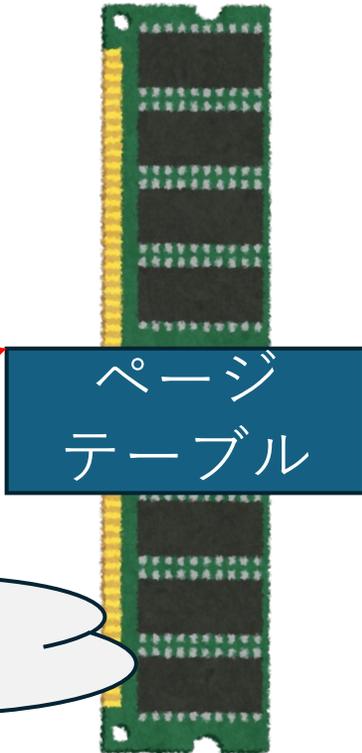
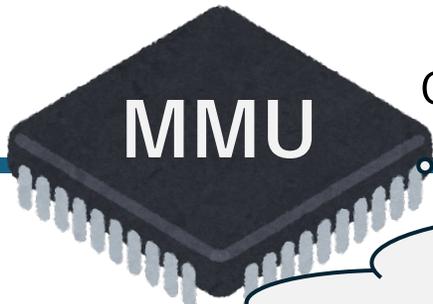
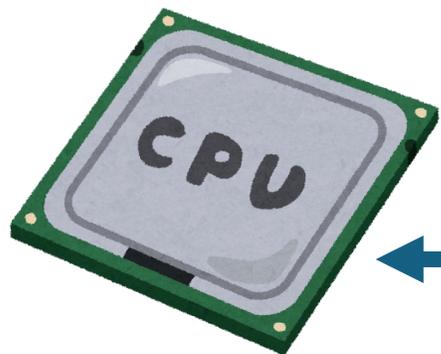
MMU が参照すべきページテーブルは cr3 レジスタを通してソフトウェアにより設定される

MMU はアドレス変換に際して cr3 で示されるページテーブルを参照

仮想アドレス 0x1000000000 は物理アドレスのどこだろう？

rax:0x1000000000

cr3: ページテーブルの物理アドレス



プログラムによるメモリアクセス

ポイント

ページテーブル自体はソフトウェアが用意できる

ページテーブルはソフトウェアによりメモリ上に用意される**仮想**アドレスと**物理**アドレスの対応を保持するテーブル
(テーブルのフォーマットはハードウェア依存)

仮想アドレス 0x1000000000 ~ 0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

MMU が参照すべき

ページテーブルは

cr3 レジスタを通して

ソフトウェアにより設定される

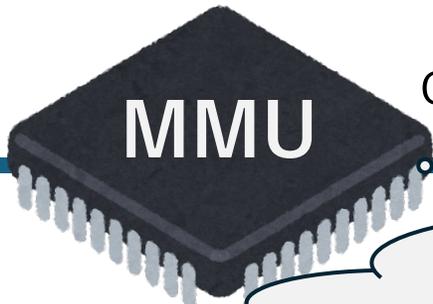
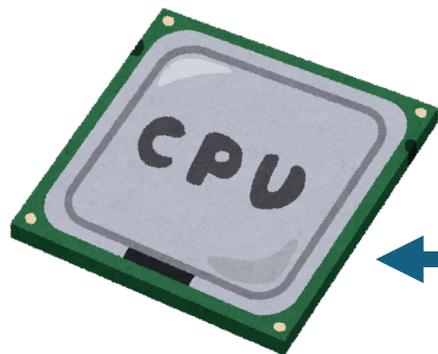
MMU はアドレス変換に際して cr3 で示されるページテーブルを参照

ページ
テーブル

仮想アドレス 0x1000000000 は物理アドレスのどこだろう？

rax:0x1000000000

cr3: ページテーブルの物理アドレス

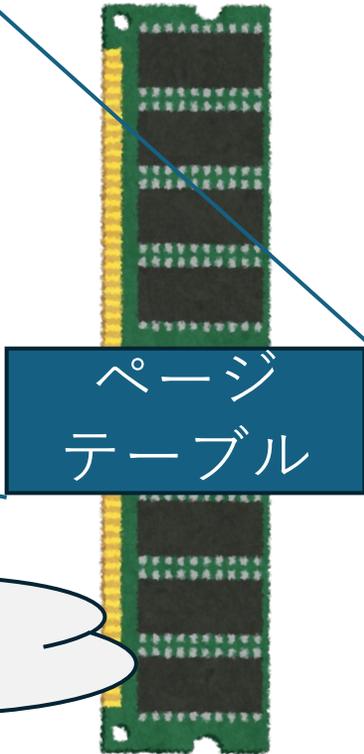
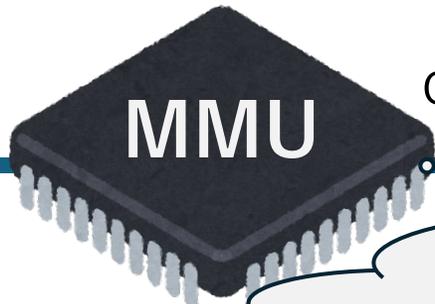
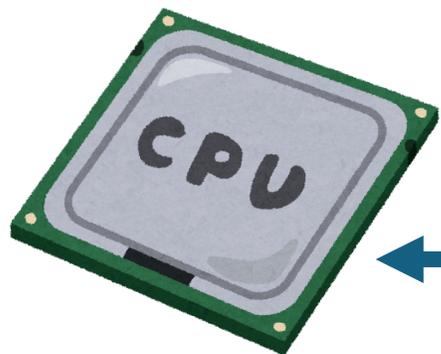


プログラムによるメモリアクセス

仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

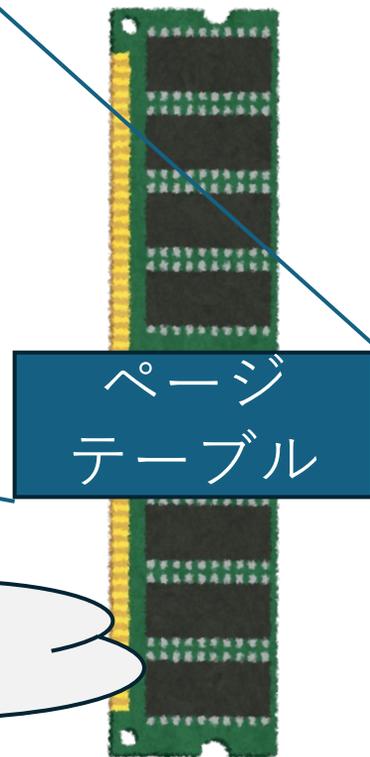
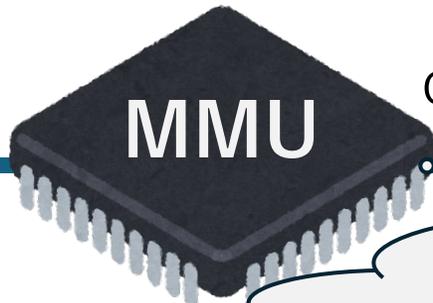
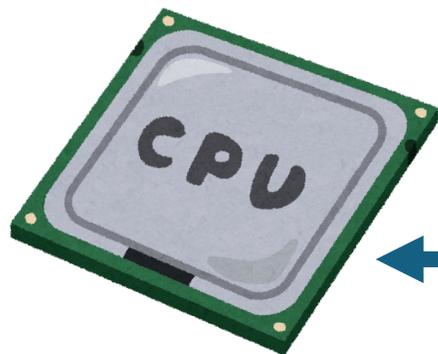
プログラムによるメモリアクセス

仮想アドレスと物理アドレスの
対応の設定は 4 KB ごとに設定
(4 KB は 0x1000 B)

仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

プログラムによるメモリアクセス

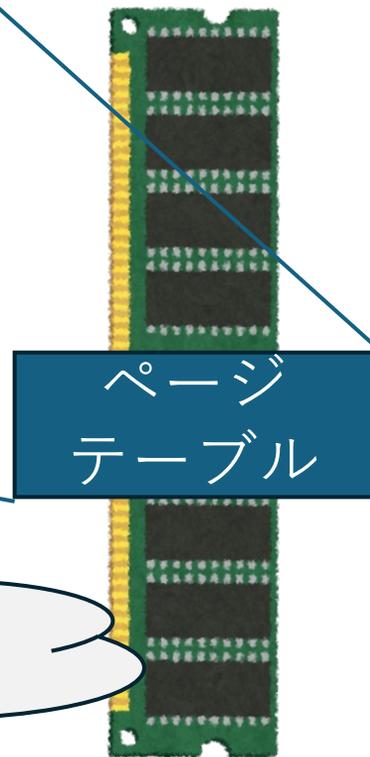
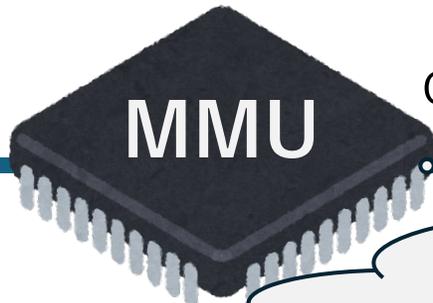
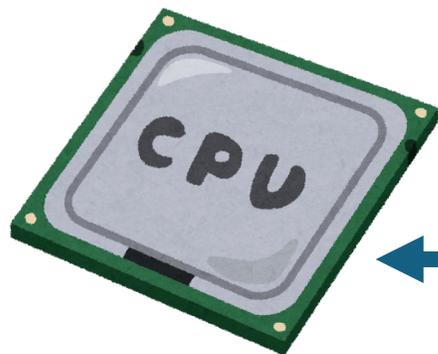
仮想アドレスと物理アドレスの
対応の設定は 4 KB ごとに設定
(4 KB は 0x1000 B)

仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス空間は設定と
ハードウェアに依存しますが
0 ~ 256 TB や 0 ~ 128 PB が
多いかもしれません

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

256 TB や 128 PB まで

仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

プログラムによるメモリアクセス

仮想アドレスと物理アドレスの
対応の設定は 4 KB ごとに設定
(4 KB は 0x1000 B)

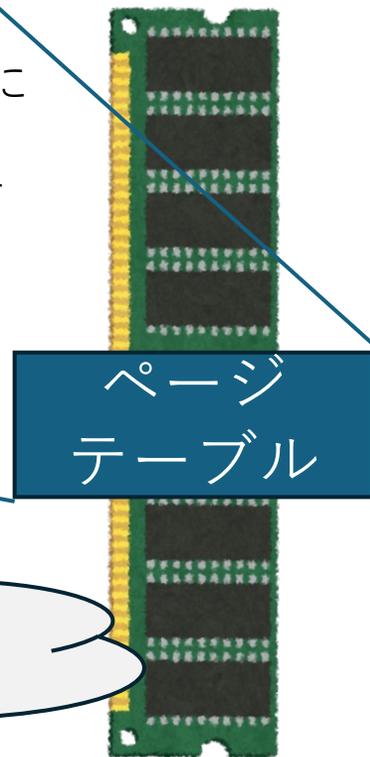
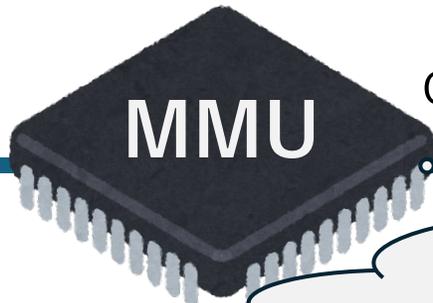
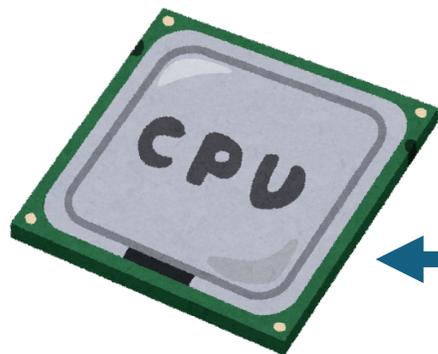
仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス空間は設定と
ハードウェアに依存しますが
0 ~ 256 TB や 0 ~ 128 PB が
多いかもしれません

仮想アドレス空間の全体に
対応する物理アドレスを
設定する必要はないです

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

プログラムによるメモリアクセス

仮想アドレスと物理アドレスの

仮想

物理

仮想アドレス空間は設定と
アに依存しますが
や 0 ~ 128 PB が
おられません

ページテーブルのフォーマットはハードウェア依存ですが
x86-64 では 4 KB のページを物理メモリアドレスの参照を
通して接続した木構造のようなデータ構造になっています

cr3 レジスタに設定するのは木構造の root として扱う
ページの物理メモリアドレスです

より詳細な説明はこちらをご参照ください

<https://yasukata.hatenablog.com/entry/2023/04/10/085714>

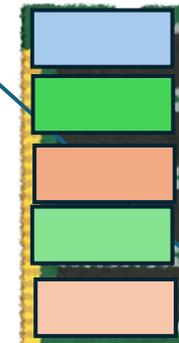
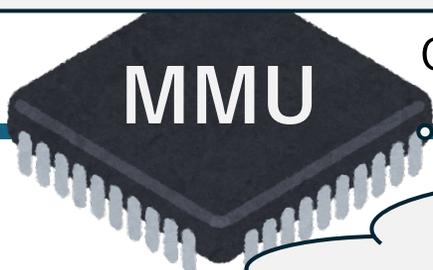
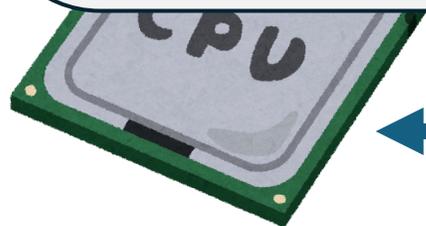
cr3 で示されるページテーブルを参照

ページ
テーブル

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

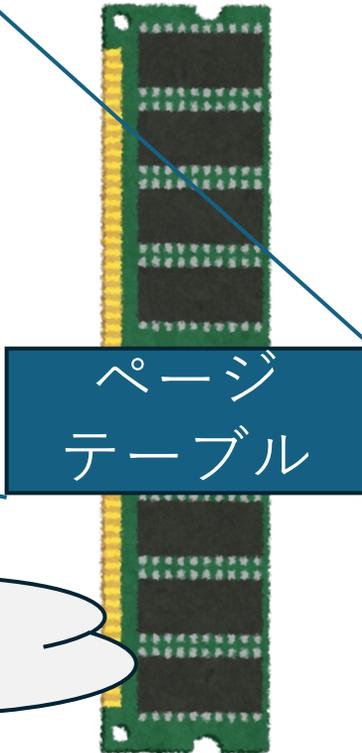
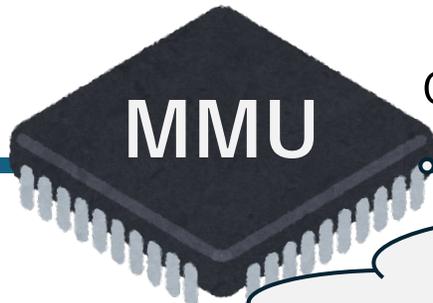
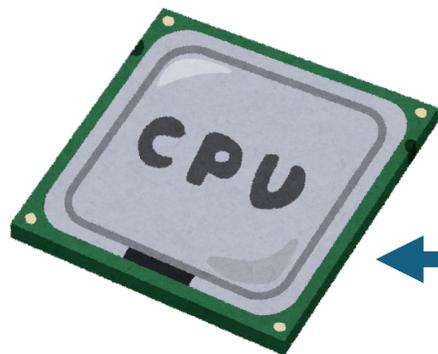


プログラムによるメモリアクセス

仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス 0x100000000 ~ 0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して cr3 で示されるページテーブルを参照

rax: 0x100000000
cr3: ページテーブルの物理アドレス

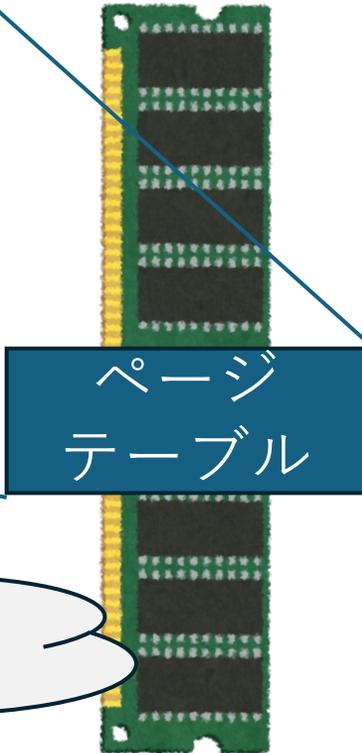
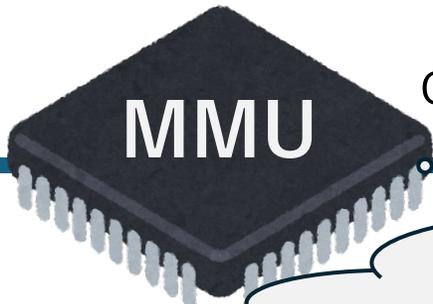
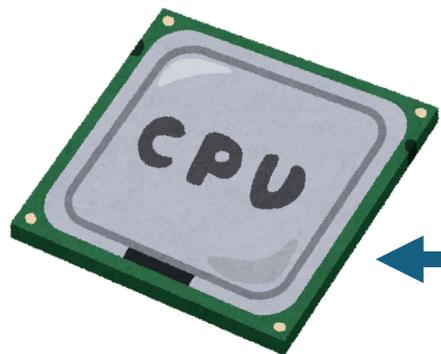
仮想アドレス 0x100000000 は物理アドレスのどこだろう？

プログラムによるメモリアクセス

仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

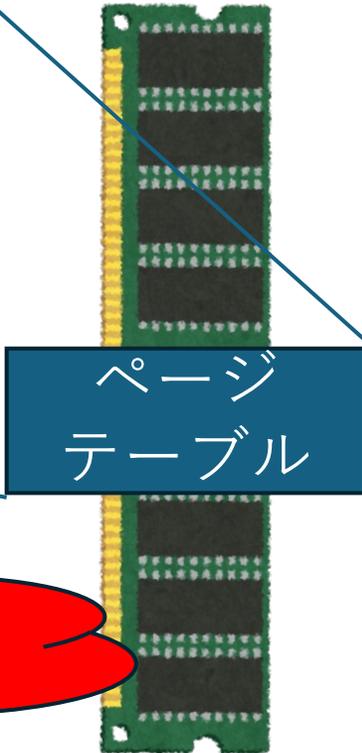
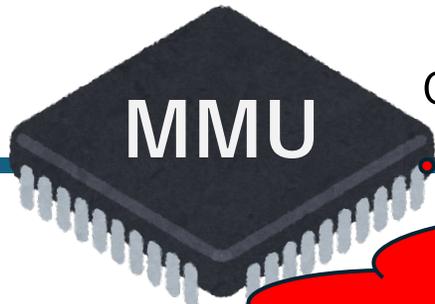
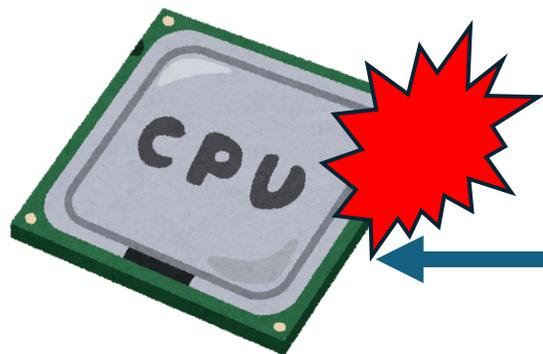
仮想アドレス 0x100000000 と
対応する物理アドレスがない

プログラムによるメモリアクセス

仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

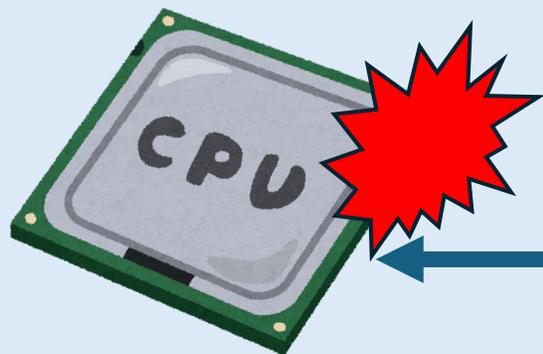
cr3: ページテーブルの物理アドレス

Segmentation fault を起動

プログラムによるメモリアクセス

0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



rax:0x100000000

実行するとおそらく
Segmentation fault という
エラーで停止して
メモリアクセスができない

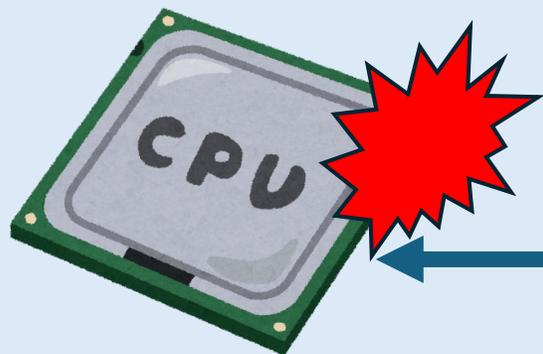
この命令の実行時に CPU から
メモリアクセスが試みられる



プログラムによるメモリアクセス

0x100000000 ~
0x12345678 を書き込み

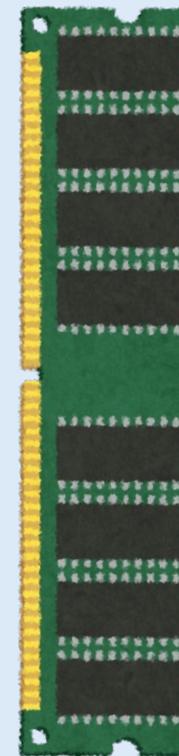
```
movl $0x12345678, (%rax)
```



rax:0x100000000

実行するとおそらく
Segmentation fault という
エラーで停止して
メモリアクセスができない

この命令の実行時に CPU から
メモリアクセスが試みられる

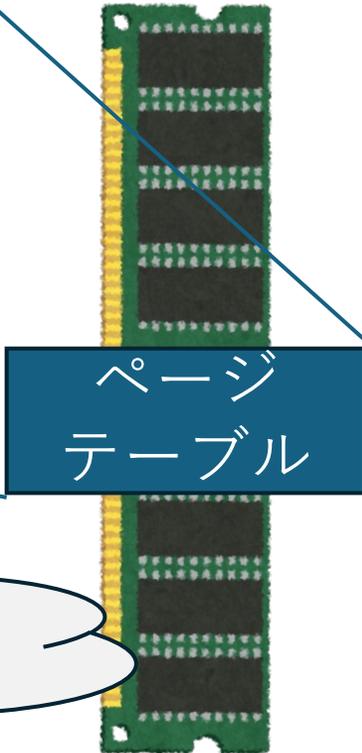
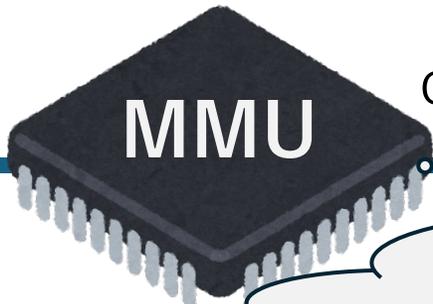
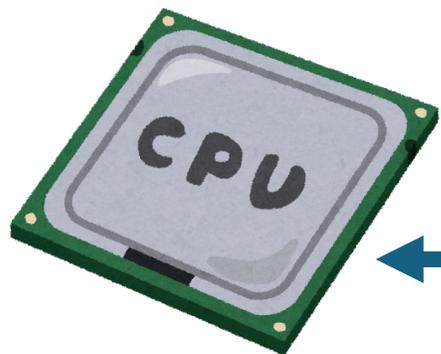


プログラムによるメモリアクセス

仮想	物理
0x0000	
0x1000	
...	
0x100000000	
...	

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

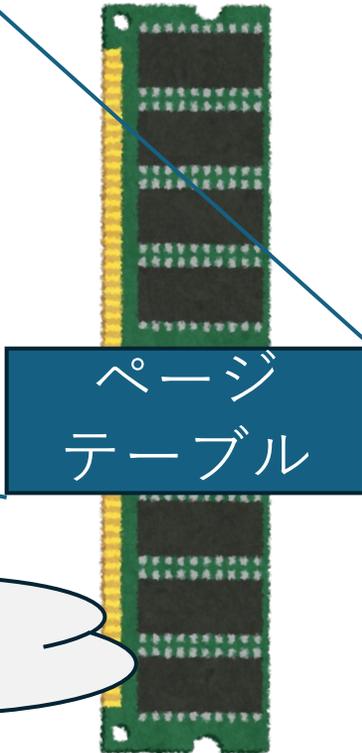
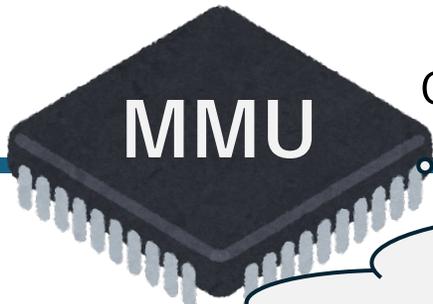
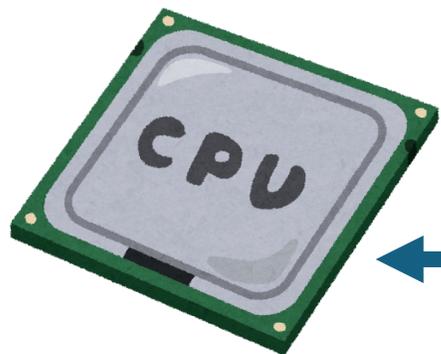
プログラムによるメモリアクセス

対応する**物理アドレス**として
0x2000 が設定されていた場合

仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

仮想	物理
0x0000	
0x1000	
...	
0x100000000	0x2000
...	



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスのどこだろう？

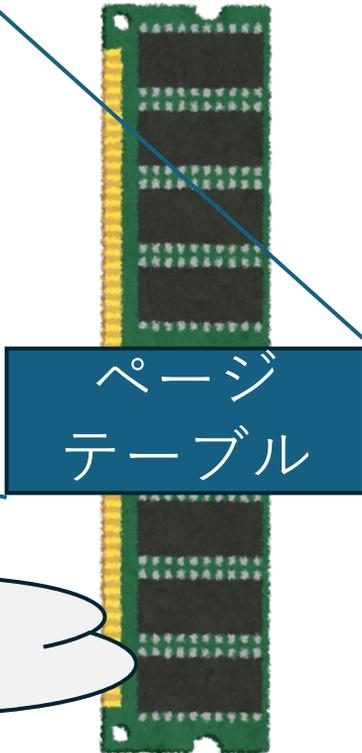
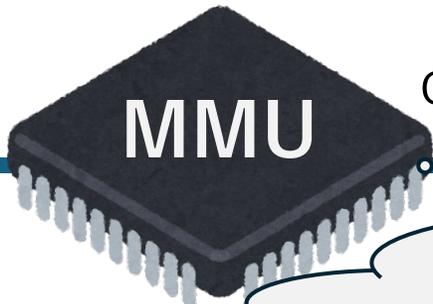
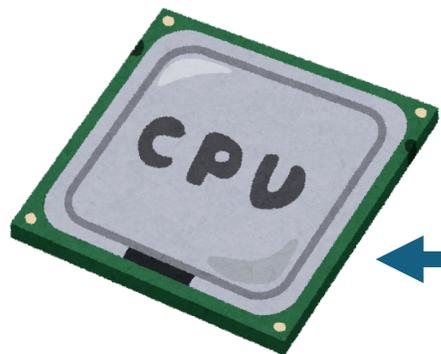
プログラムによるメモリアクセス

対応する**物理アドレス**として
0x2000 が設定されていた場合

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

仮想	物理
0x0000	
0x1000	
...	
0x100000000	0x2000
...	



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000
cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000

プログラムによるメモリアクセス

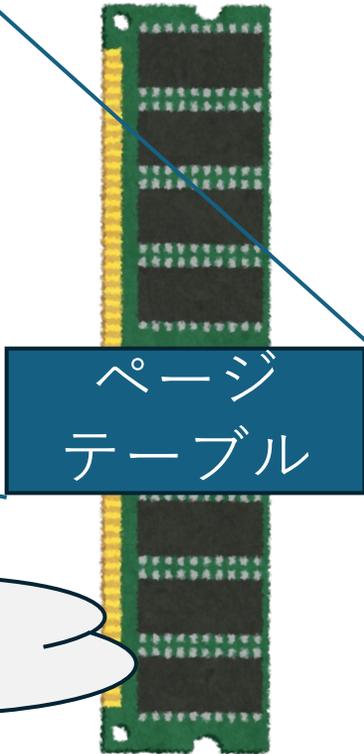
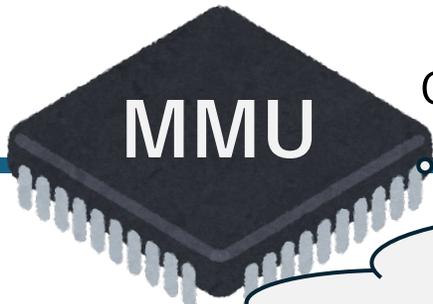
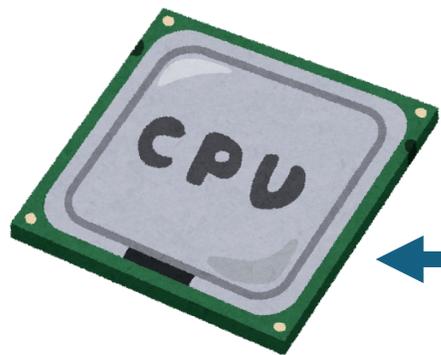
対応する**物理アドレス**として
0x2000 が設定されていた場合

仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

仮想	物理
0x0000	
0x1000	
0x100000000	0x2000
...	

物理アドレス 0x2000 へ
0x12345678 を書き込み



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000
cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000

プログラムによるメモリアクセス

対応する**物理アドレス**として
0x2000 が設定されていた場合

仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

仮想	物理
0x0000	
0x1000	
0x100000000	0x2000
...	

物理アドレス 0x2000 へ
0x12345678 を書き込み

0x2000 0x12345678

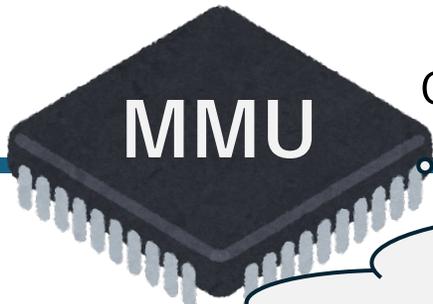
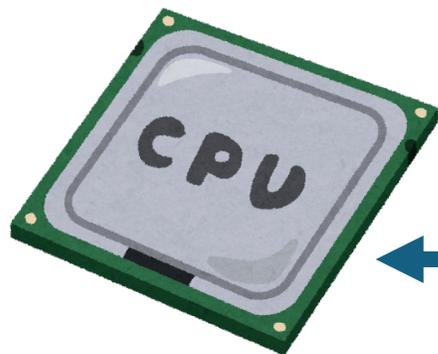
ページ
テーブル

MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

仮想アドレス 0x100000000 は
物理アドレスの 0x2000

rax:0x100000000

cr3: ページテーブルの物理アドレス

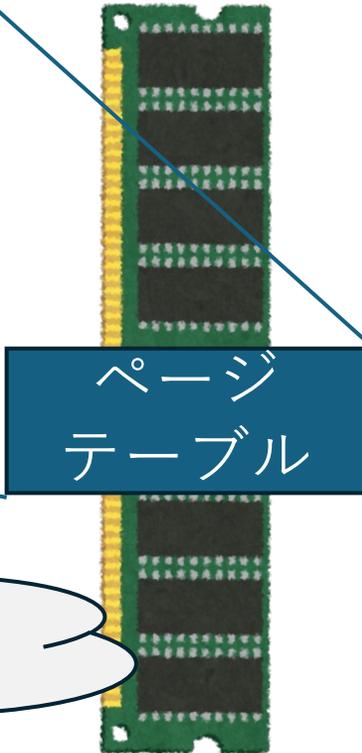
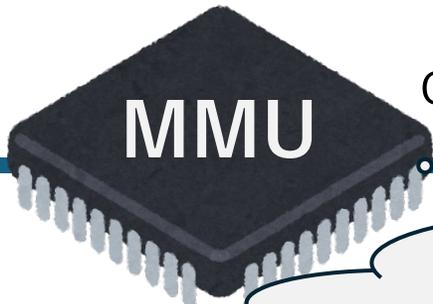
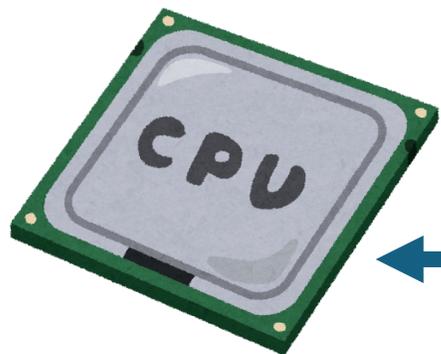


プログラムによるメモリアクセス

仮想	物理
0x0000	
0x1000	
...	
0x100000000	0x2000
...	

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000

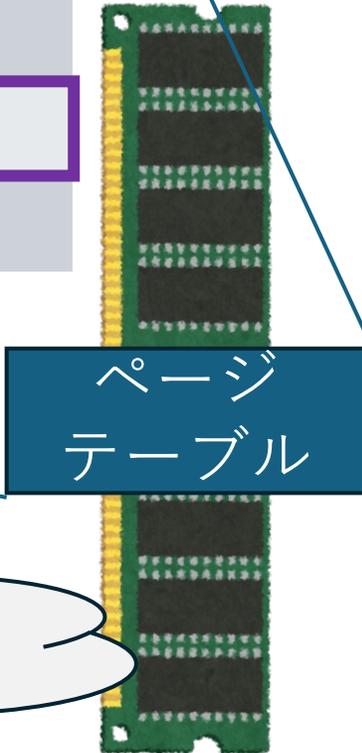
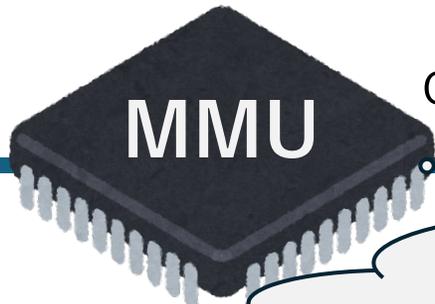
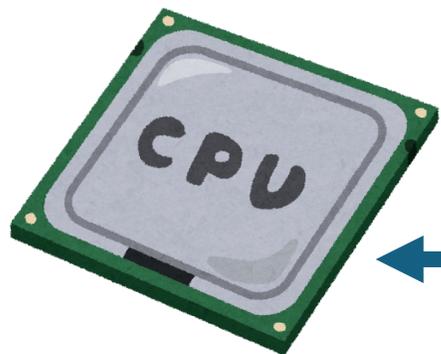
プログラムによるメモリアクセス

非特権モード時の
アクセスの可否も
設定できます

仮想	物理	非特権モード時
0x0000		
0x1000		
...		
0x100000000	0x2000	アクセス可能
...		

仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

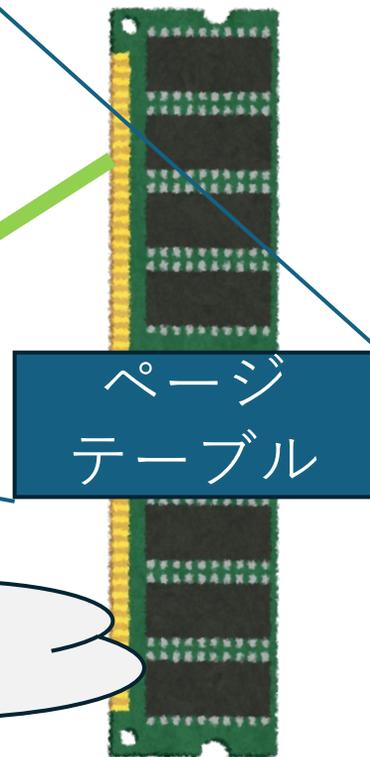
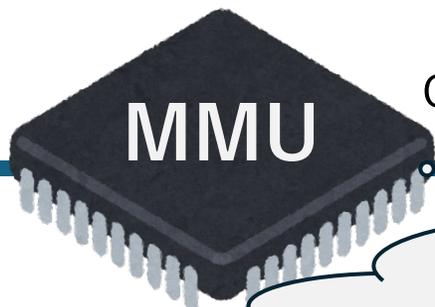
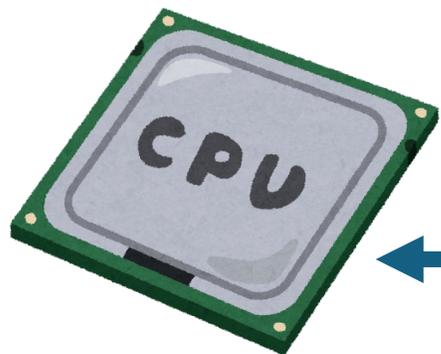
仮想アドレス 0x100000000 は
物理アドレスの 0x2000

プログラムによるメモリアクセス

仮想	物理
0x0000	
0x1000	
...	
0x100000000	0x2000
...	

仮想アドレス 0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000

プログラムによるメモリアクセス

ポイント 1

実行中のプログラムは
MMU が参照している
ページテーブルに記載のある
物理メモリアドレスしか
アクセスできない

```
movl $0x12345678, (%rax)
```

仮想	物理
0x0000	
0x1000	
...	
0x100000000	0x2000
...	

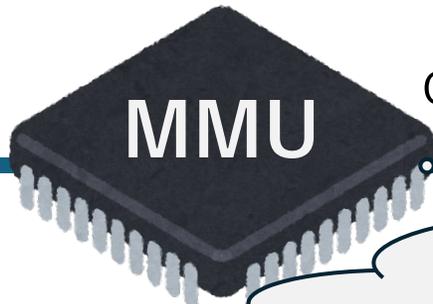
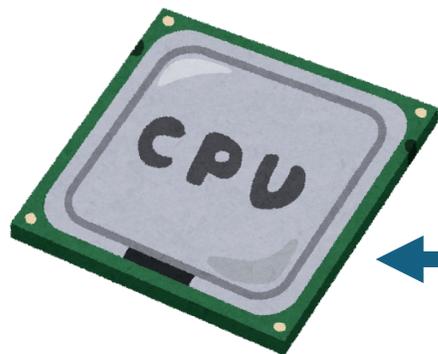
MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

ページ
テーブル

rax:0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000



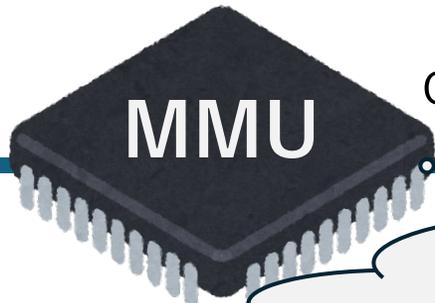
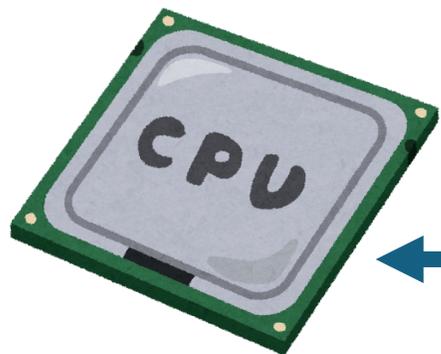
プログラムによるメモリアクセス

ポイント 1

実行中のプログラムは
MMU が参照している
ページテーブルに記載のある
物理メモリアドレスしか
アクセスできない

```
movl $0x12345678, (%rax)
```

仮想	物理
0x0000	例えば物理メモリアドレスが一つも設定されていなければプログラムが仮想メモリアドレス 0 ~ 上限 (256 TB や 128 PB) のどこにアクセスしても Segmentation fault になる
0x1000	
...	
0x100000000	
...	



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

ページ
テーブル

rax: 0x100000000
cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000

プログラムによるメモリアクセス

ポイント 1

実行中のプログラムは MMU が参照しているページテーブルに記載のある物理メモリアドレスしかアクセスできない

仮想	物理
0x0000	
0x1000	
...	
0x100000000	0x2000

ポイント 2

非特権モードでは cr3 レジスタの値を変更する CPU 命令を実行できない (**特権モード**であれば可能)

はアドレス変換に際してされるページテーブルを参照

ページ
テーブル

rax: 0x100000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000



プログラムによるメモリアクセス

ポイント 1

実行中のプログラムは MMU が参照しているページテーブルに記載のある物理メモリアドレスしかアクセスできない

仮想	物理
0x0000	
0x1000	
...	
0x100000000	

特権モードで動作するカーネルはページテーブルの操作を通して
非特権モードで動作するプログラムがアクセス可能なメモリ領域を制限できる

ポイント 2

非特権モードでは cr3 レジスタの値を変更する CPU 命令を実行できない
(**特権モード**であれば可能)

はアドレス変換に際して参照されるページテーブルを参照

ページ
テーブル

rax: 0x1000000000

cr3: ページテーブルの物理アドレス

仮想アドレス 0x1000000000 は物理アドレスの 0x2000

プログラムによるメモリアクセス

ポイント 1

実行中のプログラムは MMU が参照しているページテーブルに記載のある物理メモリアドレスしかアクセスできない

仮想	物理
0x0000	
0x1000	
...	
0x10000000	

特権モードで動作するカーネルはページテーブルの操作を通して**非特権モード**で動作するプログラムがアクセス可能なメモリ領域を制限できる

ポイント 2

非特権モードでは cr3 レジスタの値を変更する CPU 命令を実行できない (**特権モード**であれば可能)

制限に必要なポイント

ページテーブルが置かれた物理メモリアドレスを非特権モード時に参照されるページテーブルに記載してアクセス可能にしないこと

rax: 0x1000000000

cr3: ページテーブルの物理アドレス

仮想アドレス

物理アドレス

はア
される

プログラムによるメモリアクセス

ポイント 1

実行中のプログラムは MMU が参照しているページテーブルに記載のある物理メモリアドレスしかアクセスできない

仮想	物理
0x0000	
0x1000	
...	
0x10000000	

特権モードで動作するカーネルはページテーブルの操作を通して**非特権モード**で動作するプログラムがアクセス可能なメモリ領域を制限できる

ポイント 2

非特権モードでは cr3 レジスタの値を

cr3 は変更できなくてもページテーブルを編集できる場合にはアクセスしたい物理メモリアドレスを書き込めばアクセスできてしまうので

制限に必要なポイント

ページテーブルが置かれた物理メモリアドレスを非特権モード時に参照されるページテーブルに記載してアクセス可能にしないこと

mov

rax:
cr3:

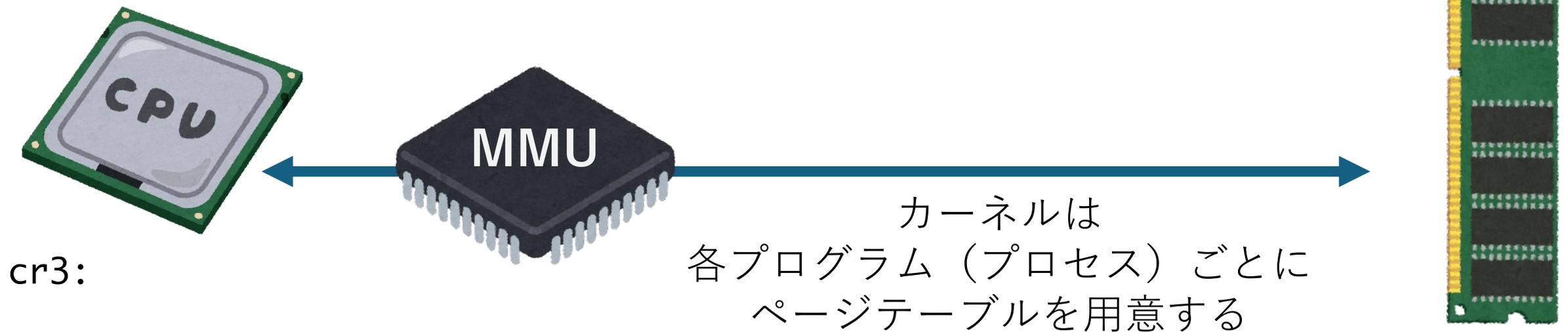
mmap システムコールでカーネルに
仮想アドレス 0x1000000000 へ対応する物理
メモリの確保と設定をリクエストした場合

```
#include <stdio.h>
#include <stdlib.h>
#include <sys/mman.h>
int main(void)
{
    void *mem = mmap((void *) 0x100000000,
                     0x1000, PROT_READ|PROT_WRITE,
                     MAP_PRIVATE|MAP_ANONYMOUS|MAP_FIXED,
                     -1, 0);

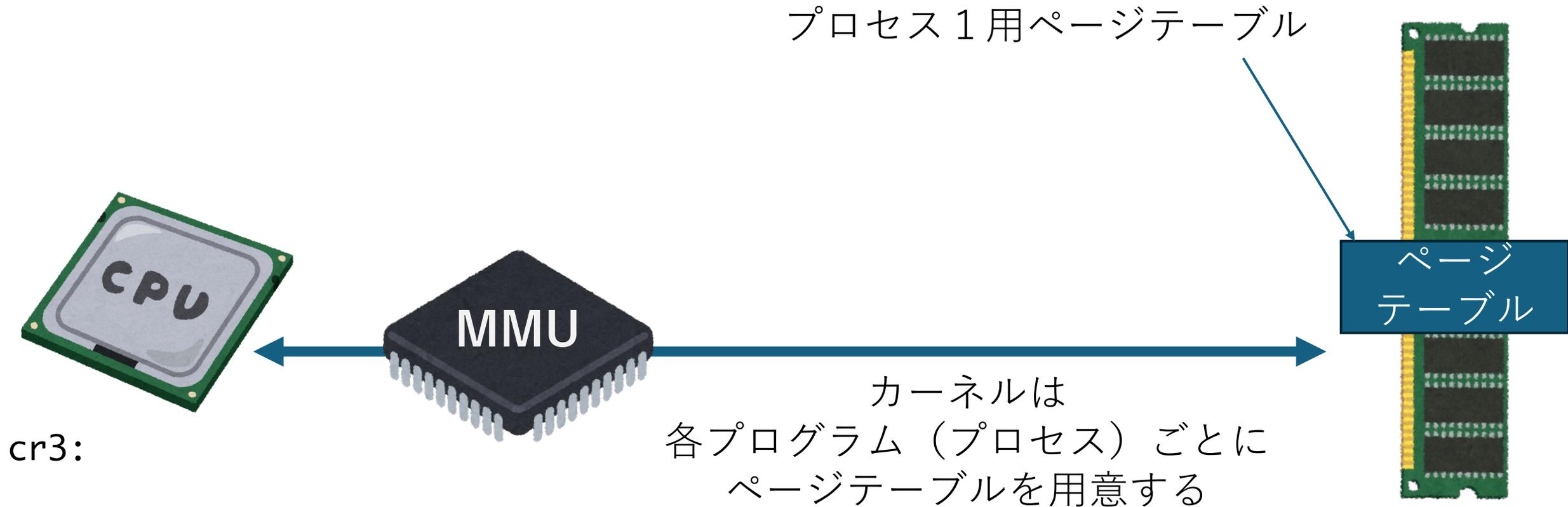
    if (mem != (void *) 0x100000000 || mem == MAP_FAILED) {
        printf("mmap failed %p¥n", mem);
        exit(1);
    }
    *((int *) 0x100000000) = 0x12345678;
    printf("success¥n");
}
```

セグメンテーション違反が発生せず
success と表示されれば成功

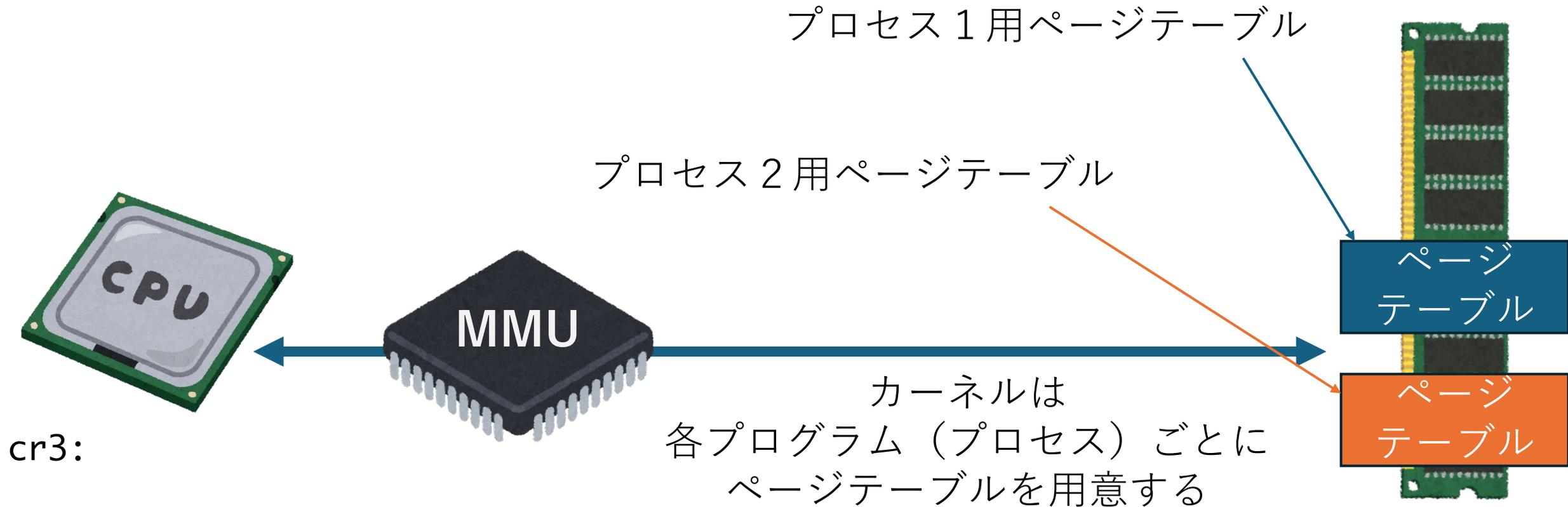
一般的な OS での運用



一般的な OS での運用

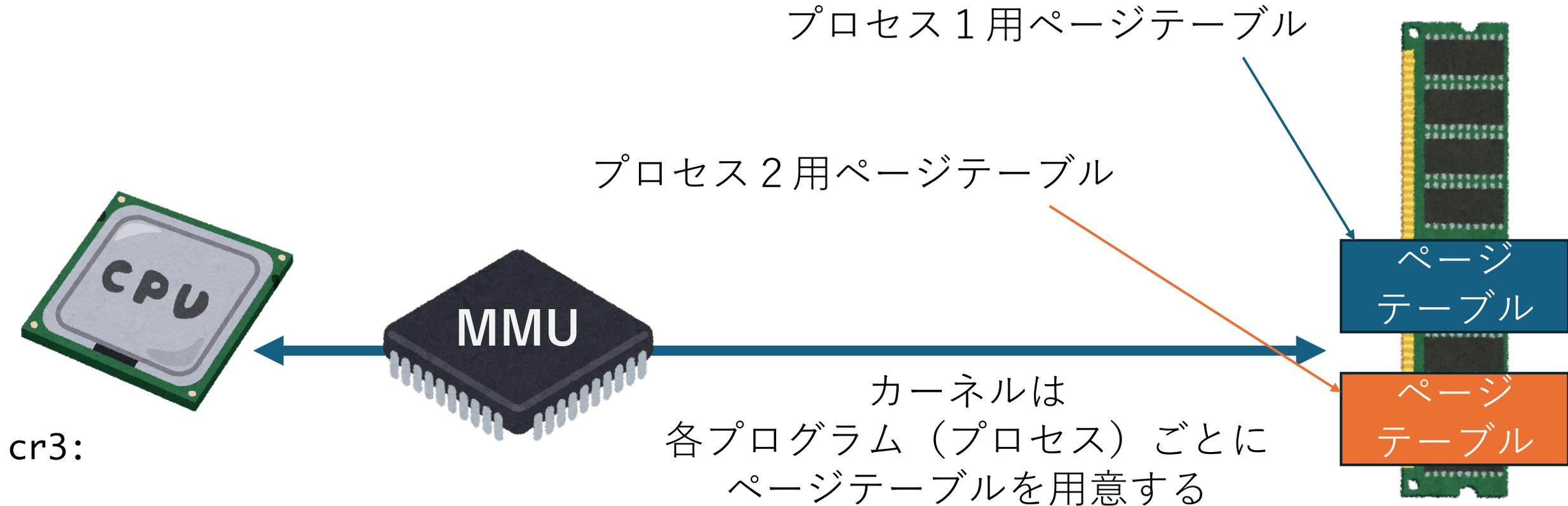


一般的な OS での運用



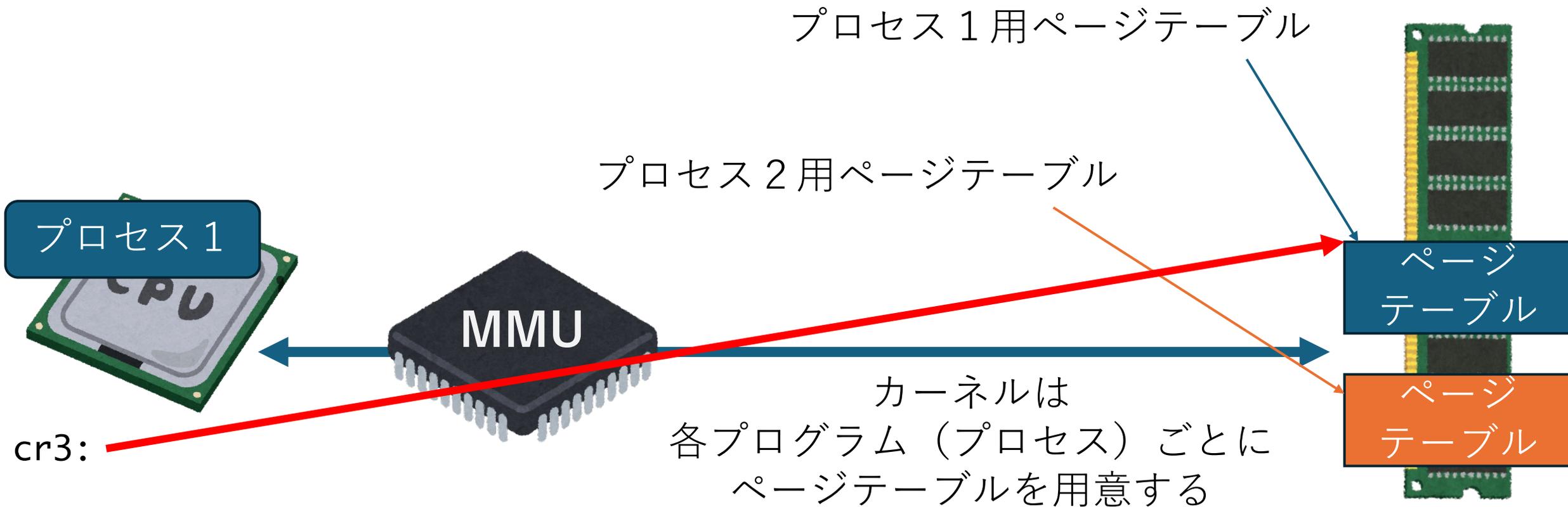
一般的な OS での運用

カーネルは実行しているプロセスの切り替え毎に
cr3 の値を書き換えてMMU が参照するページテーブルを切り替える



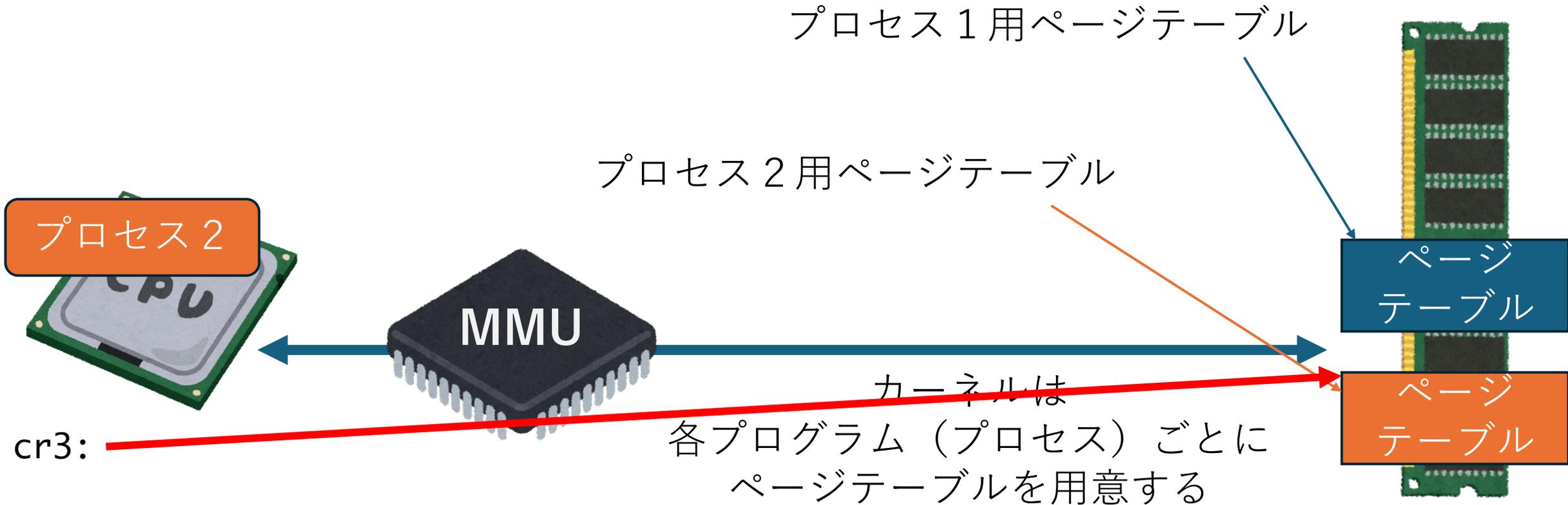
一般的な OS での運用

カーネルは実行しているプロセスの切り替え毎に
cr3 の値を書き換えてMMU が参照するページテーブルを切り替える



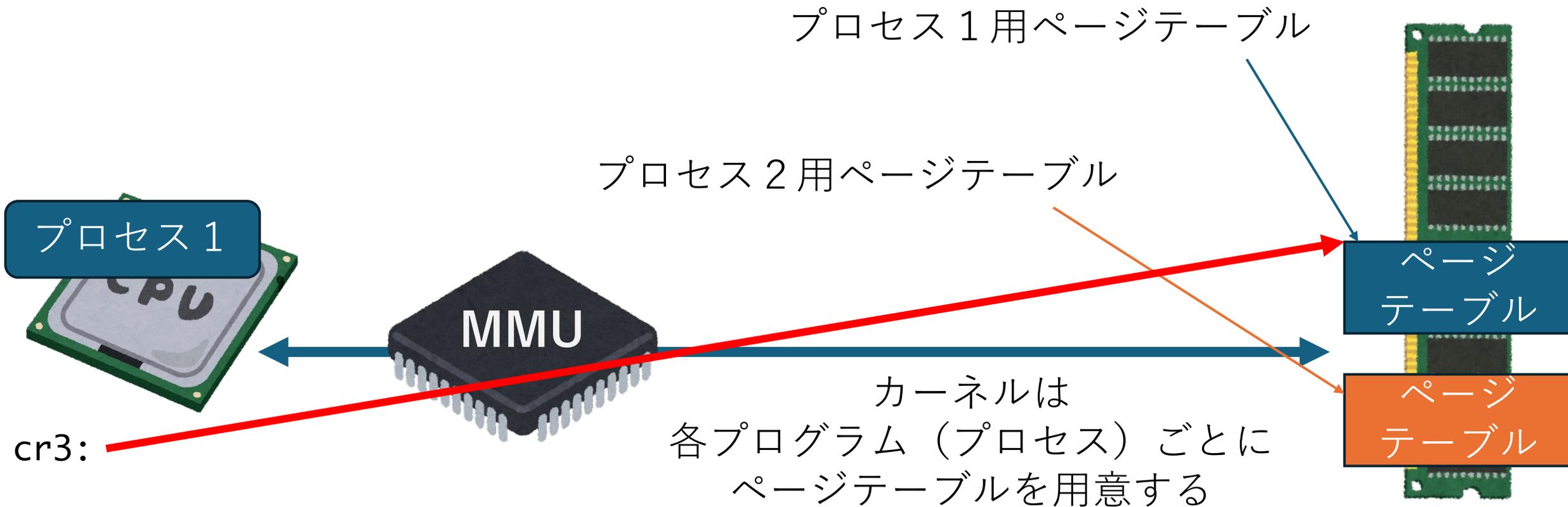
一般的な OS での運用

カーネルは実行しているプロセスの切り替え毎に
cr3 の値を書き換えてMMU が参照するページテーブルを切り替える



一般的な OS での運用

カーネルは実行しているプロセスの切り替え毎に
cr3 の値を書き換えてMMU が参照するページテーブルを切り替える



一般的な OS での運用

仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 1

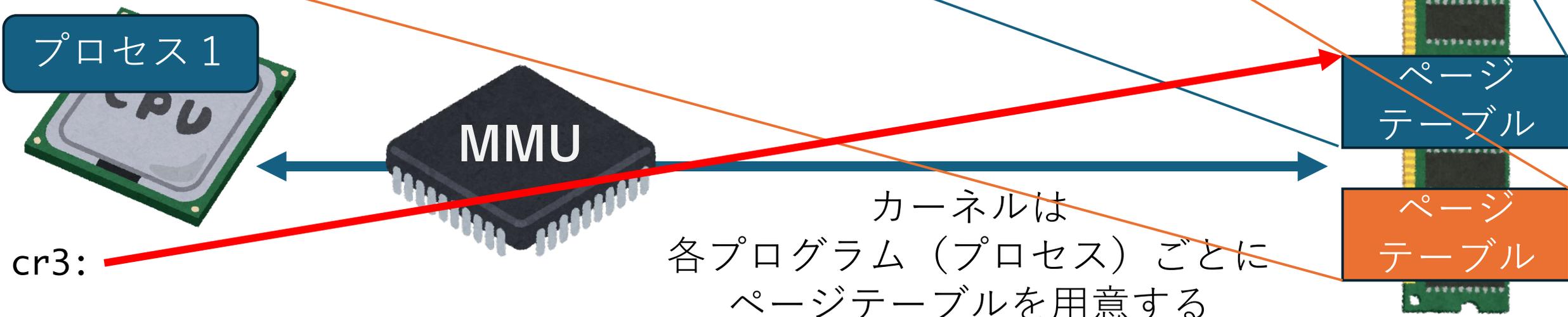
MMU

ページ
テーブル

ページ
テーブル

cr3:

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する



一般的な OS での運用

仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 1

MMU

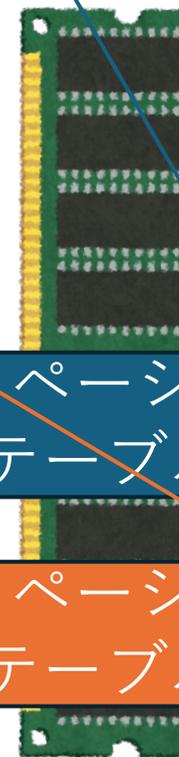
ページ
テーブル

ページ
テーブル

cr3:

プロセス 1 が実行中

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する



一般的な OS での運用

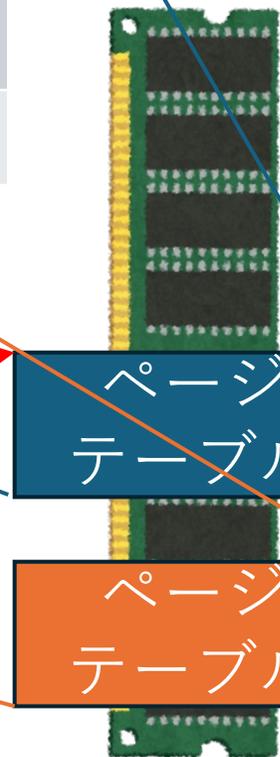
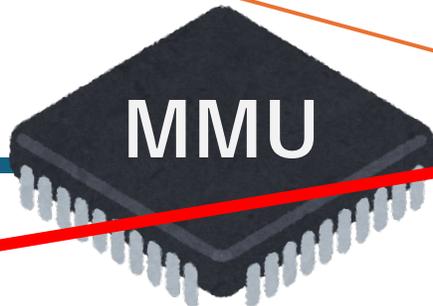
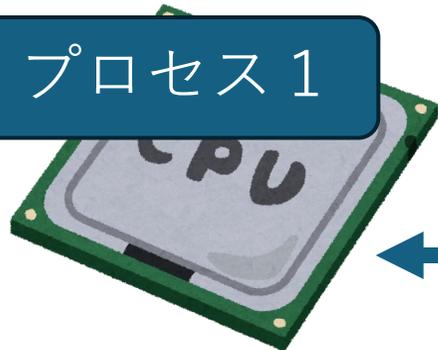
仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 1



cr3: **仮想** アドレス 0x1000 へアクセス

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

一般的な OS での運用

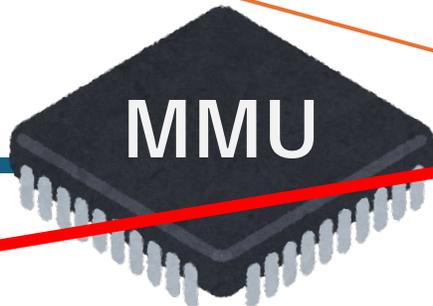
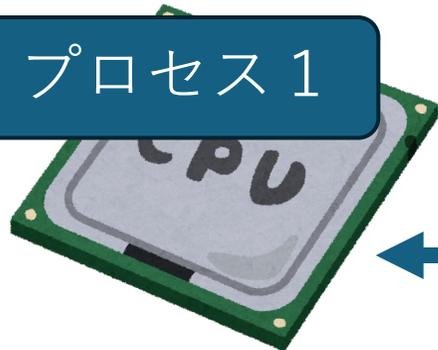
仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 1



ページ
テーブル

ページ
テーブル

cr3:

仮想 アドレス 0x1000 へアクセス

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する



一般的な OS での運用

物理アドレス 0x2000 へアクセス

仮想	物理
0x0000	
0x1000	0x3000
...	
...	

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 2 用ページテーブル

プロセス 1 用ページテーブル

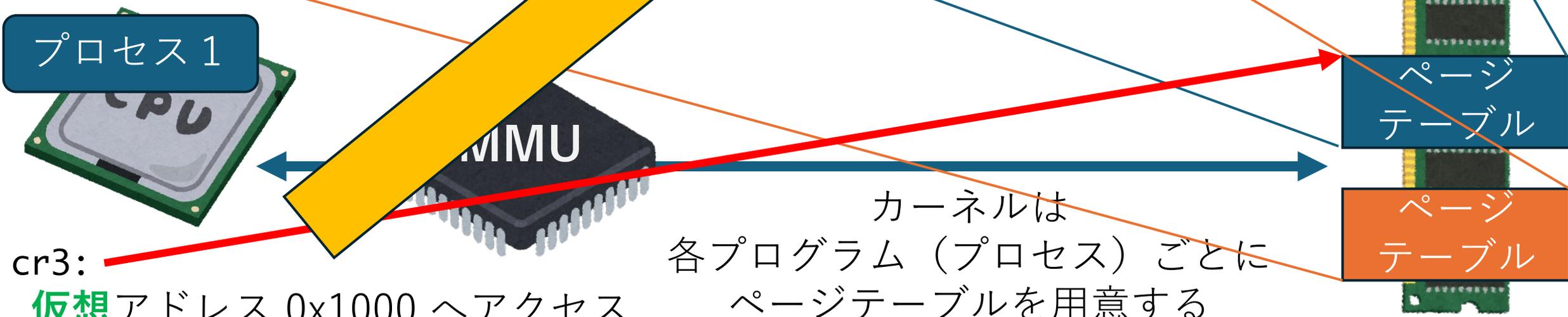
プロセス 1

ページ
テーブル

ページ
テーブル

cr3: 仮想アドレス 0x1000 へアクセス

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する



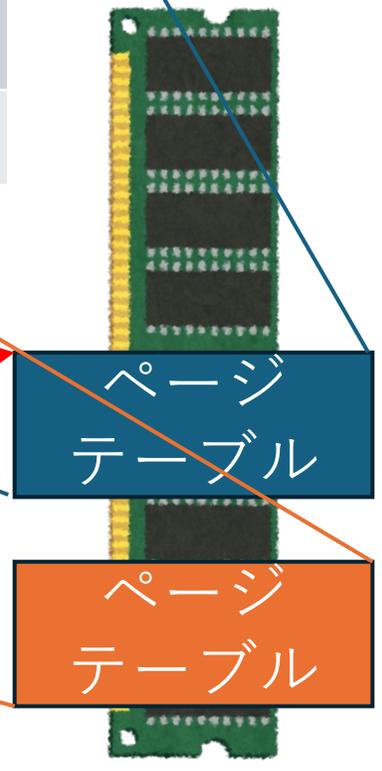
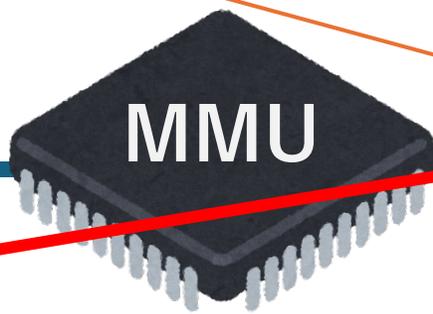
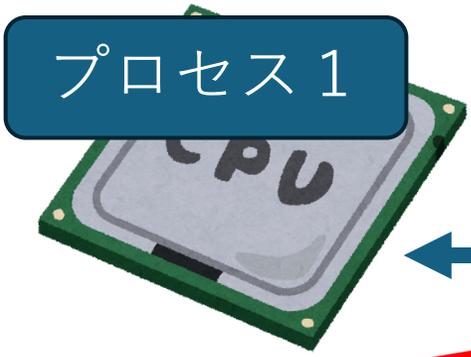
一般的な OS での運用

仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル



cr3:

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

一般的な OS での運用

仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 2

MMU

ページ
テーブル

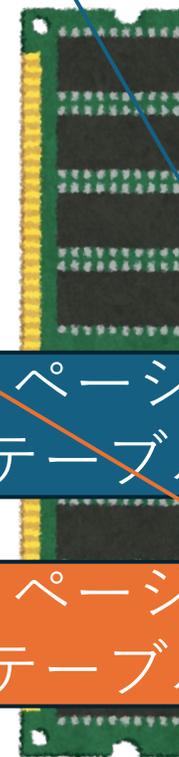
ページ
テーブル

cr3:

プロセスが切り替えられた

カーネルは

各プログラム（プロセス）ごとに
ページテーブルを用意する



一般的な OS での運用

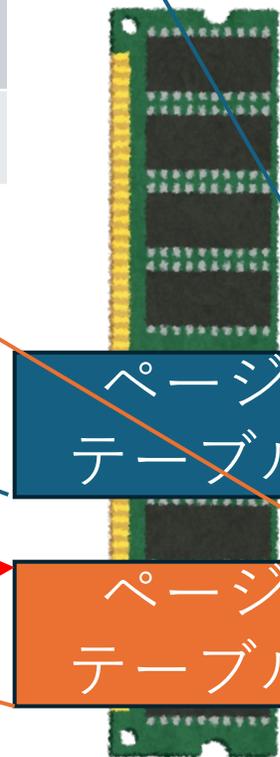
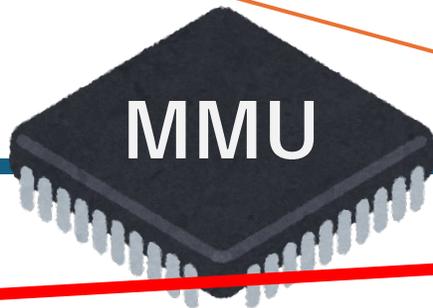
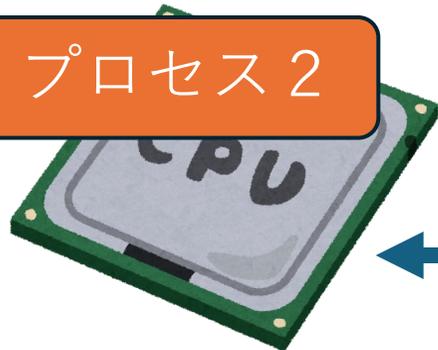
仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 2



ページ
テーブル

ページ
テーブル

cr3:

仮想 アドレス 0x1000 へアクセス

カーネルは

各プログラム（プロセス）ごとに
ページテーブルを用意する

一般的な OS での運用

仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 2

MMU

ページ
テーブル

ページ
テーブル

cr3:

仮想 アドレス 0x1000 へアクセス

カーネルは

各プログラム（プロセス）ごとに
ページテーブルを用意する



一般的な OS での運用

物理アドレス 0x3000 へアクセス

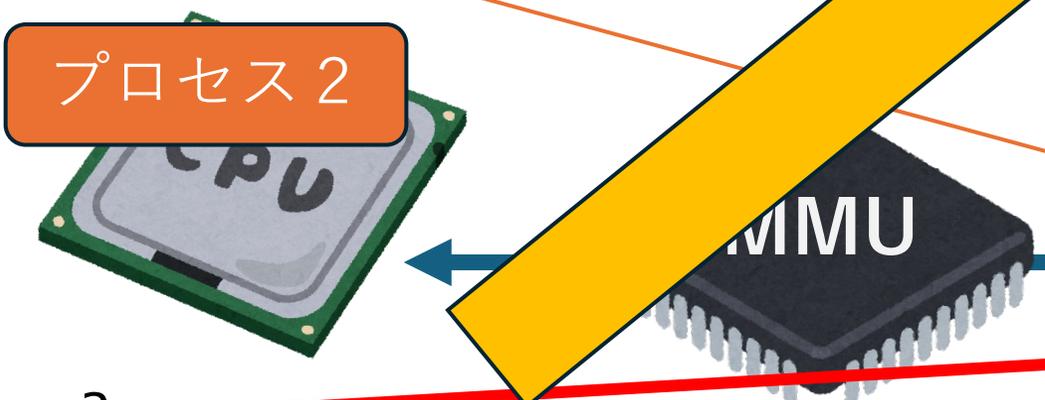
仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

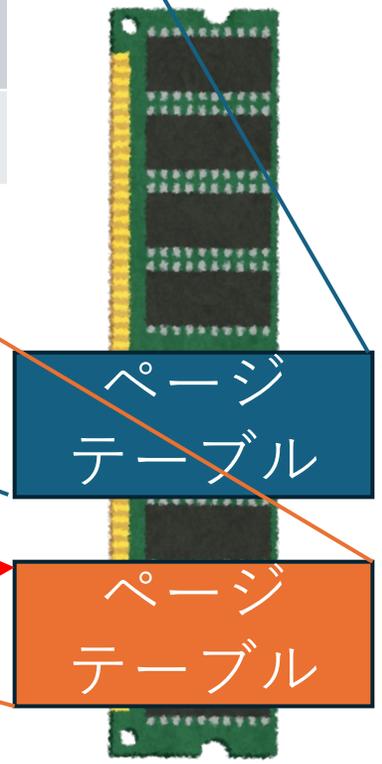
仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 2



cr3: **仮想** アドレス 0x1000 へアクセス



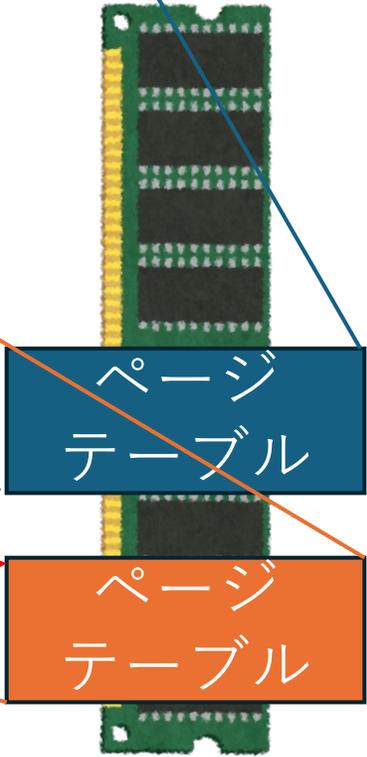
カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

一般的な OS での運用

仮想	物理	仮想	物理
0x0000		0x0000	
0x1000	0x3000	0x1000	0x2000
...		...	
...		...	

同じ**仮想**アドレス 0x1000 へのアクセスも
ページテーブルが違えば
別の**物理**アドレスへのアクセスになります

プロセス 2



cr3: **仮想**アドレス 0x1000 へアクセス

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

ページ
テーブル

ページ
テーブル

一般的な OS での運用

仮想	物理	仮想	物理
0x0000		0x0000	
0x1000	0x3000	0x1000	0x2000
...			
...			

同じ**仮想**アドレス 0x1000 へのアクセスも
ページテーブルが違えば
別の**物理**アドレスへのアクセスになります

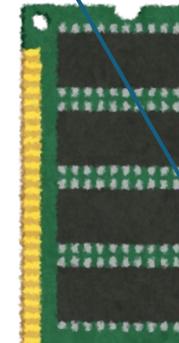
プロセス 2

ポイント

カーネルは異なるプロセスが同じ物理メモリ領域へアクセスできないように
注意しながらページテーブルを用意します
これによりプロセス間の分離 (isolation) が担保されます

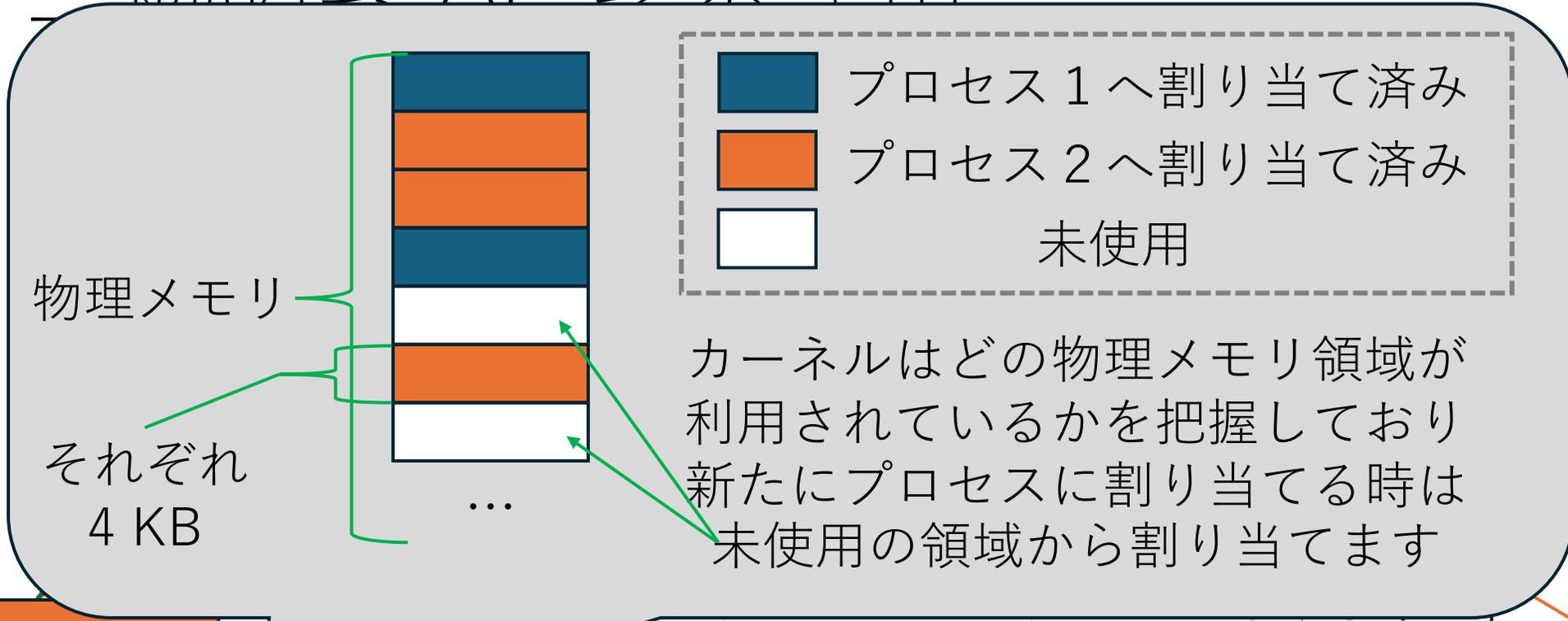
cr

仮想



ページ
ブル

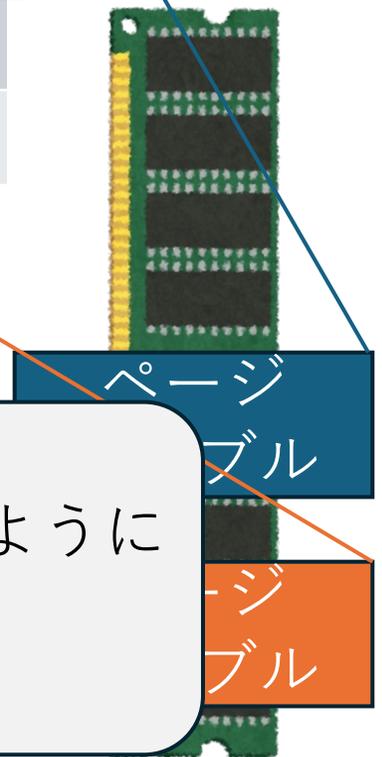
ページ
ブル



プロセス 2

ポイント

カーネルは異なるプロセスが同じ物理メモリ領域へアクセスできないように
注意しながらページテーブルを用意します
これによりプロセス間の分離 (isolation) が担保されます



一般的な OS での運用

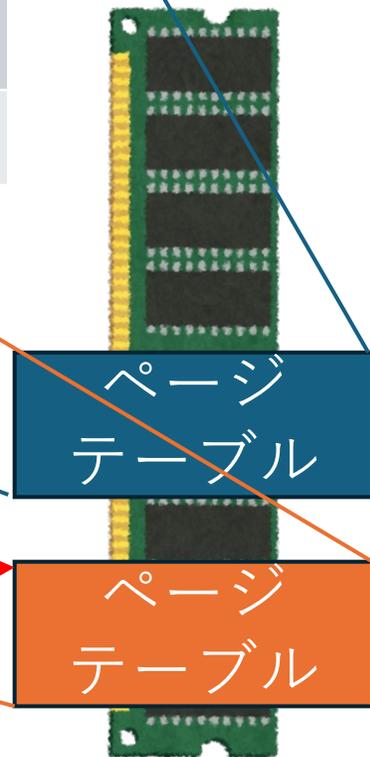
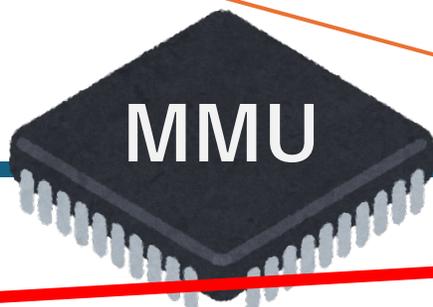
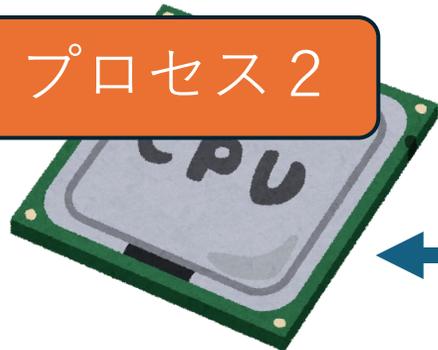
仮想	物理
0x0000	
0x1000	0x3000
...	
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
...	
...	

プロセス 1 用ページテーブル

プロセス 2



cr3: **仮想** アドレス 0x1000 へアクセス

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

ページ
テーブル

ページ
テーブル

一般的な OS での運用

ちなみに、敢えて異なるプロセスが
同じ物理アドレスを参照できるように
設定するのがプロセス間の**共有メモリ**です

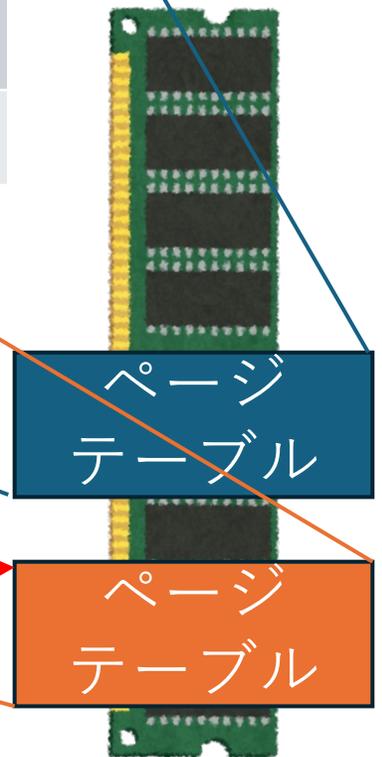
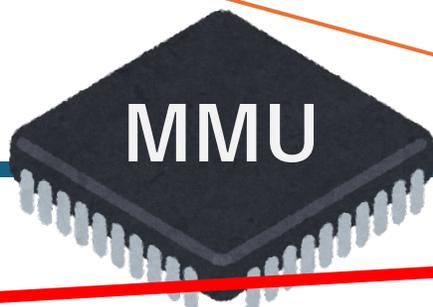
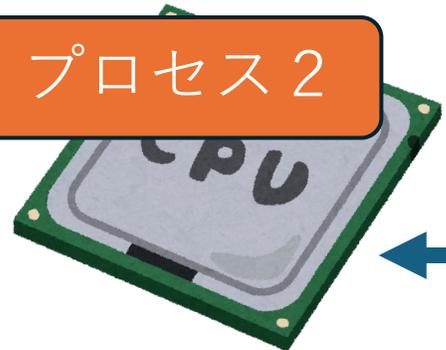
仮想	物理
0x0000	0x4000
0x1000	0x3000
...	仮想アドレスはプロセス1とプロセス2で 一致する必要はありません
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
0x2000	0x4000
...	

プロセス 1 用ページテーブル

プロセス 2



この場合
物理メモリアドレス 0x4000 ~ 0x4fff が
プロセス 1 とプロセス 2 で共有されます

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

cr3: **仮想** アドレス 0x1000 へアクセス

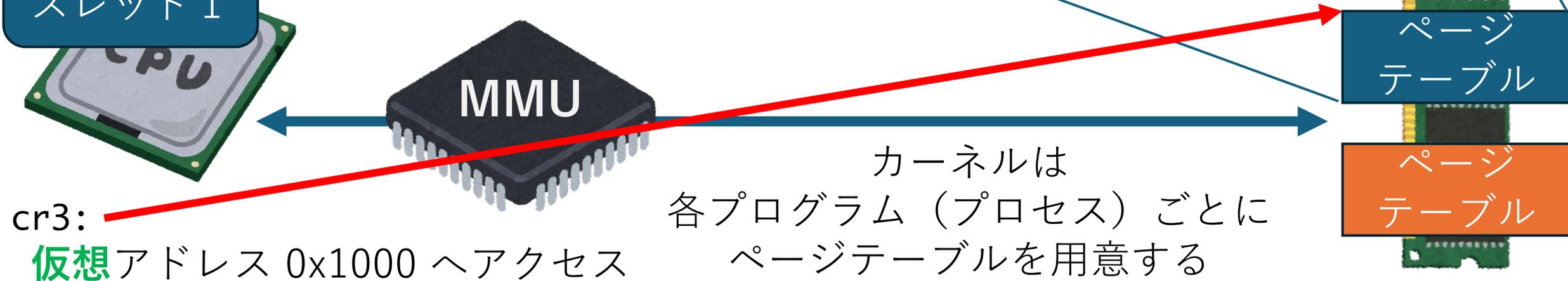
一般的な OS での運用

ちなみに、同じプロセスから生成されたスレッドは
実行時に MMU が同じページテーブルを参照するため
スレッド間でメモリ空間が共有されます

仮想	物理
0x0000	
0x1000	0x2000
0x2000	0x4000
...	

プロセス 1 用ページテーブル

プロセス 1
スレッド 1



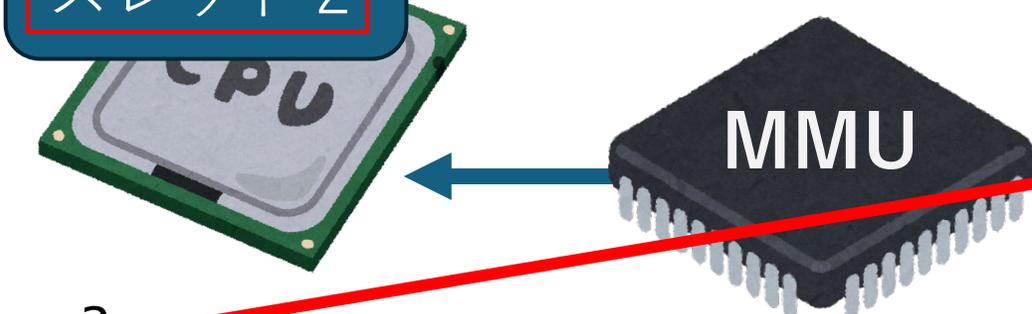
一般的な OS での運用

ちなみに、同じプロセスから生成されたスレッドは
実行時に MMU が同じページテーブルを参照するため
スレッド間でメモリ空間が共有されます

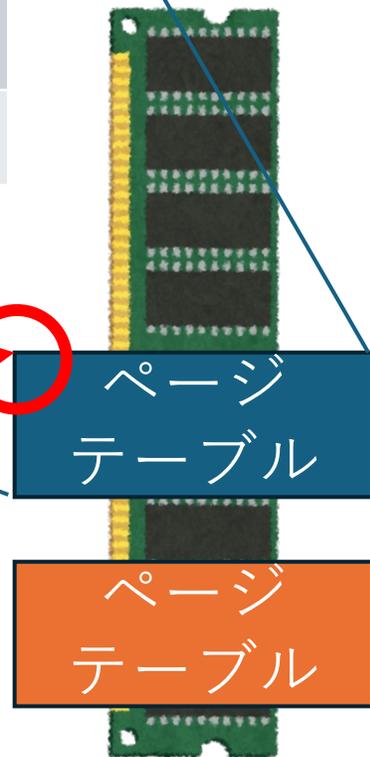
仮想	物理
0x0000	
0x1000	0x2000
0x2000	0x4000
...	

プロセス 1 用ページテーブル

プロセス 1
スレッド 2



同じプロセス 1 から生成された
スレッド 1 とスレッド 2 の
切り替え時には cr3 は書き変えない



cr3: **仮想** アドレス 0x1000 へアクセス

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのため

Q. デバイスへの

アクセスとは

具体的に何？

本能的には

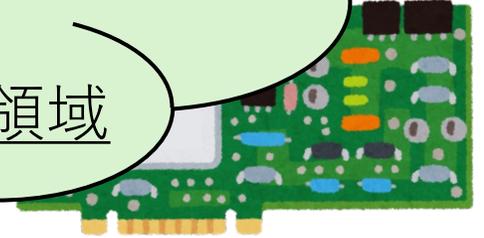
メモリの読み書き

Q. ユーザー空間とは何？

A. CPU が非特権モードで

動作している間に

アクセス可能なメモリ領域



デバイス

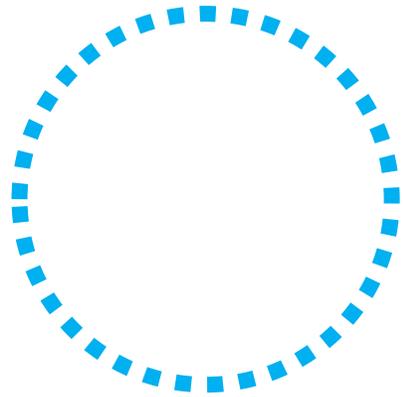
アクセスの可否？

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する



カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセス

Q. デバイスへの

アクセスとは

具体的に何？

本能的には

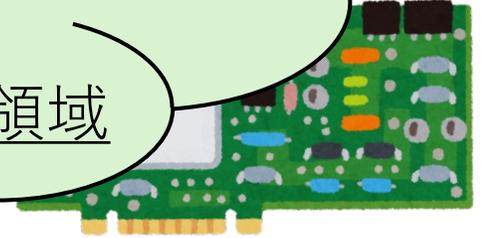
の読み書き

Q. ユーザー空間とは何？

A. CPU が非特権モードで

動作している間に

アクセス可能なメモリ領域



デバイス

アクセスの可否？

Q.

A.

一般的にカーネルがページテーブルを通じて制限を適用している

る？
常では

Q. このプログラムはバイパスする？

A. ユーザー空間プログラムへデバイス操作用メモリ領域へのアクセスを許可する

アクセスが許可されていないから

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセス

Q. デバイスへのアクセスとは具体的に何？

基本的にはメモリの読み書き

カーネルがページテーブルを通じてアクセスを許可した物理メモリ領域および対応する仮想アドレス

Q. ユーザー空間とは何？

A. CPU が非特権モードで動作している間にアクセス可能なメモリ領域



デバイス

アクセスの可否？

Q. 一般的にカーネルがページテーブルを通じて制限を適用している

る？
常では
への

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作用メモリ領域へのアクセスを許可する

アクセスが許可されていないから

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

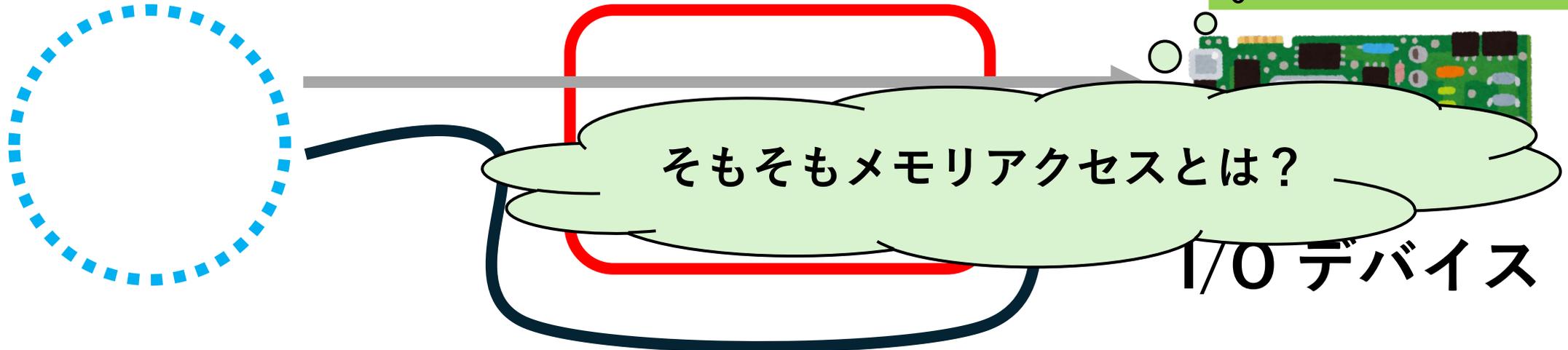
A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

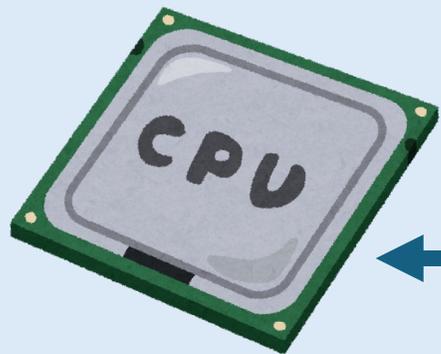
Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

プログラムによるメモリアクセス

0x100000000 ~
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```



この命令の実行時に CPU から
メモリアクセスが試みられる

rax:0x100000000



プログラムによるメモリアクセス

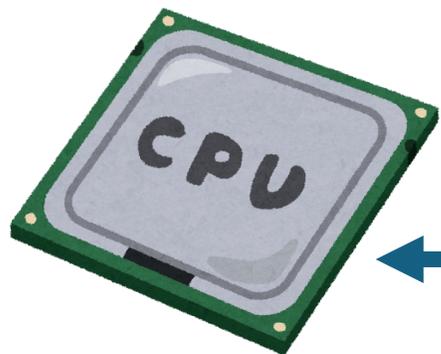
ポイント

このような CPU 命令による
メモリの読み書きで

デバイスを操作することができます

0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

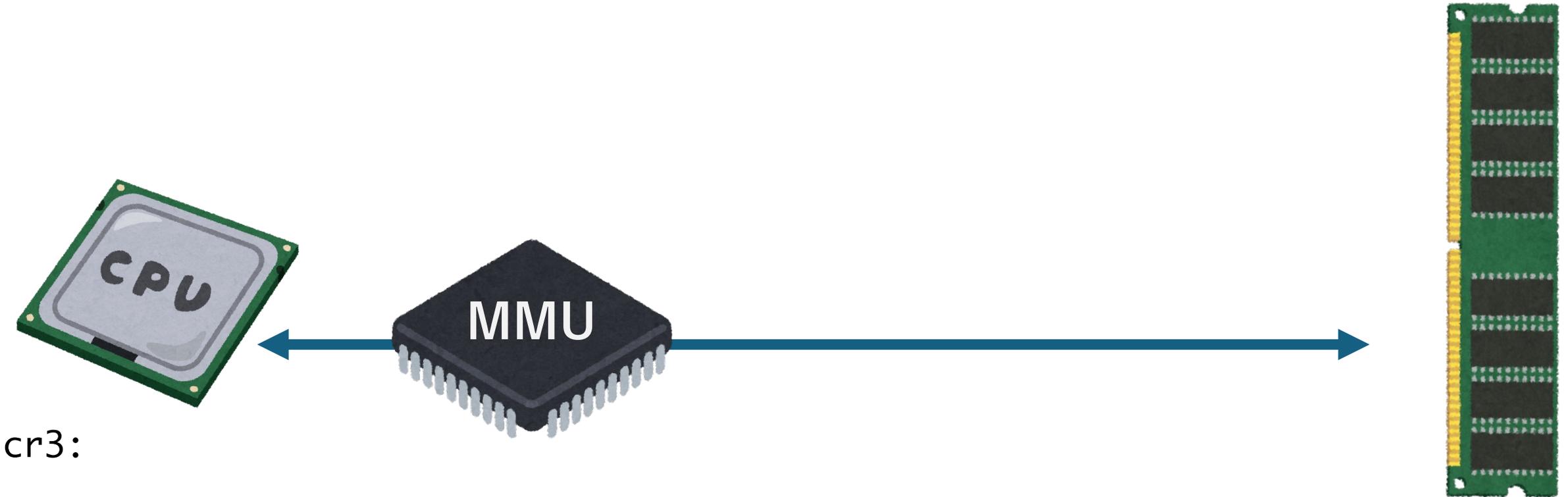


この命令の実行時に CPU から
メモリアクセスが試みられる

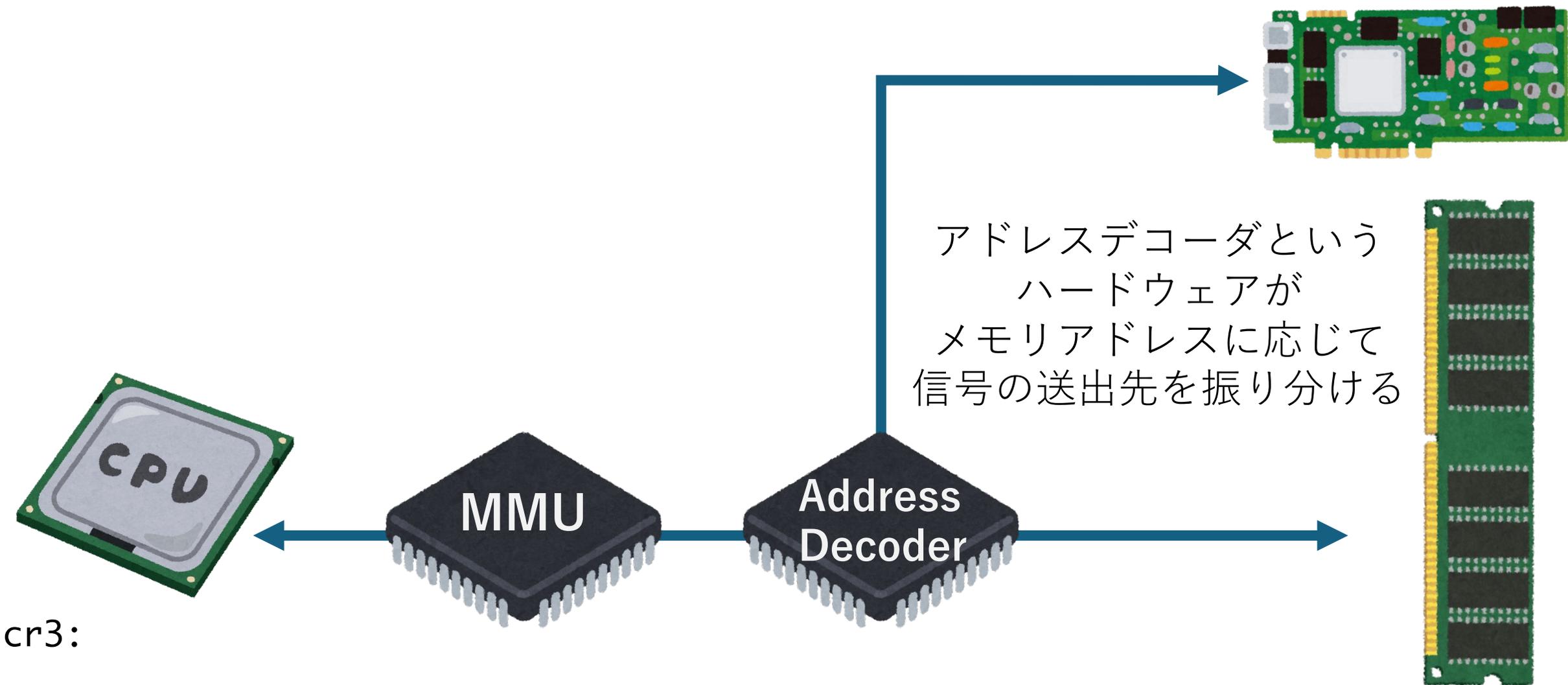
```
rax:0x100000000
```



メモリの読み書きを通じたデバイス操作

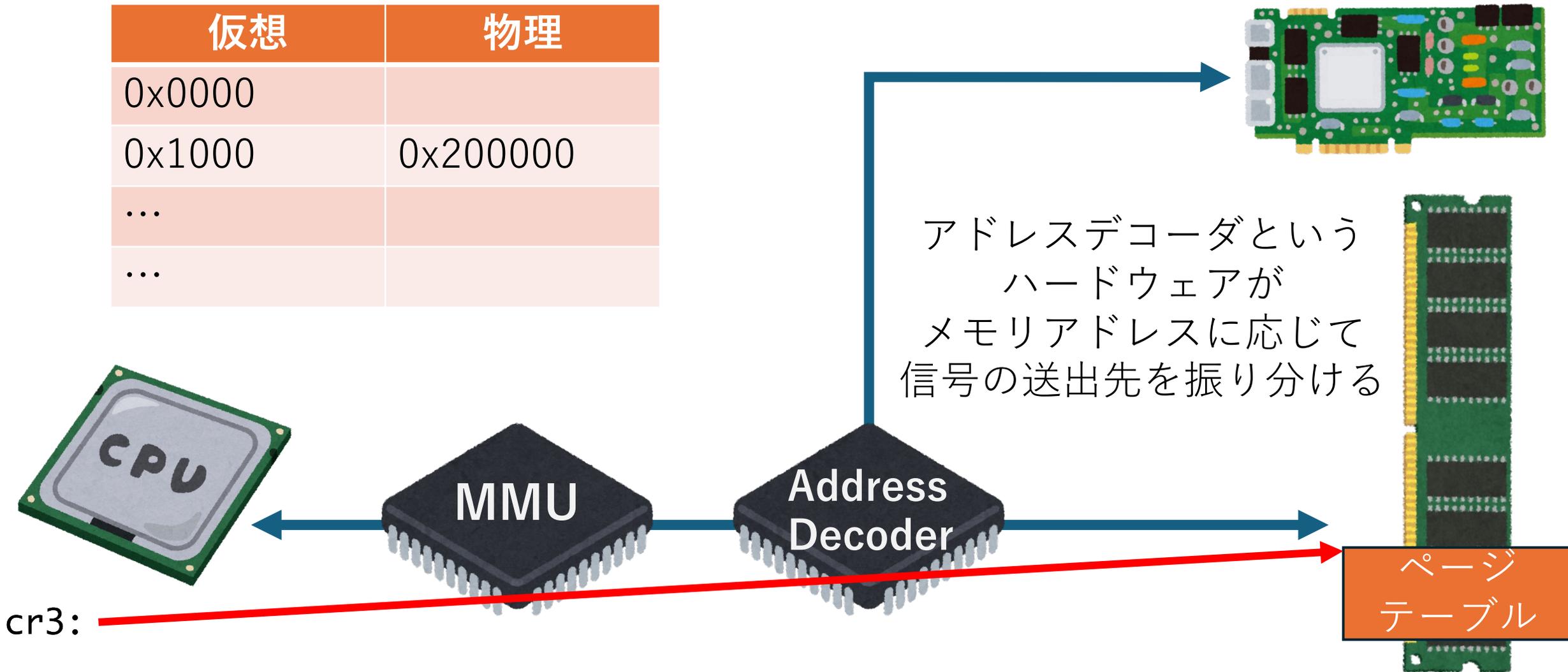


メモリの読み書きを通じたデバイス操作



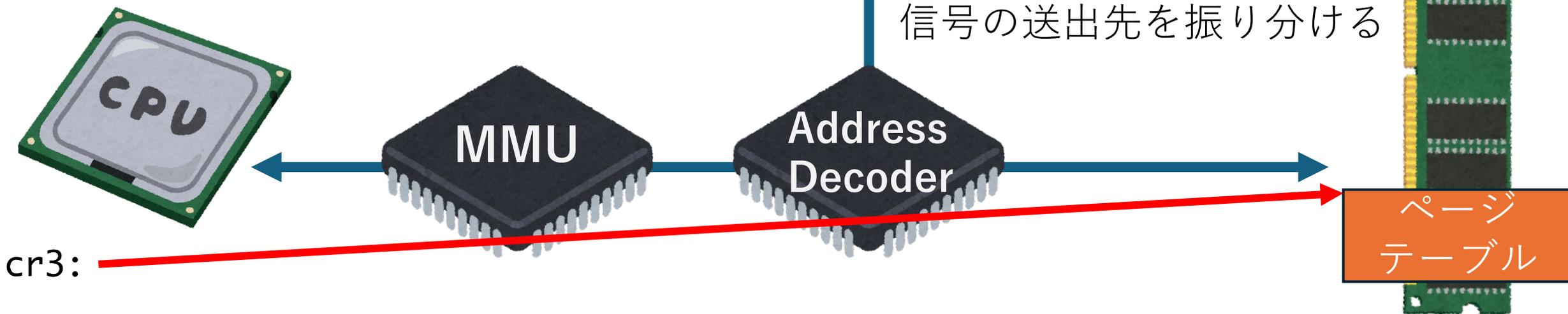
メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
...	
...	



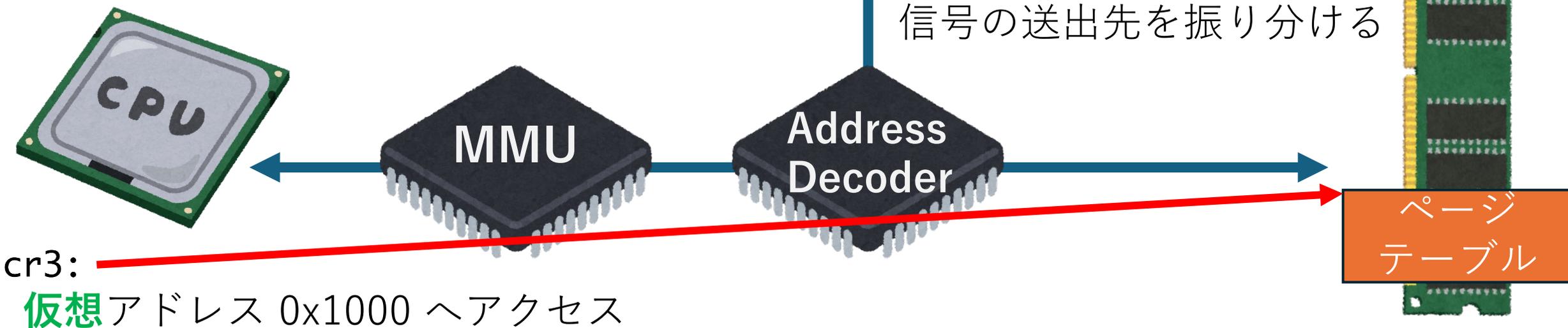
メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



メモリの読み書きを通じたデバイス操作

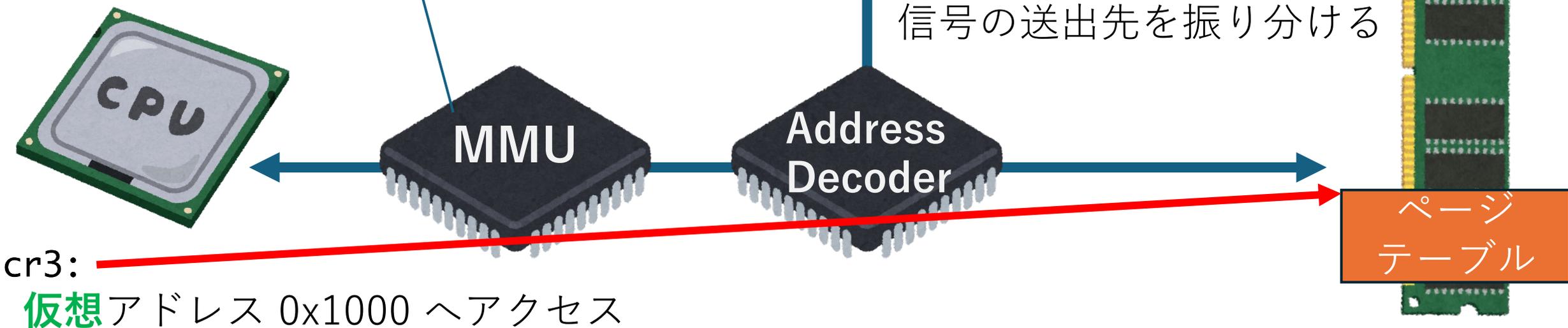
仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



cr3: **仮想** アドレス 0x1000 へアクセス

メモリの読み書きを通じたデバイス操作

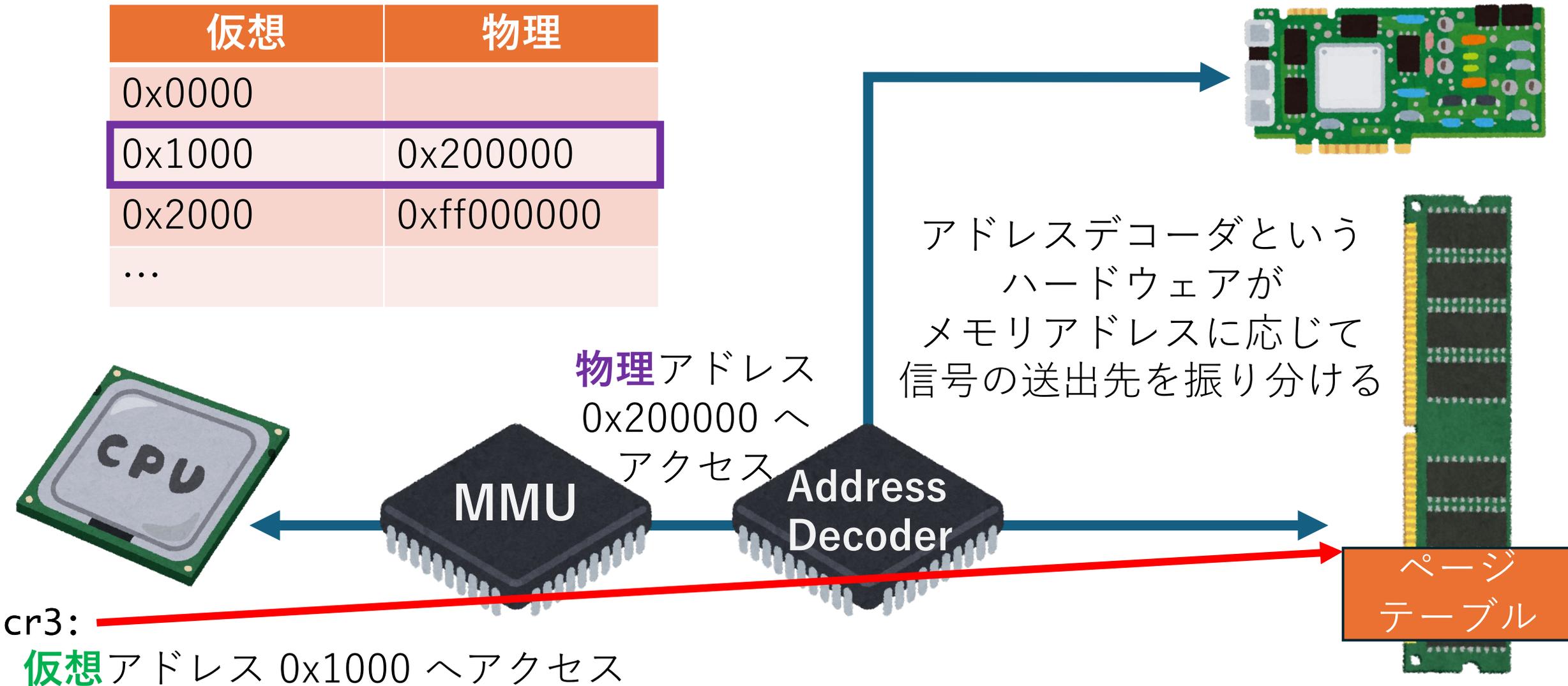
仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



cr3: **仮想** アドレス 0x1000 へアクセス

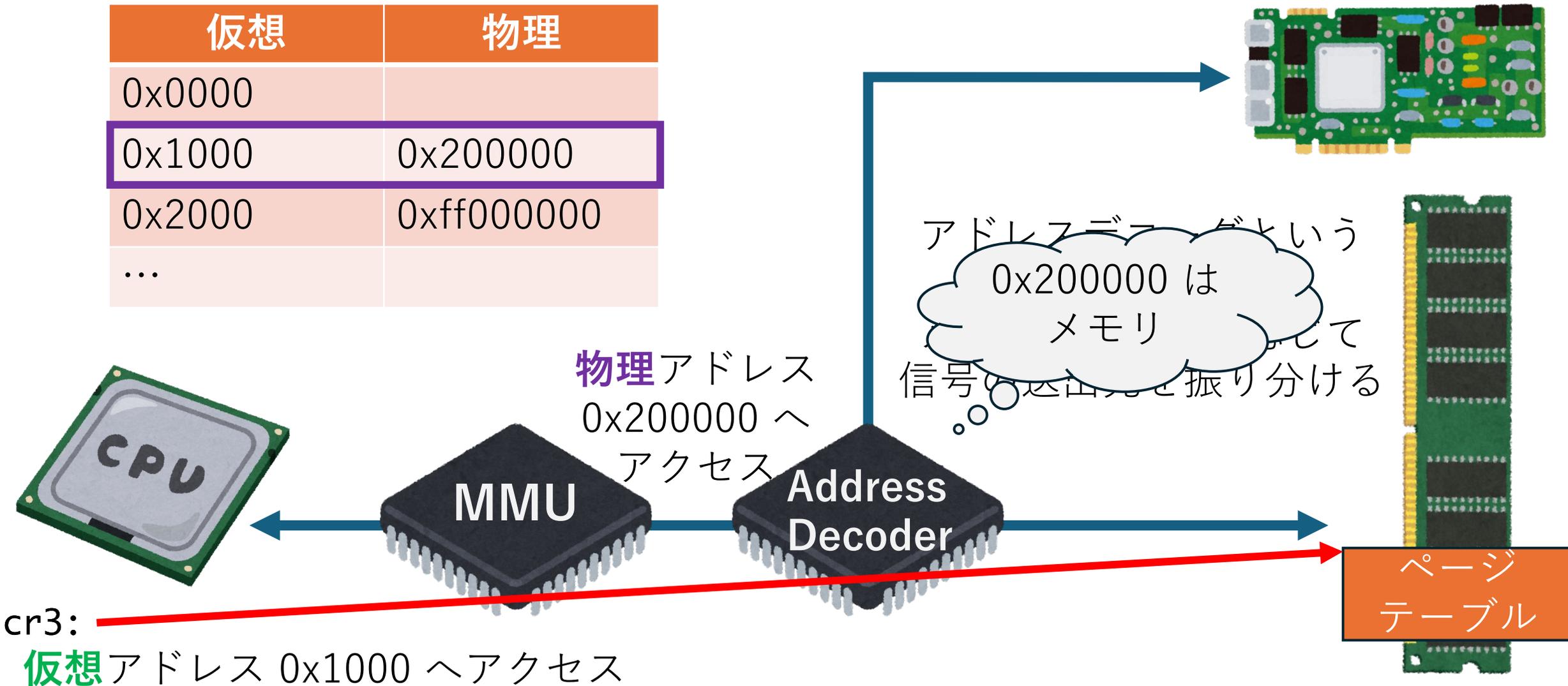
メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



メモリの読み書きを通じたデバイス操作

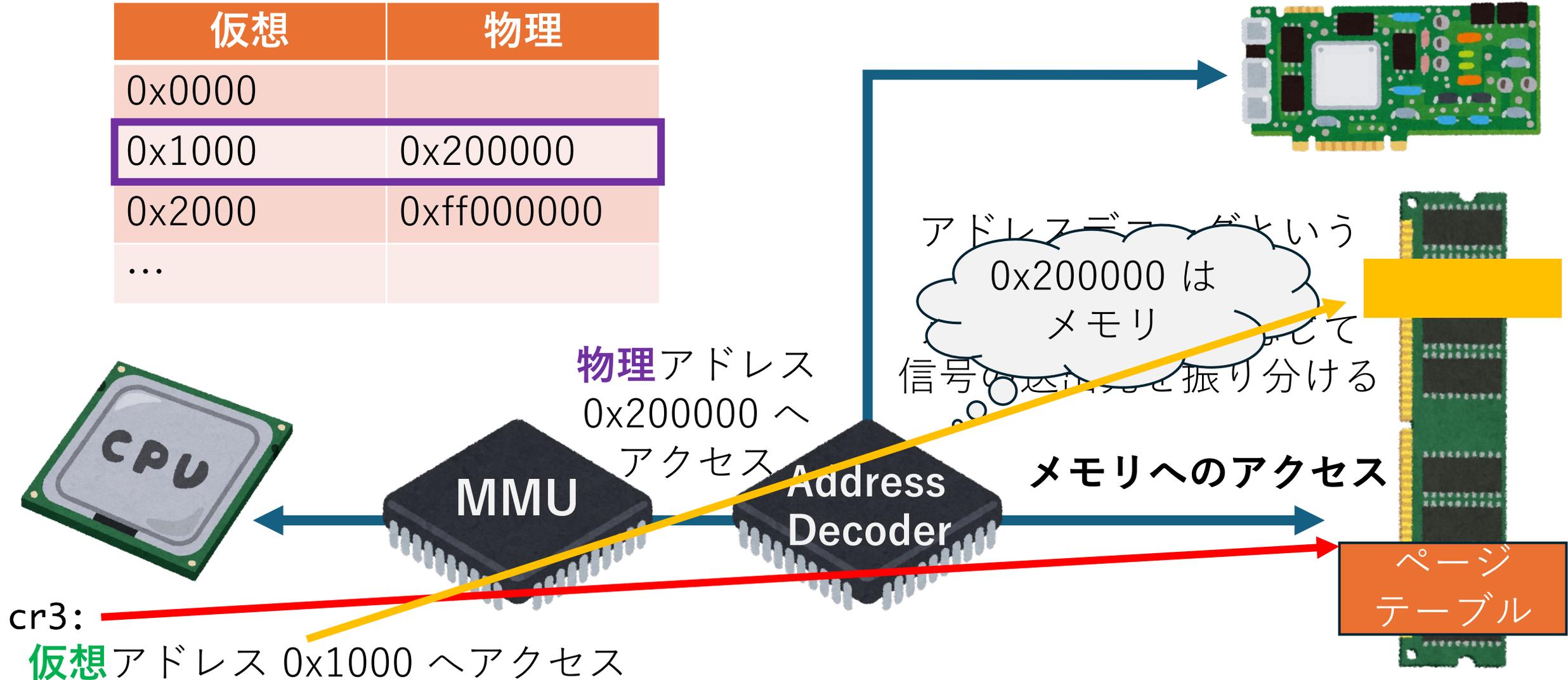
仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



cr3: 仮想アドレス 0x1000 へアクセス

メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



アドレスミューザという
0x200000 は
メモリ
信号の送付先を振り分ける

物理アドレス
0x200000 へ
アクセス

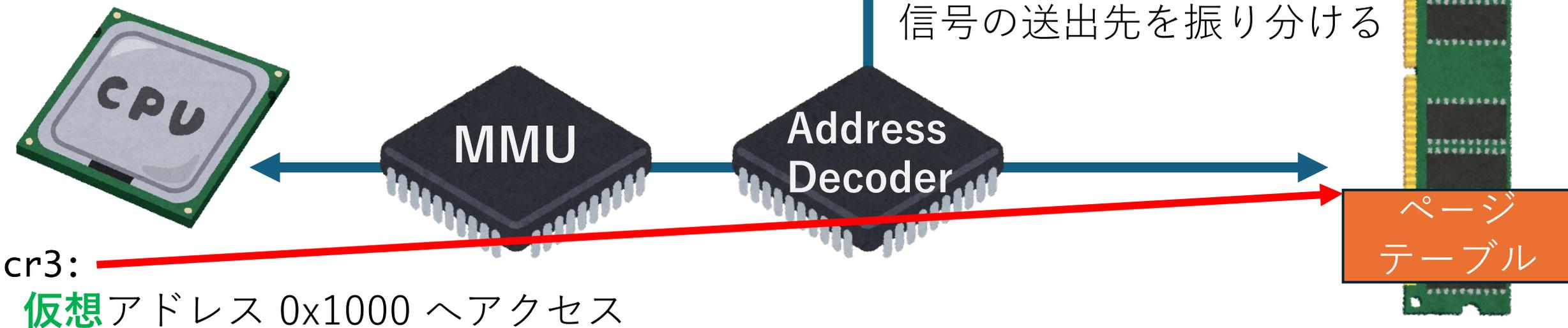
メモリへのアクセス

cr3: 仮想アドレス 0x1000 へアクセス

ページ
テーブル

メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	

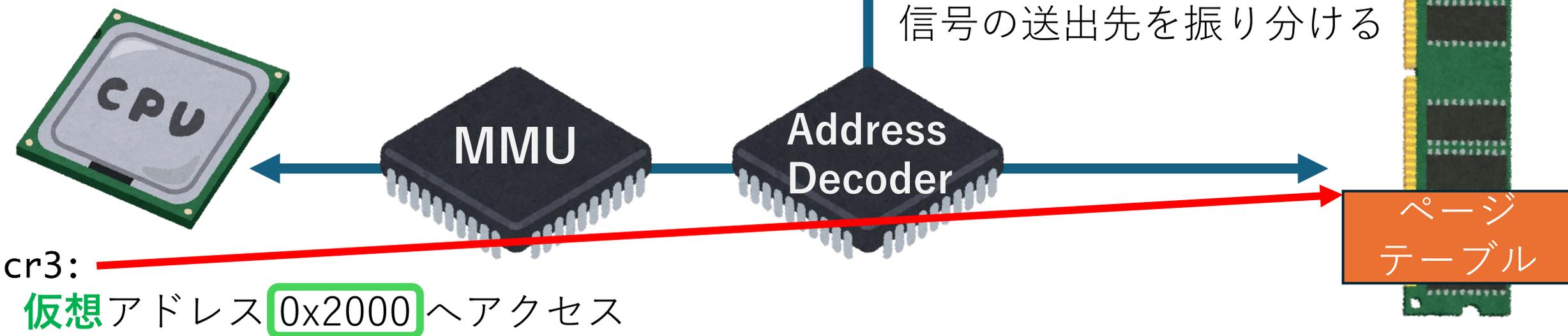


cr3:

仮想アドレス 0x1000 へアクセス

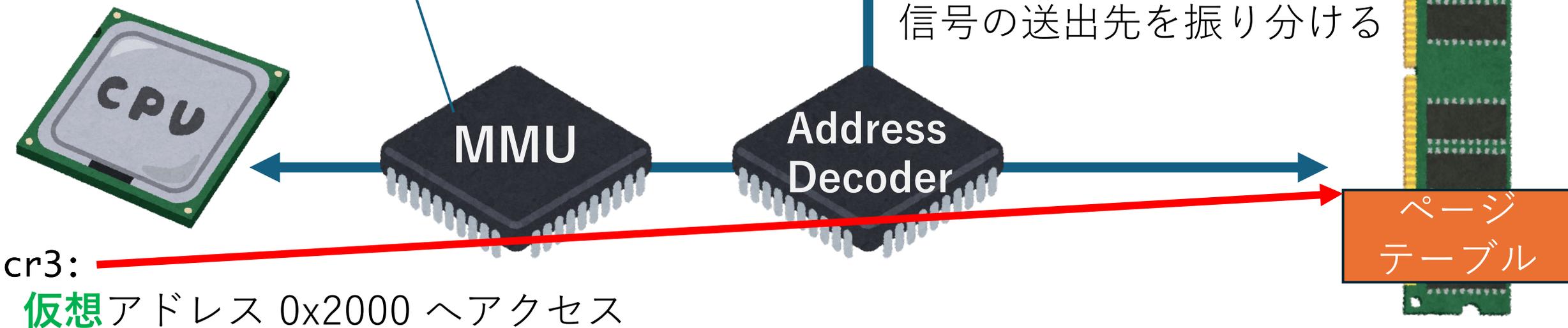
メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	



メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	

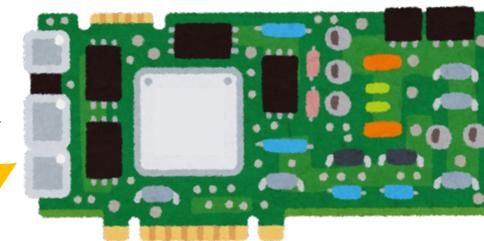


cr3: **仮想** アドレス 0x2000 へアクセス

メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	

NIC へのアクセス

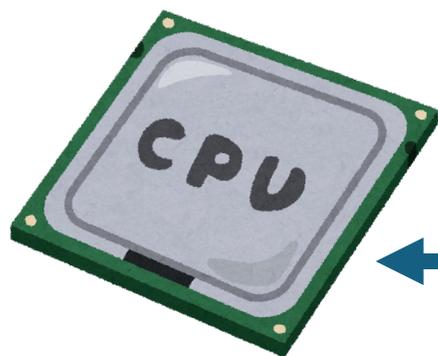


アドレスごまかすという
0xff000000 は
NIC
信号の出入りを振り分ける

物理アドレス
0xff000000 へ
アクセス

Address
Decoder

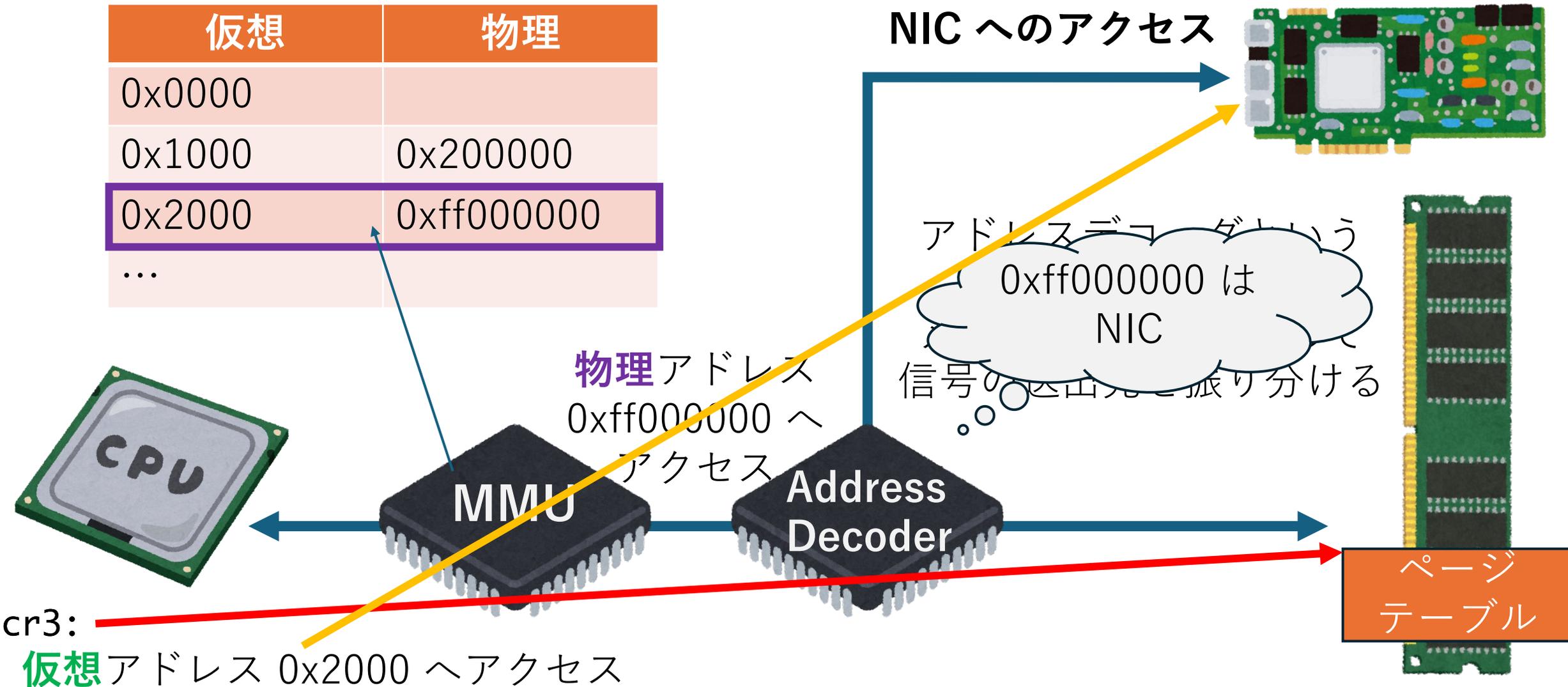
MMU



ページ
テーブル

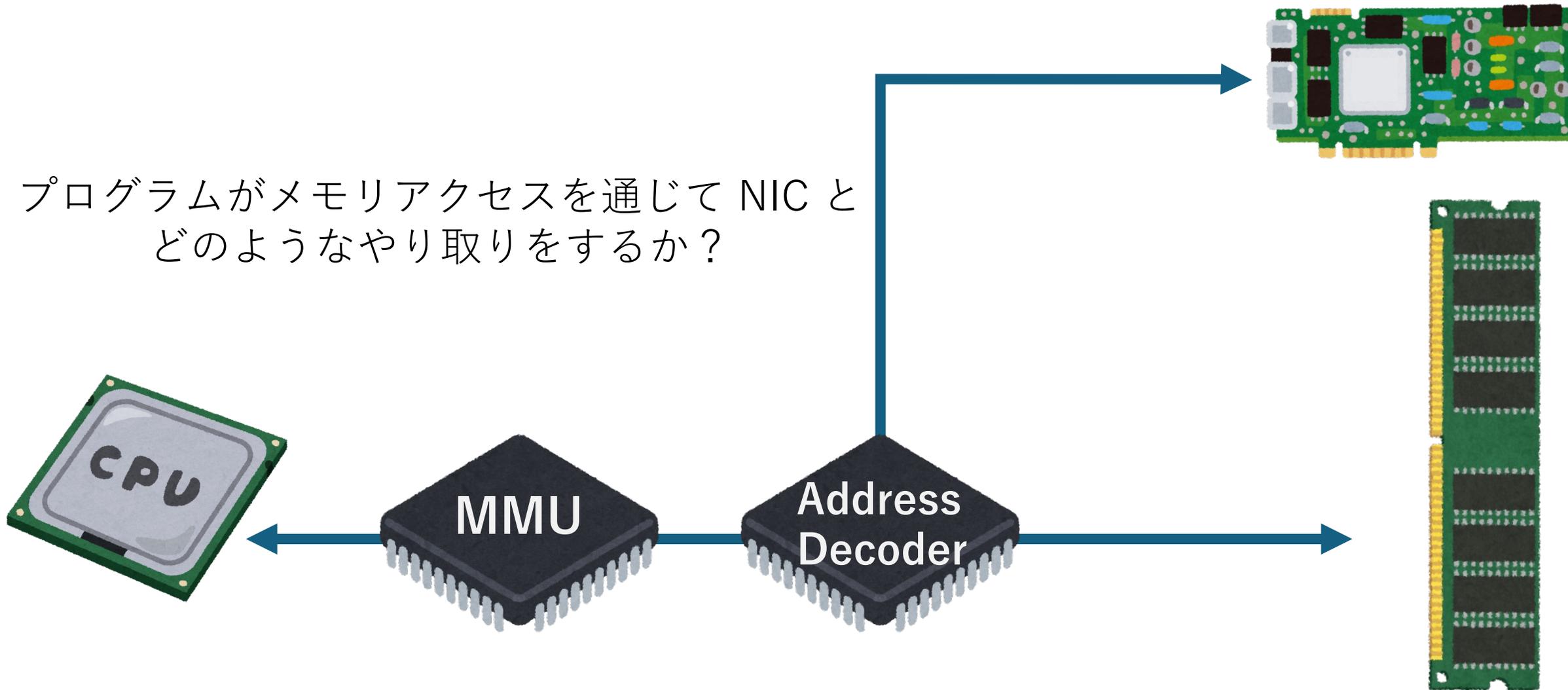
cr3:

仮想アドレス 0x2000 へアクセス



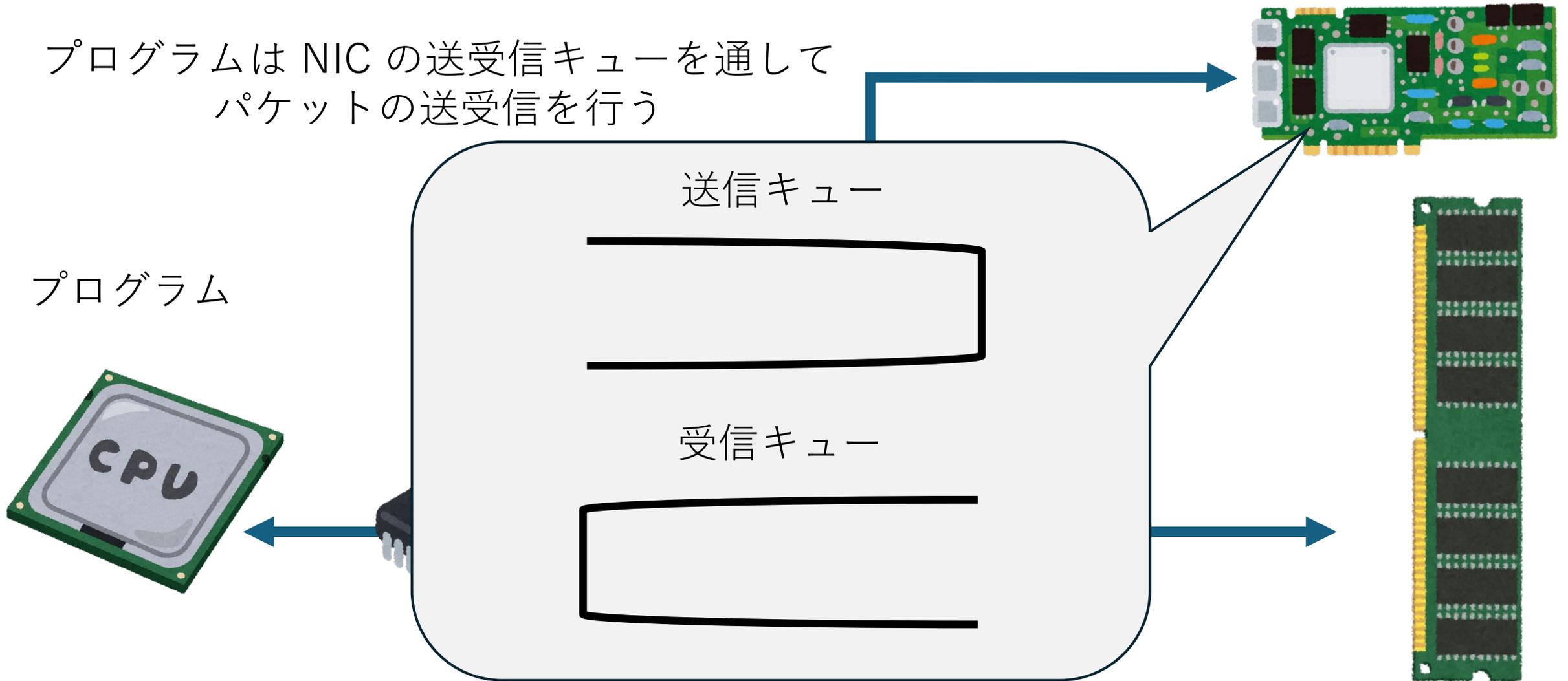
メモリの読み書きを通じたデバイス操作

プログラムがメモリアクセスを通じて NIC と
どのようなやり取りをするか？



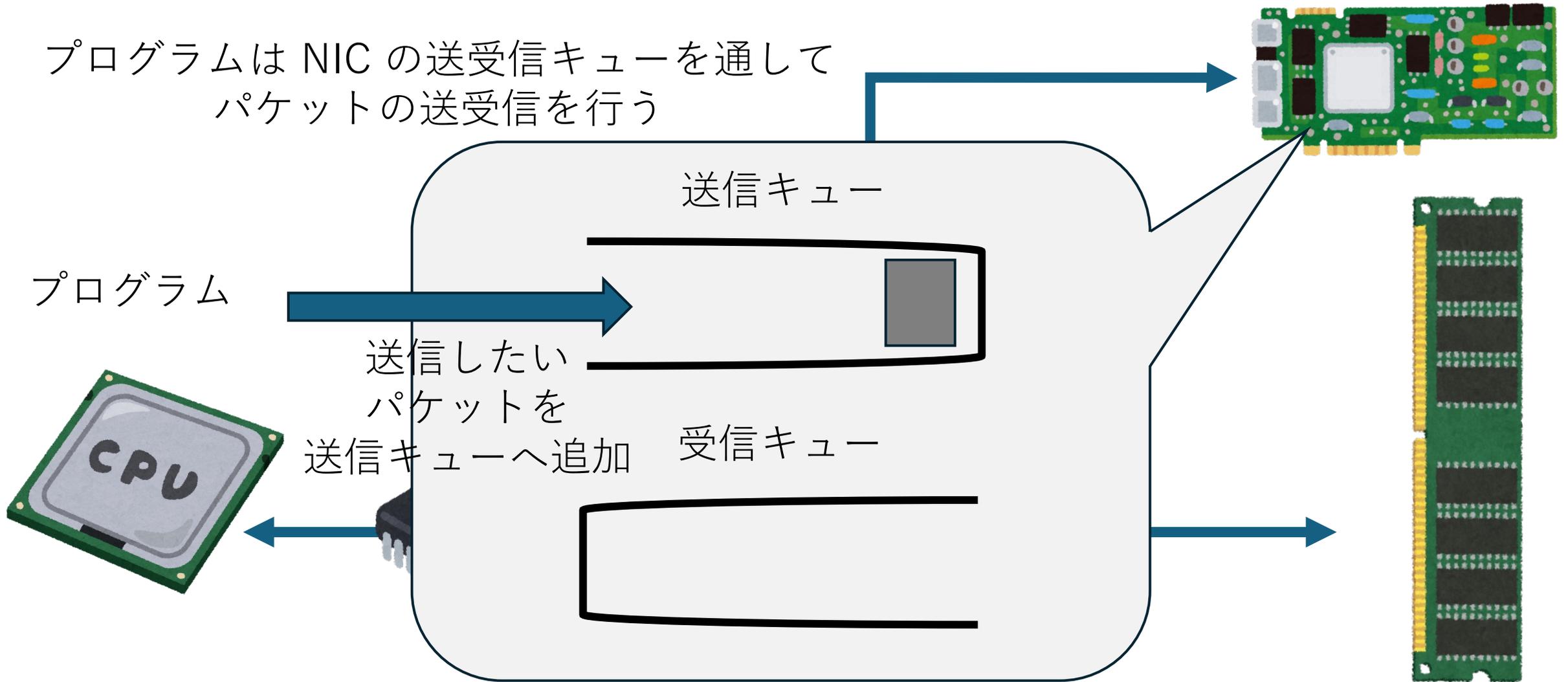
プログラムによる NIC の取り扱い

プログラムは NIC の送受信キューを通して
パケットの送受信を行う



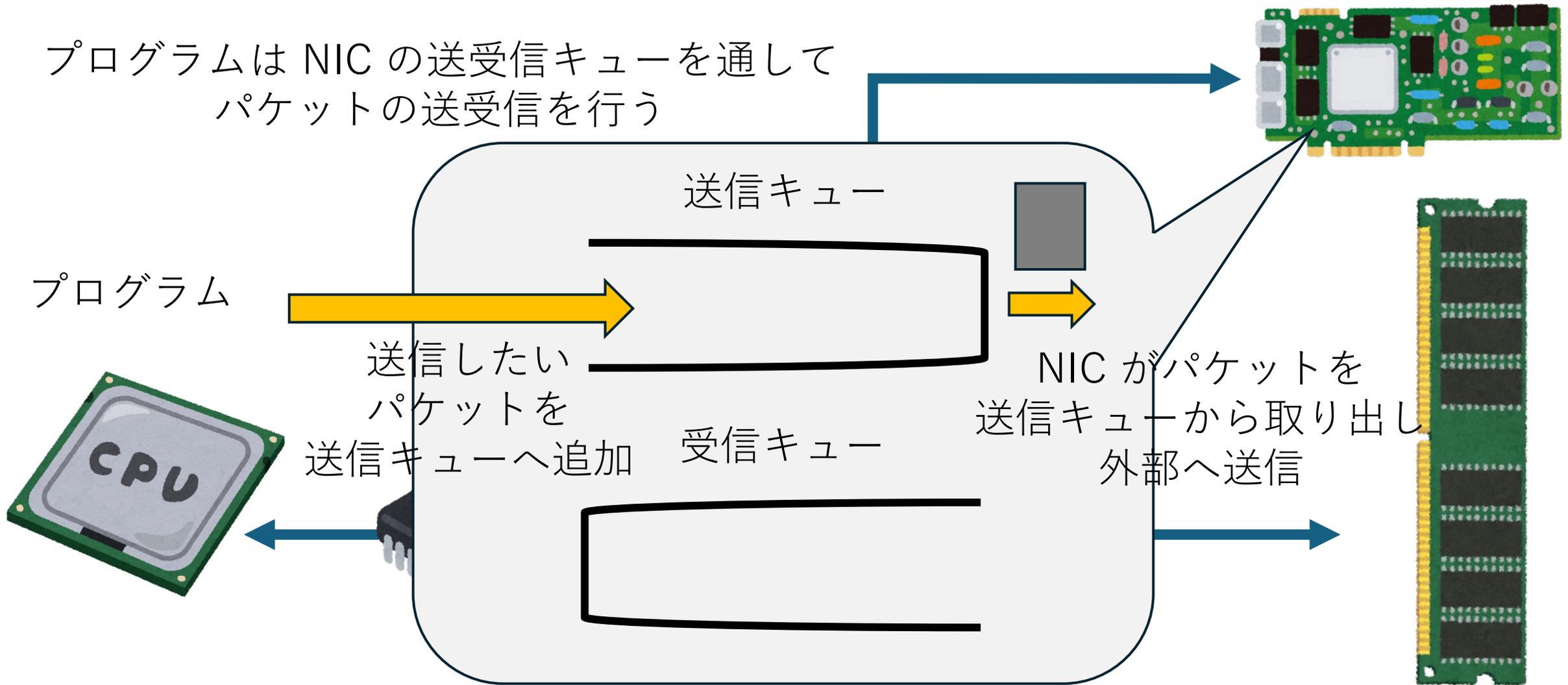
プログラムによる NIC の取り扱い

プログラムは NIC の送受信キューを通して
パケットの送受信を行う



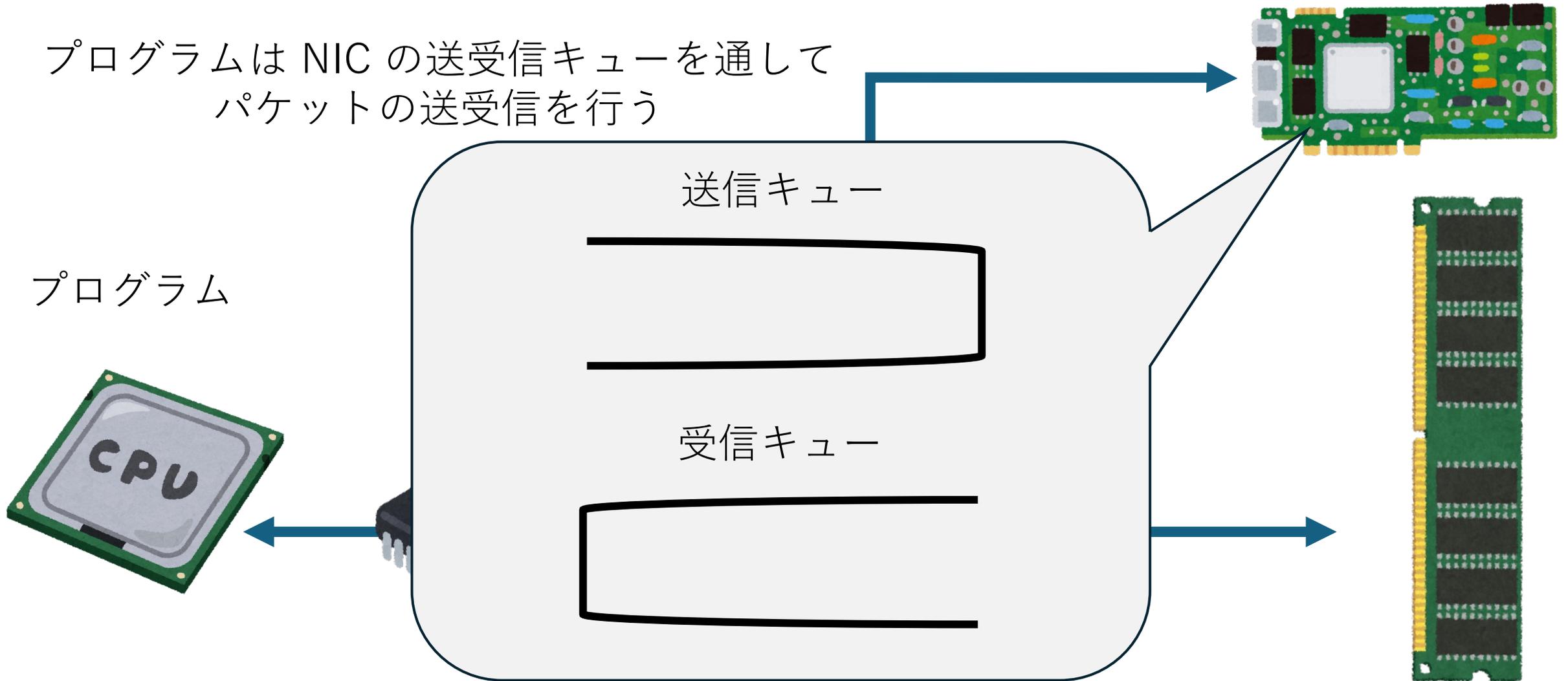
プログラムによる NIC の取り扱い

プログラムは NIC の送受信キューを通して
パケットの送受信を行う



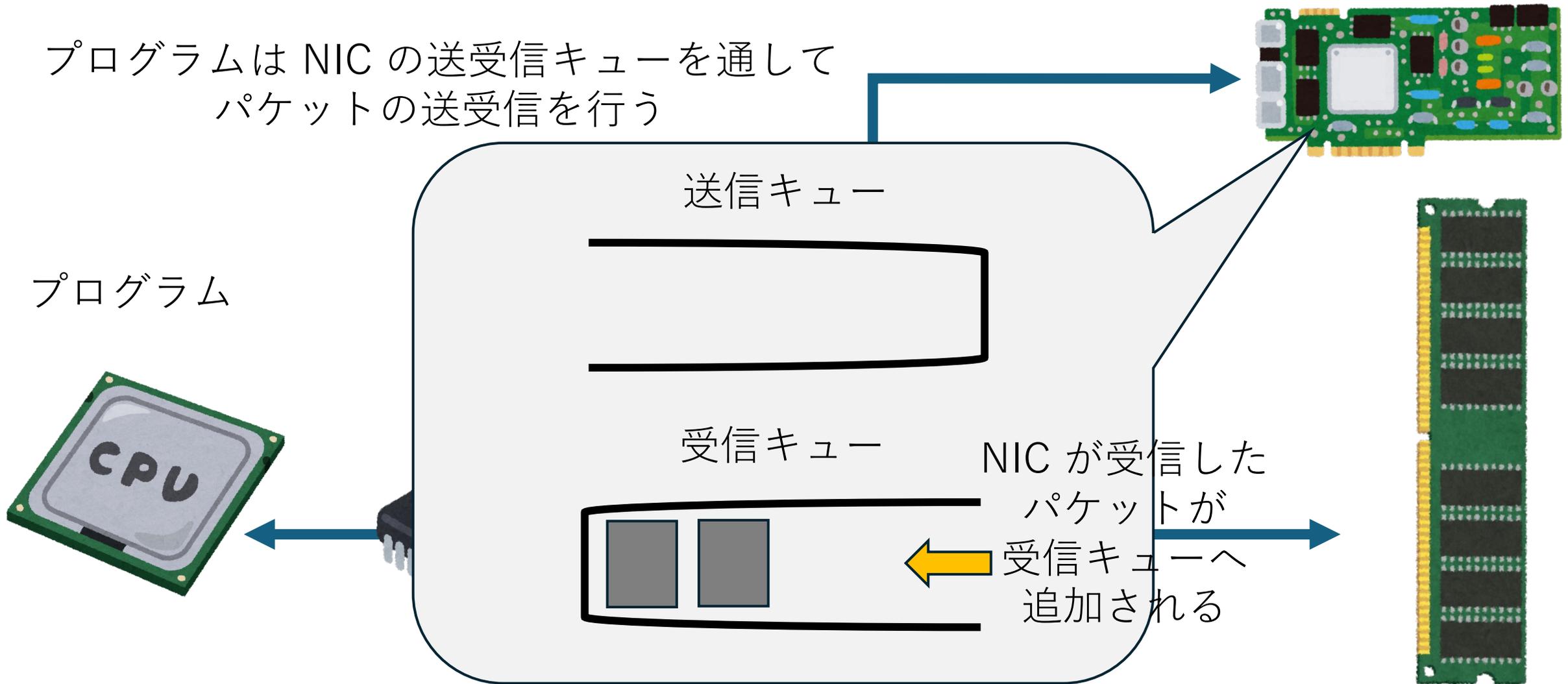
プログラムによる NIC の取り扱い

プログラムは NIC の送受信キューを通して
パケットの送受信を行う



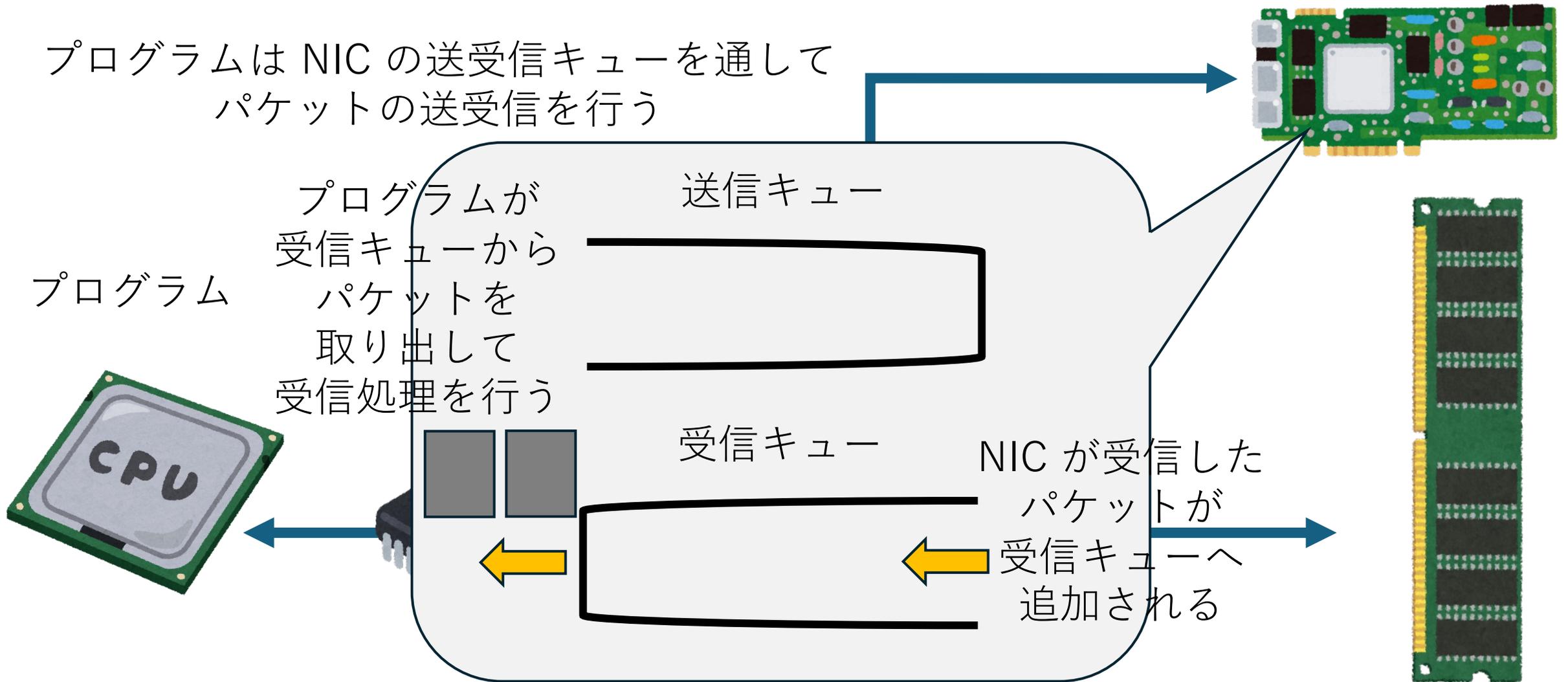
プログラムによる NIC の取り扱い

プログラムは NIC の送受信キューを通して
パケットの送受信を行う



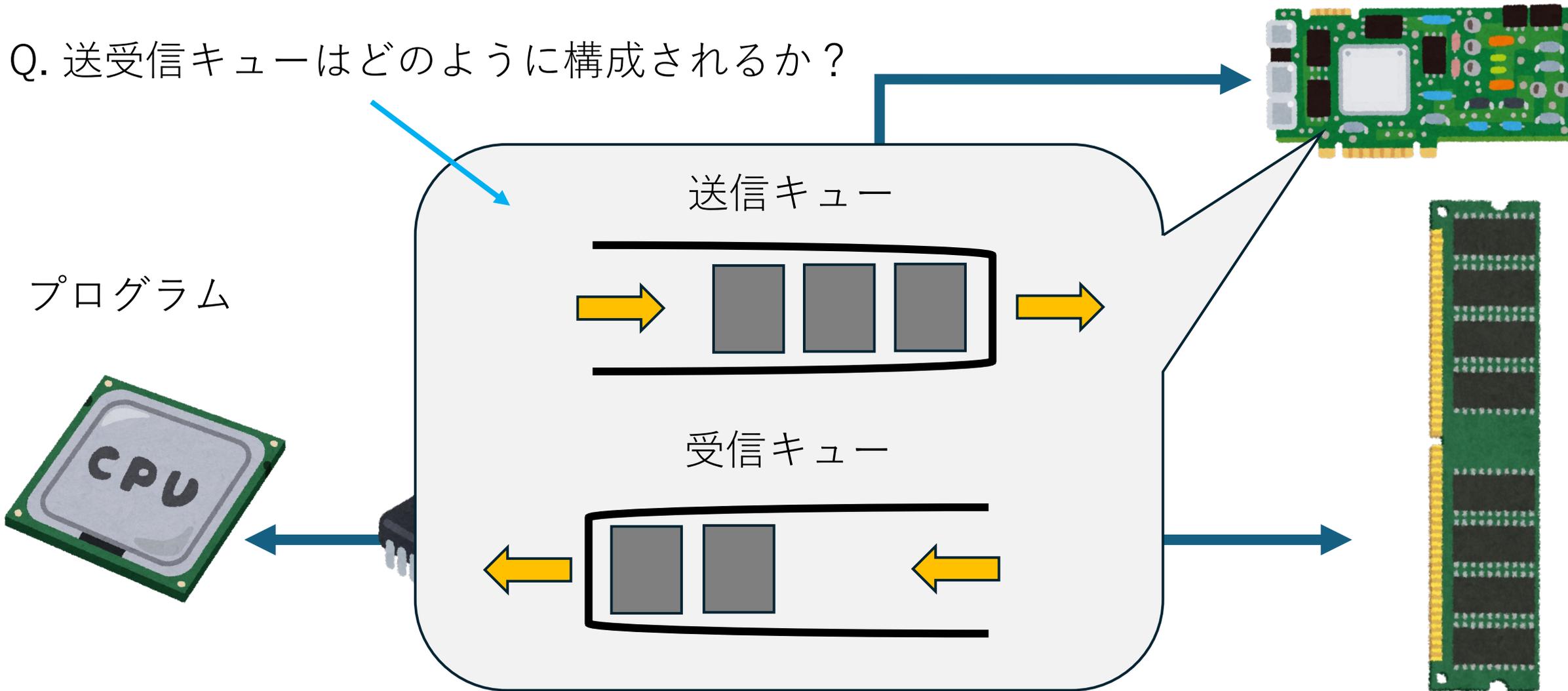
プログラムによる NIC の取り扱い

プログラムは NIC の送受信キューを通して
パケットの送受信を行う



プログラムによる NIC の取り扱い

Q. 送受信キューはどのように構成されるか？

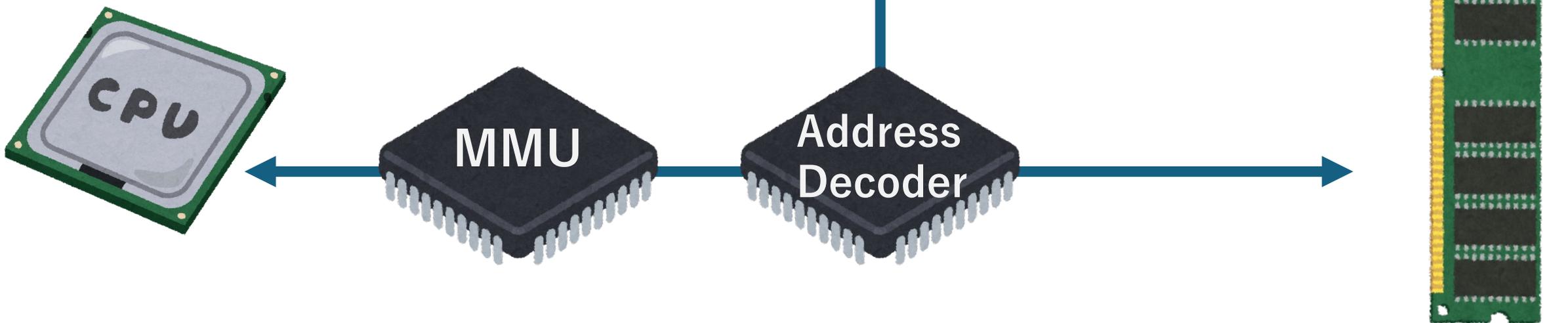


プログラムによる NIC の取り扱い

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム



プログラムによる NIC の取り扱い

Q. 送受信キューはどのように構成されるか？

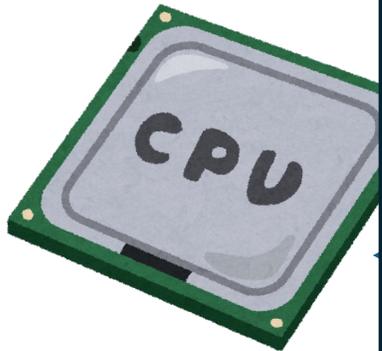
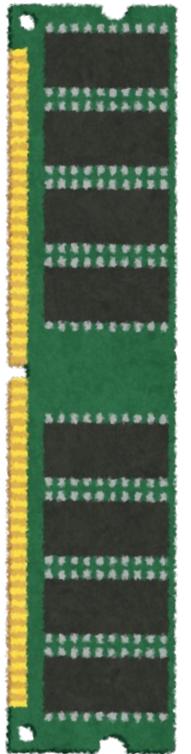
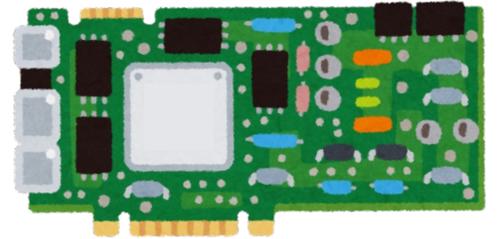
A. キューは NIC のレジスタとメモリ上のデータで構成される リングバッファ

プログラム

簡単な実装

```
int head;  
int tail;  
#define NUM_SLOT 4  
struct {  
    void *address;  
    int length;  
} slot[NUM_SLOT];
```

Address Decoder



プログラムによる NIC の取り扱い

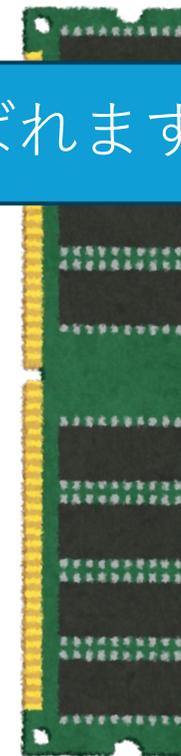
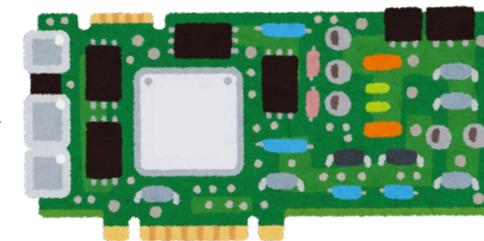
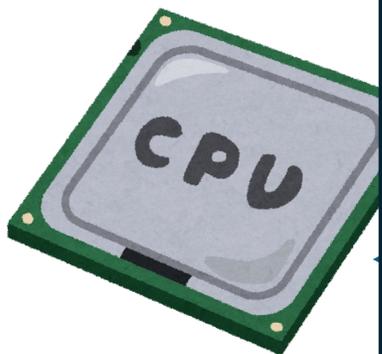
Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```



デスクリプタリングとも呼ばれます

デスクリプタとも呼ばれます

プログラムによる NIC の取り扱い

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

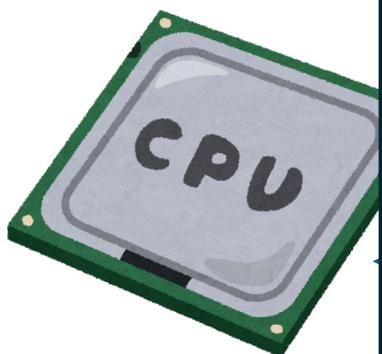
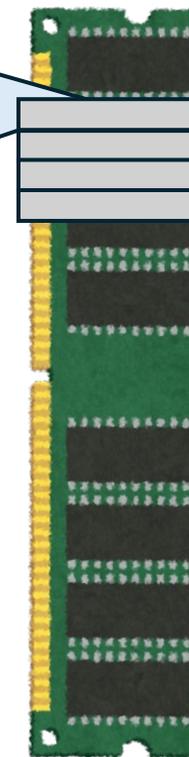
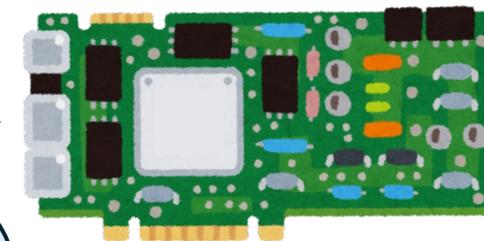
プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length



プログラムによる NIC

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

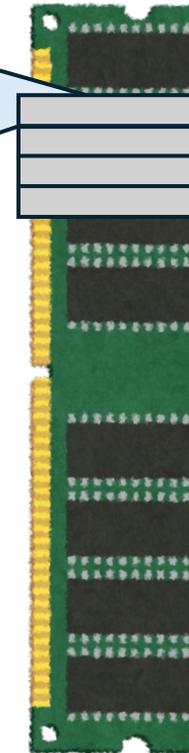
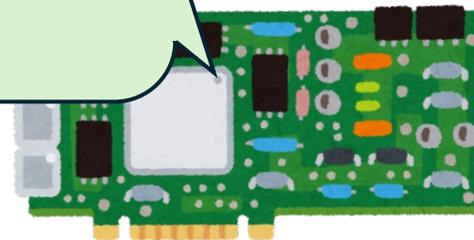
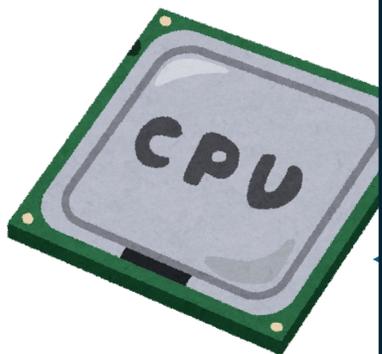
簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

NIC のレジスタ
ring_head
ring_tail
ring_address
ring_size

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length



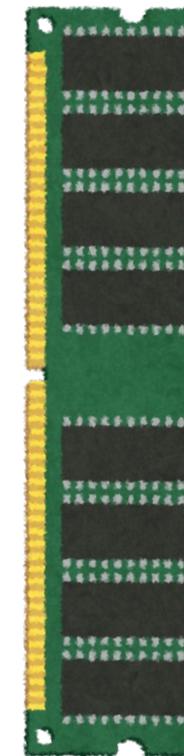
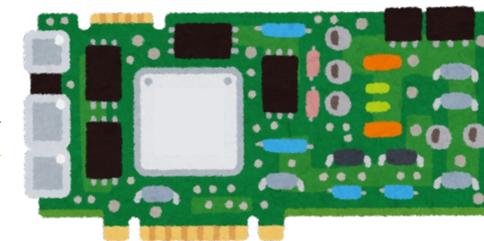
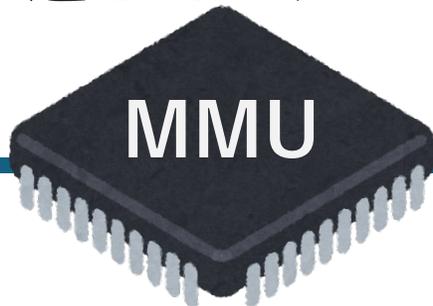
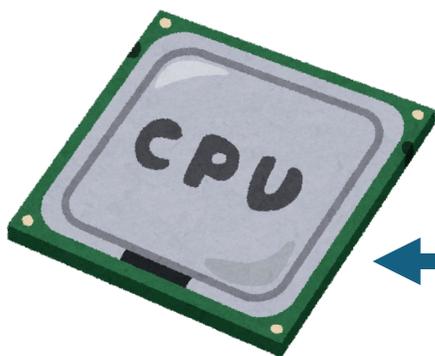
プログラムによる NIC の取り扱い

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラムから
**NIC のレジスタと DRAM は
MMU / Address Decoder を
通じてアクセス可能**

アドレスデコーダという
ハードウェアが
メモリアドレスに応じて
信号の送出先を振り分ける



(どのアドレスが NIC のどのレジスタに対応するかは NIC の仕様に依存)

これらレジスタは NIC が扱うことのできる
キュー（リングバッファ）の数だけあります

NIC のレジスタ
ring_head
ring_tail
ring_address
ring_size

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上の
データで構成されるリングバッファ

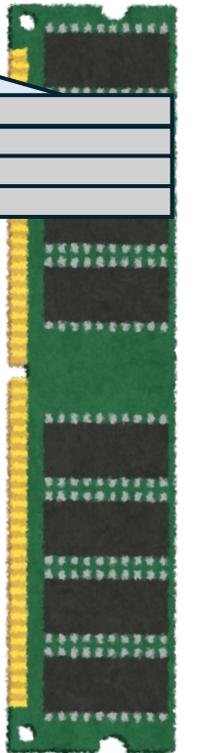
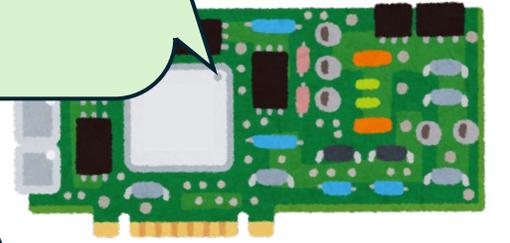
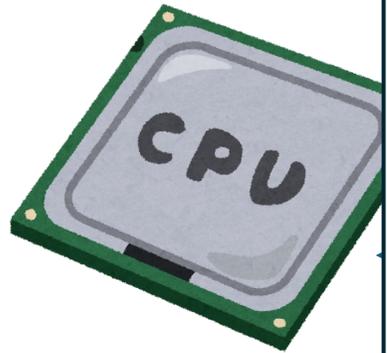
プログラム

簡単な実装

```
int head; }  
int tail; }  
#define NUM_SLOT 4  
struct {  
    void *address;  
    int length;  
} slot[NUM_SLOT];
```

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length



これらレジスタは NIC が扱うことのできる
キュー（リングバッファ）の数だけあります

NIC のレジスタ
ring_head
ring_tail
ring_address
ring_size

Q 送受信キューはどのように構成されるか？

A キューは基本的に
送信キューもしくは受信キューです

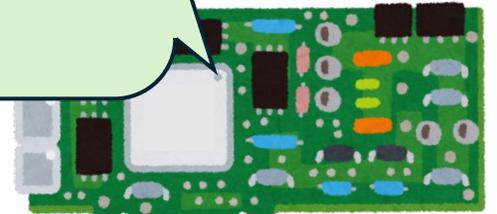
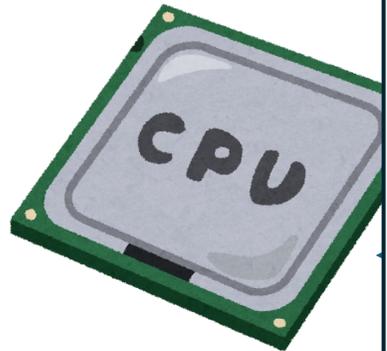
プログラム

簡単な実装

```
int head; }  
int tail; }  
#define NUM_SLOT 4  
struct {  
    void *address; }  
    int length; }  
} slot[NUM_SLOT];
```

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length



プログラムによる NIC

プログラムによる初期設定

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

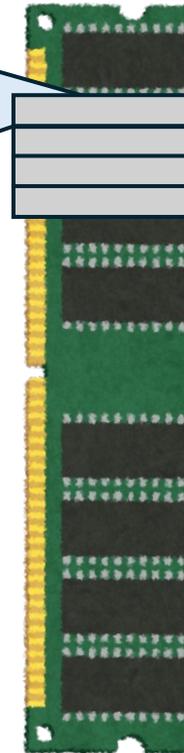
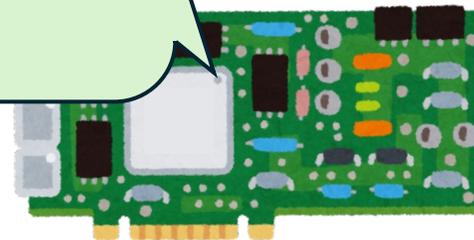
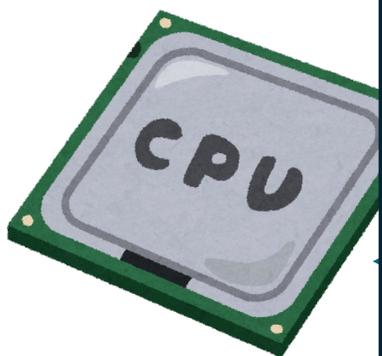
```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

NIC のレジスタ

```
ring_head
ring_tail
ring_address
ring_size
```

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length



プログラムによる NIC の初期設定

プログラムによる初期設定

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

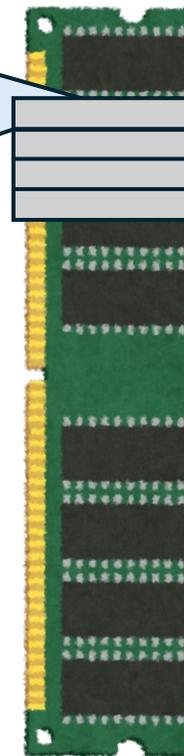
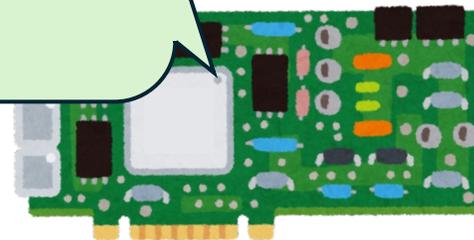
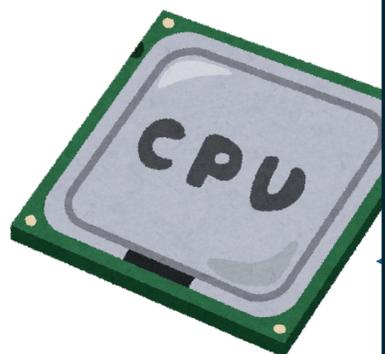
NIC のレジスタ

```
ring_head
ring_tail
ring_address
ring_size: 4
```

リングサイズを設定

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length



プログラムによる NIC

プログラムによる初期設定

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

NIC のレジスタ

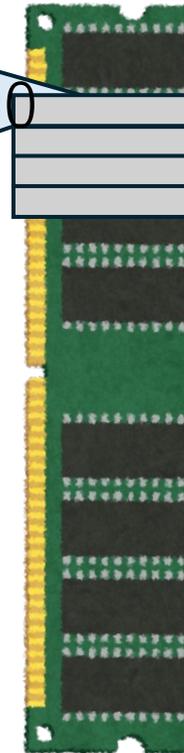
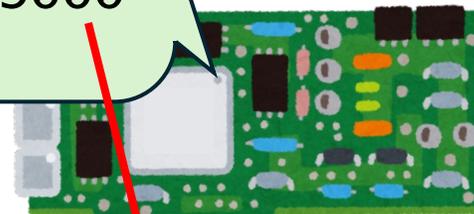
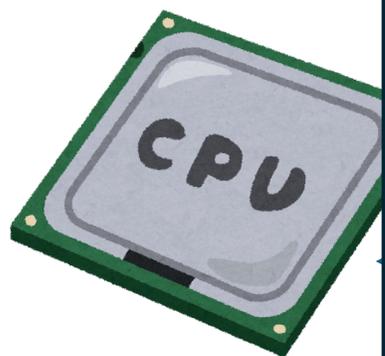
```
ring_head
ring_tail
ring_address:0x5000
ring_size: 4
```

DRAM 上の配列の先頭の物理アドレスを登録

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length

0x5000



プログラムによる NIC

プログラムによる初期設定

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head
tail

DRAM 上の配列

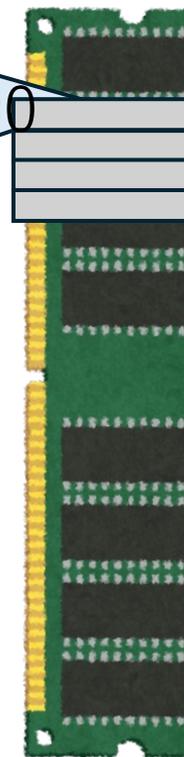
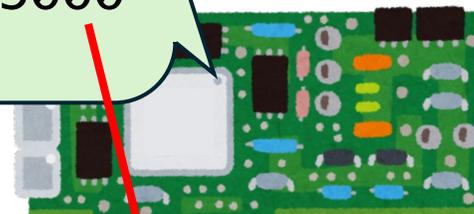
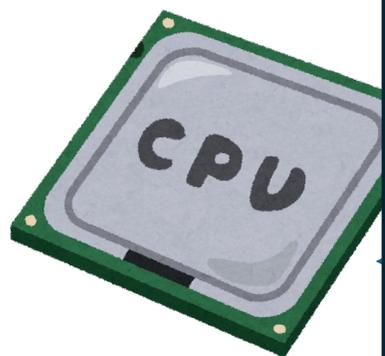
[0]	address length
[1]	address length
[2]	address length
[3]	address length

0x5000

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

head と tail は
0 に設定



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head
tail

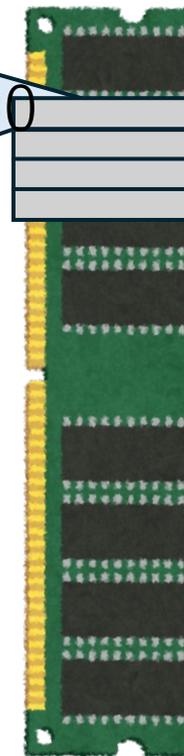
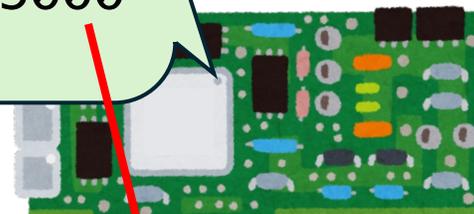
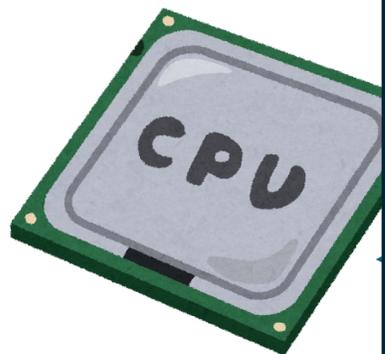
DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length

0x5000

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

head
tail

DRAM 上の配列

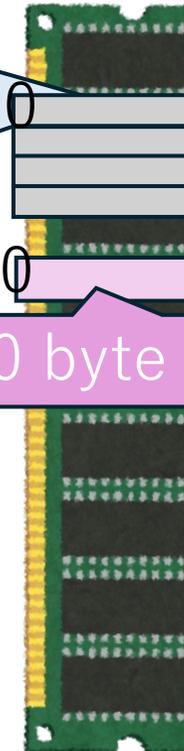
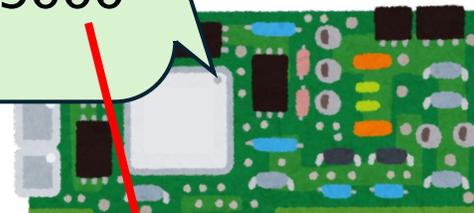
[0]	address length	0x5000
[1]	address length	0x30000
[2]	address length	
	address length	

500 byte

プログラムが DRAM 上に送信データを用意：例
物理アドレス 0x30000 に 500 byte のデータ

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

head
tail

DRAM 上の配列

[0]	address length	0x5000
[1]	address length	0x30000
[2]	address length	0x40000
	address length	

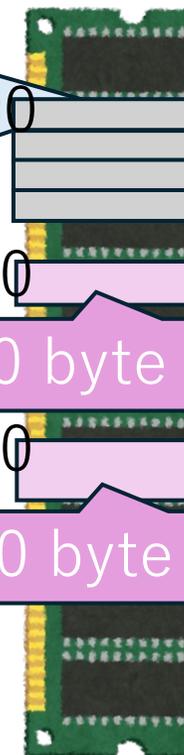
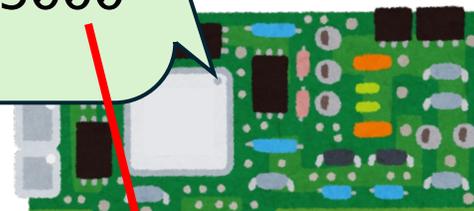
500 byte

800 byte

プログラムが DRAM 上に送信データを用意：例
物理アドレス 0x30000 に 500 byte のデータ
物理アドレス 0x40000 に 800 byte のデータ

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通して配列に送信データへの参照を設定

```
} slot[NUM_SLOT];
```

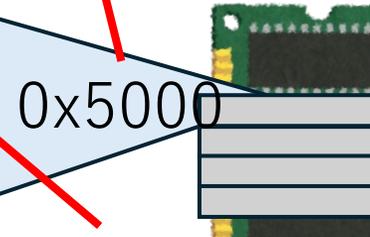
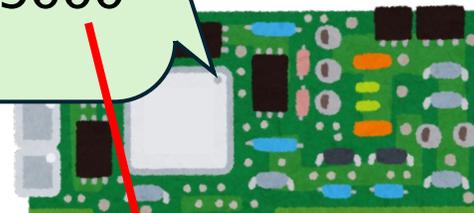
head
tail

DRAM 上の配列

[0]	address: 0x30000 length
[1]	address length
[2]	address length
[3]	address length

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



500 byte

800 byte

0x5000

0x30000

0x40000

プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通して配列に送信データへの参照を設定

```
} slot[NUM_SLOT];
```

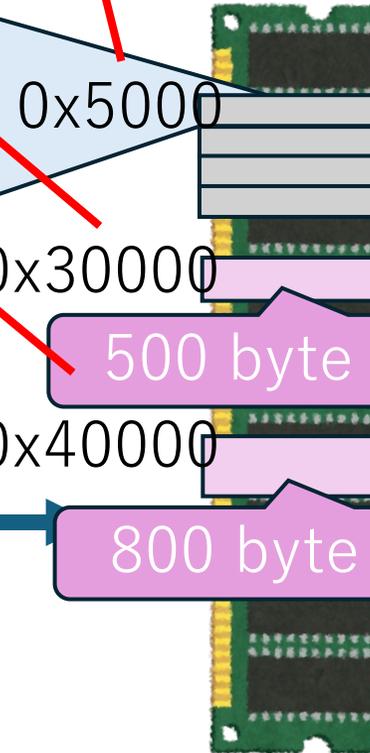
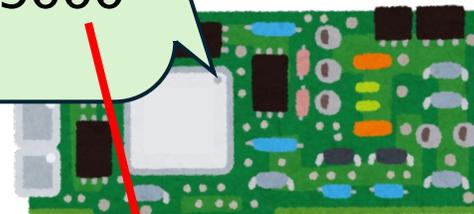
head
tail

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address length
[2]	address length
[3]	address length

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通して配列に送信データへの参照を設定

```
} slot[NUM_SLOT];
```

head
tail

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length
[2]	address length
[3]	address length

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

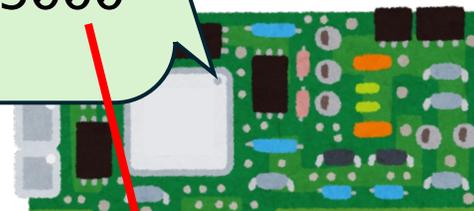
0x5000

0x30000

500 byte

0x40000

800 byte



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通して配列に送信データへの参照を設定

```
} slot[NUM_SLOT];
```

head
tail

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

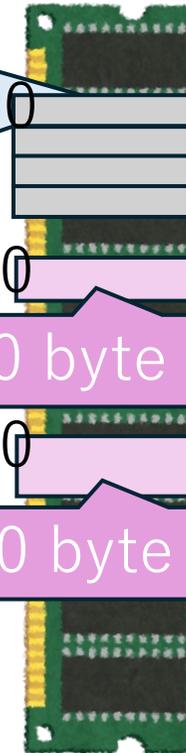
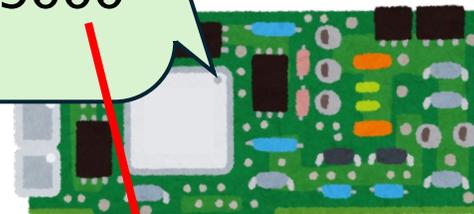
0x5000

0x30000

500 byte

0x40000

800 byte



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ値を書き込むことでパケットの送信開始をリクエストする

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

head

tail

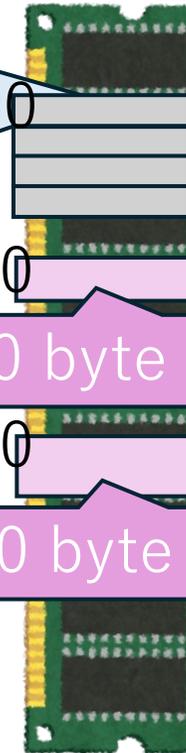
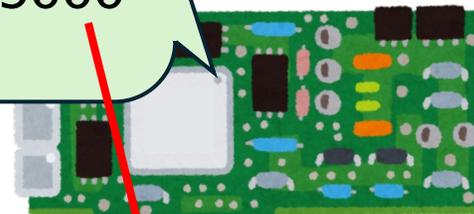
0x5000

0x30000

500 byte

0x40000

800 byte



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ値を書き込むことでパケットの送信開始をリクエストする

head

tail

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

NIC のレジスタ

```
ring_head: 0
ring_tail: 2
ring_address: 0x5000
ring_size: 4
```

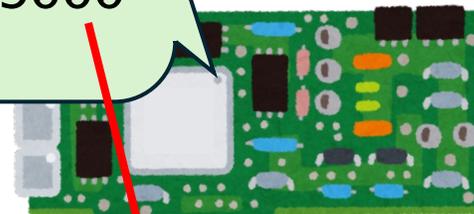
0x5000

0x30000

500 byte

0x40000

800 byte



このキュー

Q. 送受信

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ値を書き込むことでパケットの送信開始をリクエストする

NIC は配列中の head と tail の間の区間が参照する DRAM 上のデータを外部へ送信する

NIC のレジスタ

```
ring_head: 0
ring_tail: 2
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

head

tail

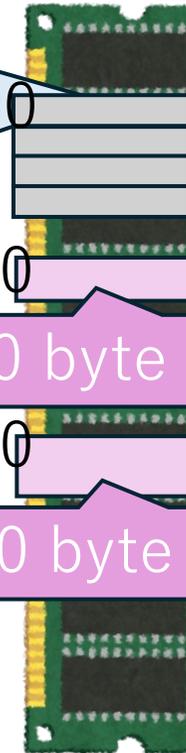
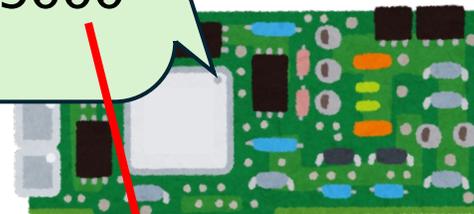
0x5000

0x30000

500 byte

0x40000

800 byte



このキュー

Q. 送受信

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ値を書き込むことでパケットの送信開始をリクエストする

NIC は配列中の head と tail の間の区間が参照する DRAM 上のデータを外部へ送信する

NIC のレジスタ

```
ring_head: 0
ring_tail: 2
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

head

tail

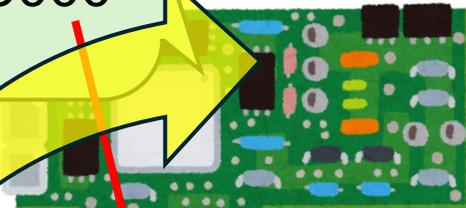
0x5000

0x30000

500 byte

0x40000

800 byte



このキュー

Q. 送受信

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ値を書き込むことでパケットの送信開始をリクエストする

NIC は配列中の head と tail の間の区間が参照する DRAM 上のデータを外部へ送信する

NIC のレジスタ

```
ring_head: 0
ring_tail: 2
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

head

tail

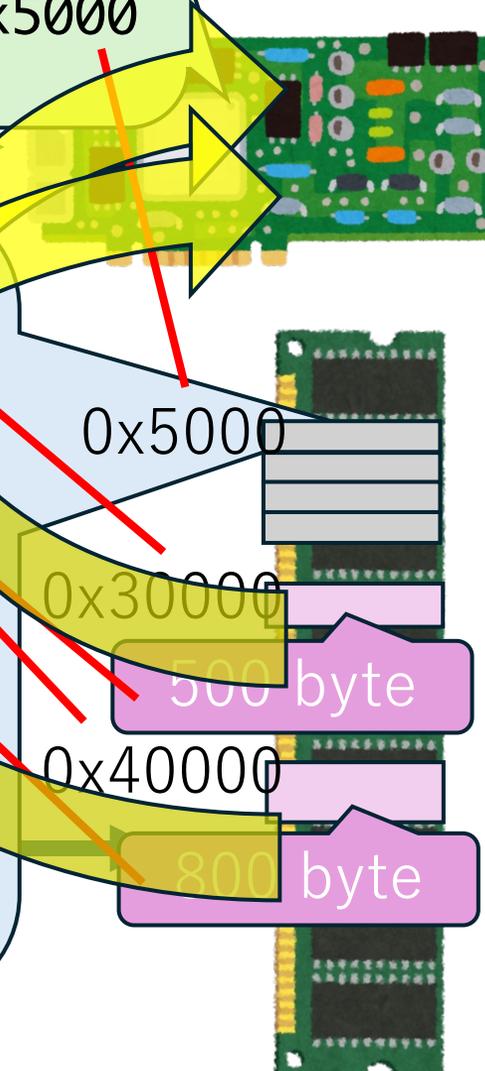
0x5000

0x30000

500 byte

0x40000

800 byte



プログラマー
このキューが送

NIC は送信完了後に
レジスタの値を更新する

NIC のレジスタ
ring_head: 0
ring_tail: 2
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上の
データで構成されるリングバッファ

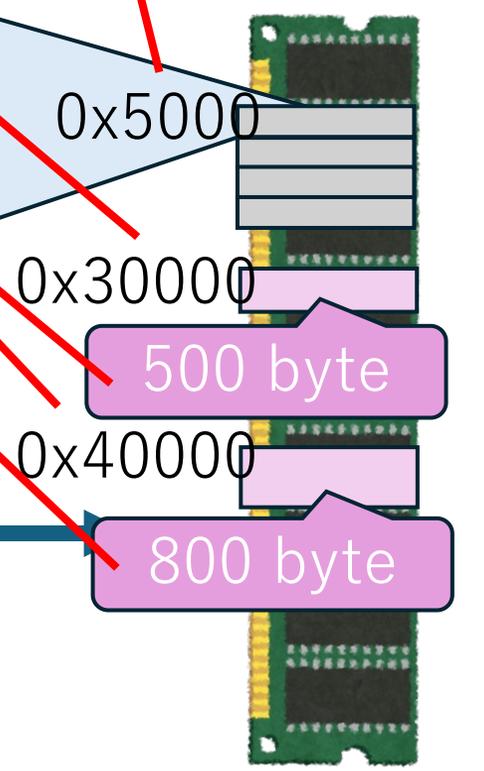
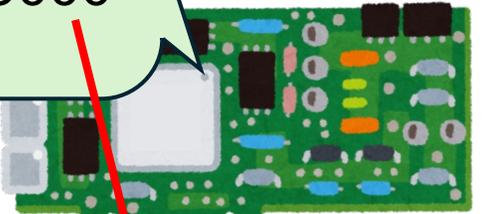
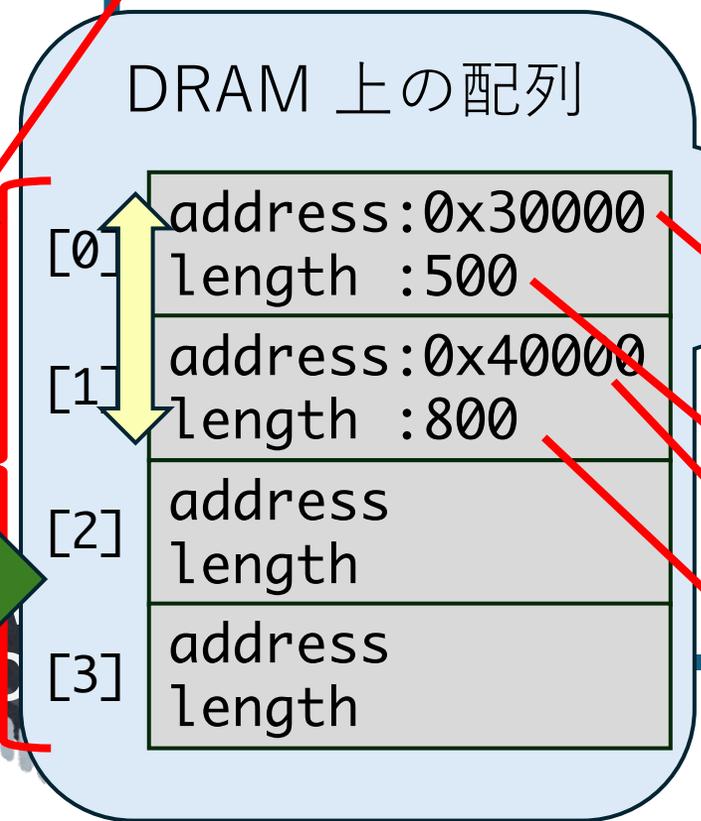
プログラム

簡単な実装
int head;
int tail;
#define NUM_SLOT 4

プログラムは NIC のレジスタへ
値を書き込むことでパケットの
送信開始をリクエストする

head

tail



500 byte

800 byte

プログラマー
このキューが送

NIC は送信完了後に
レジスタの値を更新する

NIC のレジスタ
ring_head: 2
ring_tail: 2
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上の
データで構成されるリングバッファ

プログラム

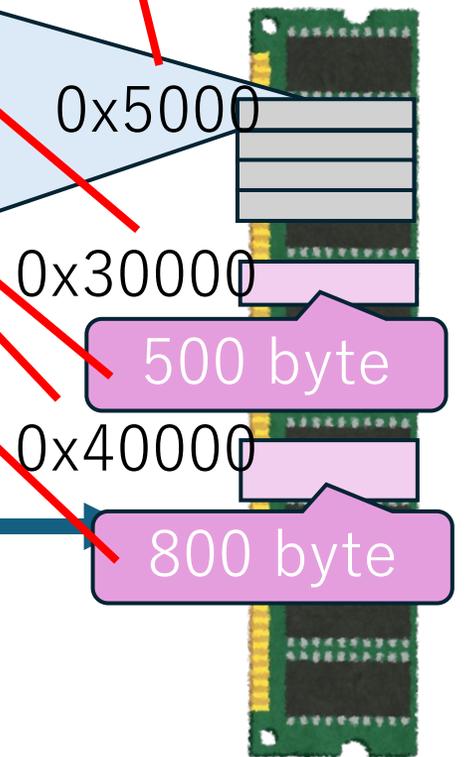
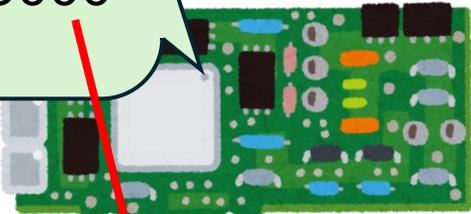
簡単な実装
int head;
int tail;
#define NUM_SLOT 4

プログラムは NIC のレジスタへ
値を書き込むことでパケットの
送信開始をリクエストする

head
tail

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length



0x5000

0x30000

0x40000

500 byte

800 byte

プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタの値を読み込むことでパケットの送信完了を確認する

NIC のレジスタ

```
ring_head: 2
ring_tail: 2
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

head
tail

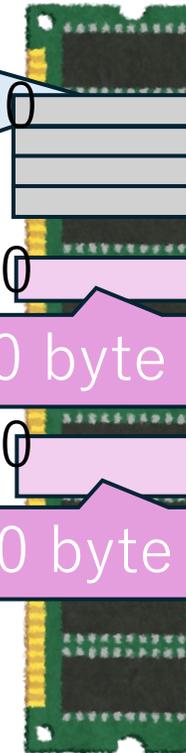
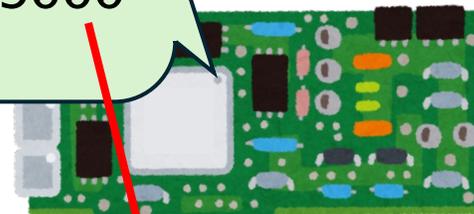
0x5000

0x30000

500 byte

0x40000

800 byte



プログラムによる NIC

このキ

Q. 送受

A. キュ

デー

プログラムは送信が完了したデータが配置されている DRAM を別のデータを配置するために利用して良いと判断できる

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタの値を読み込むことでパケットの送信完了を確認する

head
tail

NIC のレジスタ

```
ring_head: 2
ring_tail: 2
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

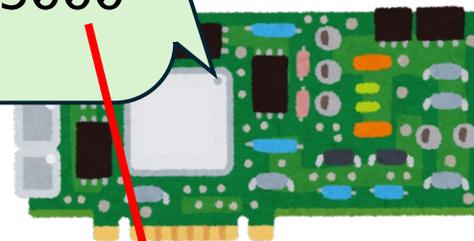
0x5000

0x30000

0x40000

500 byte

800 byte



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head
tail

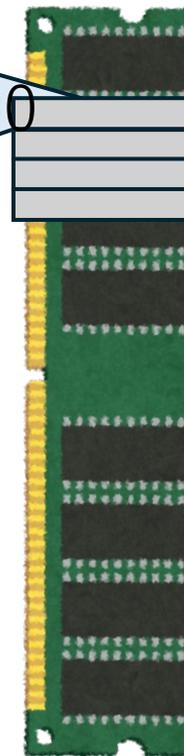
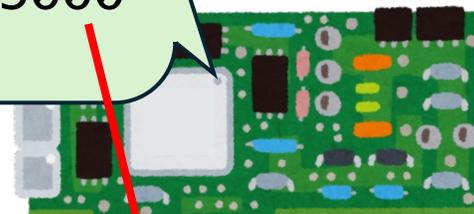
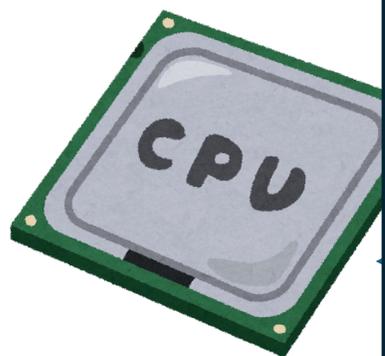
DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length

0x5000

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

NIC のレジスタ

```
ring_head: 0  
ring_tail: 0  
ring_address: 0x5000  
ring_size: 4
```

DRAM 上の配列

[0]	address length
[1]	address length
[2]	address length
[3]	address length

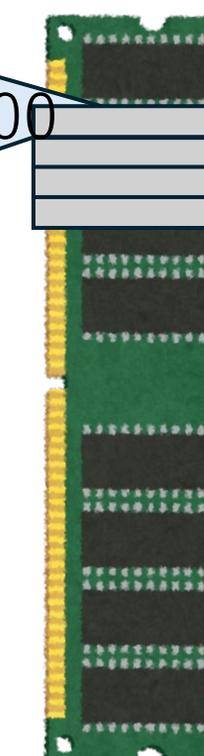
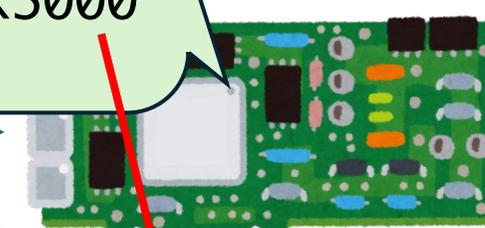
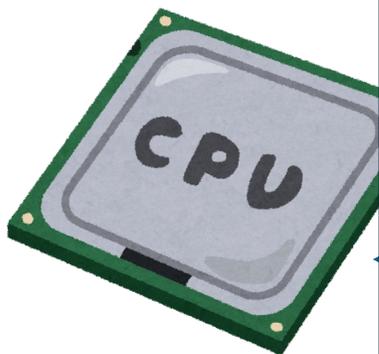
0x5000

プログラム

簡単な実装

```
int head;  
int tail;  
#define NUM_SLOT 4  
struct {  
    void *address;  
    int length;  
} slot[NUM_SLOT];
```

head
tail



プログラムによる NIC

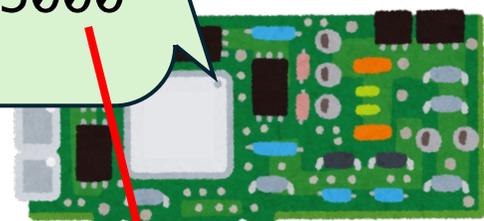
このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

NIC のレジスタ

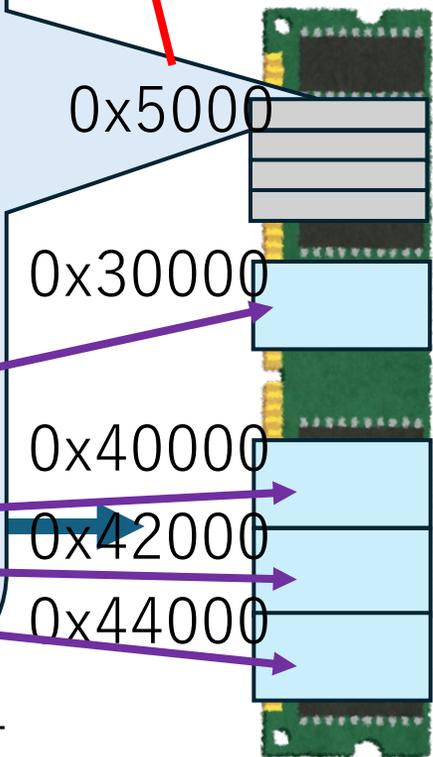
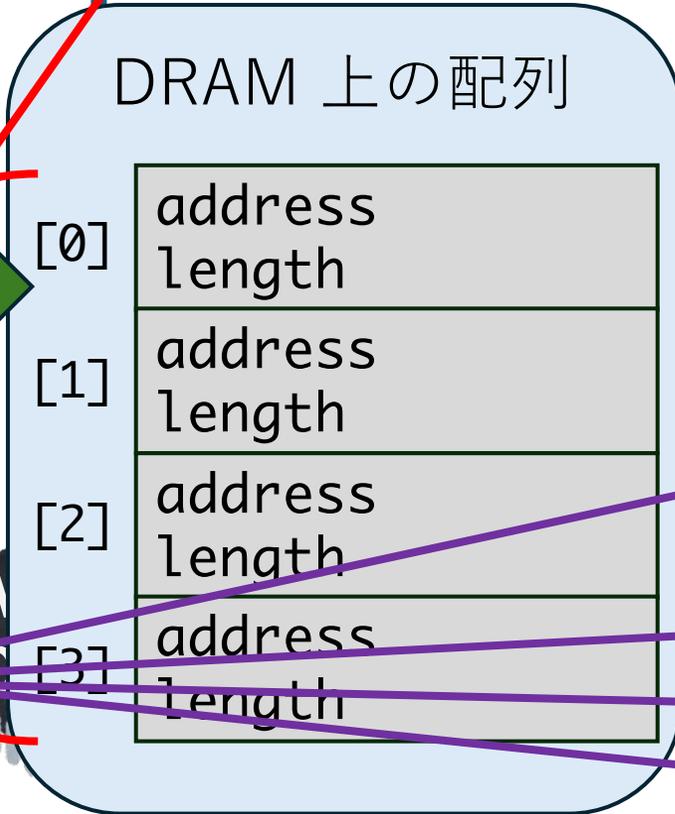
- ring_head: 0
- ring_tail: 0
- ring_address: 0x5000
- ring_size: 4



プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
} slot[NUM_SLOT];
```



プログラムは NIC が受信したデータを配置するための DRAM 領域を確保

連続的でも大丈夫です

プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通じて配列へ参照を設定することで、NIC と受信データ配置用の DRAM 領域を紐付ける

head
tail

DRAM 上の配列

[0]	address: 0x30000 length
[1]	address length
[2]	address length
[3]	address length

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

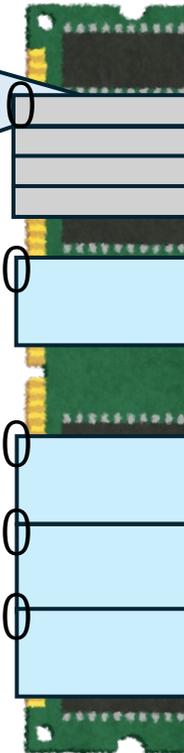
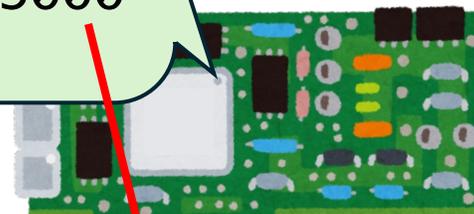
0x5000

0x30000

0x40000

0x42000

0x44000



プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
enum NUM_SLOT 4
```

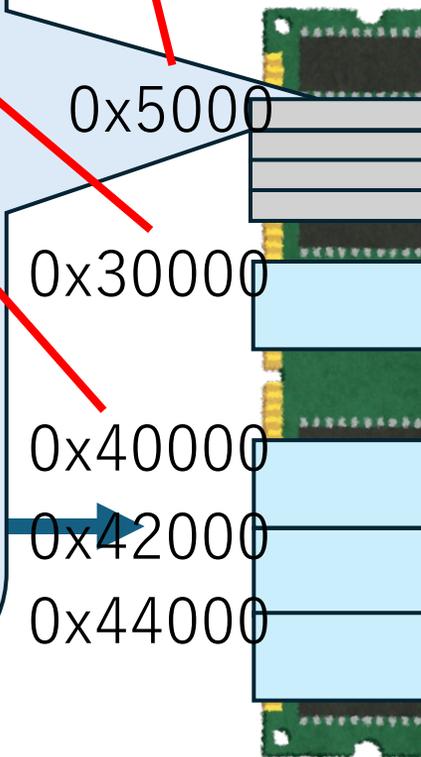
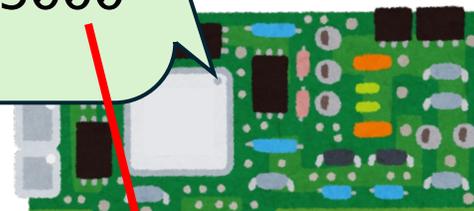
プログラムは DRAM への書き込みを通じて配列へ参照を設定することで、NIC と受信データ配置用の DRAM 領域を紐付ける

head
tail

DRAM 上の配列

[0]	address: 0x30000 length
[1]	address: 0x40000 length
[2]	address length
[3]	address length

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4



プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;  
int tail;  
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通じて配列へ参照を設定することで、NIC と受信データ配置用の DRAM 領域を紐付ける

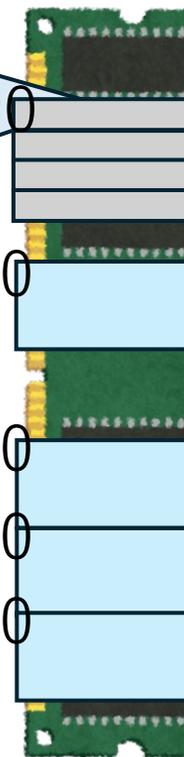
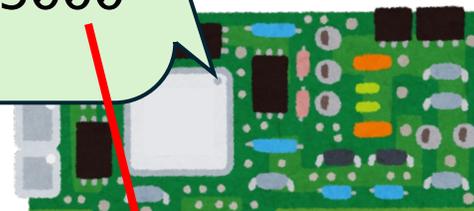
head
tail

DRAM 上の配列

[0]	address: 0x30000 length	0x5000
[1]	address: 0x40000 length	0x30000
[2]	address: 0x42000 length	0x40000
[3]	address length	0x42000 0x44000

NIC のレジスタ

```
ring_head: 0  
ring_tail: 0  
ring_address: 0x5000  
ring_size: 4
```



プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通じて配列へ参照を設定することで、NIC と受信データ配置用の DRAM 領域を紐付ける

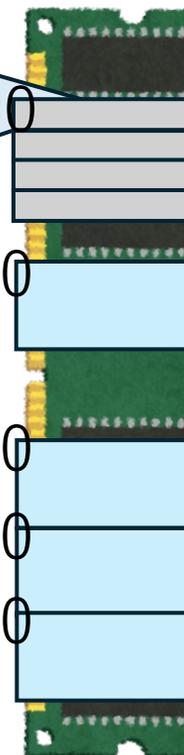
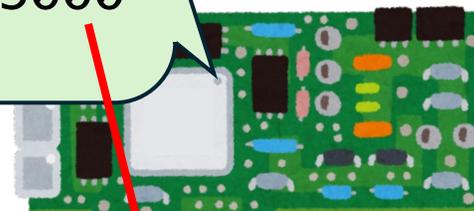
head
tail

DRAM 上の配列

[0]	address: 0x30000 length	0x5000
[1]	address: 0x40000 length	0x30000
[2]	address: 0x42000 length	0x40000
[3]	Address: 0x44000 length	0x42000 0x44000

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head

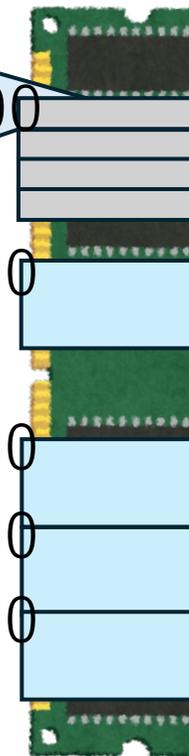
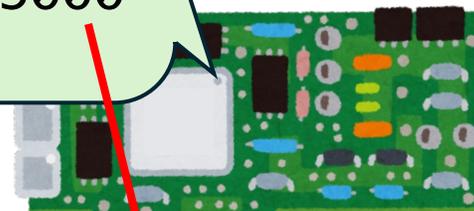
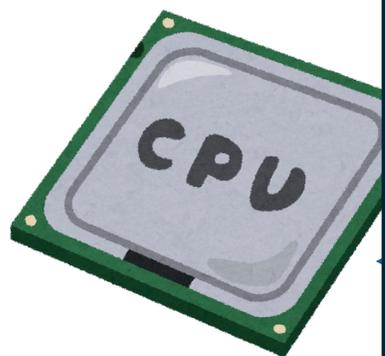
tail

DRAM 上の配列

[0]	address: 0x30000 length	0x5000
[1]	address: 0x40000 length	0x30000
[2]	address: 0x42000 length	0x40000
[3]	Address: 0x44000 length	0x42000 0x44000

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムのこのキューが受信

NIC が受信したパケットを配列の head で参照される DRAM 領域へ書き込む

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成される？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head
tail

DRAM 上の配列

[0]	address: 0x30000 length
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

0x5000

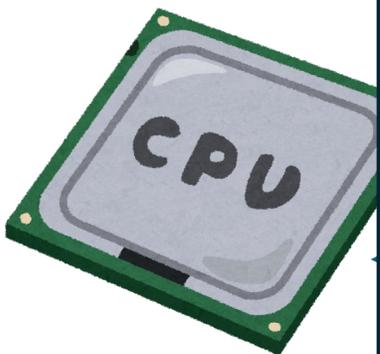
0x30000

100 byte

0x40000

0x42000

0x44000



プログラマー
このキュー

NIC は配列の長さを表すフィールド (length) に受信したデータの長さを設定

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリー上のデータで構成されるリングバッファ

プログラム

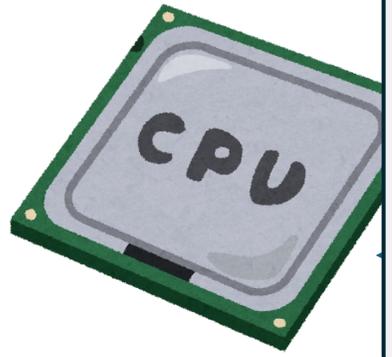
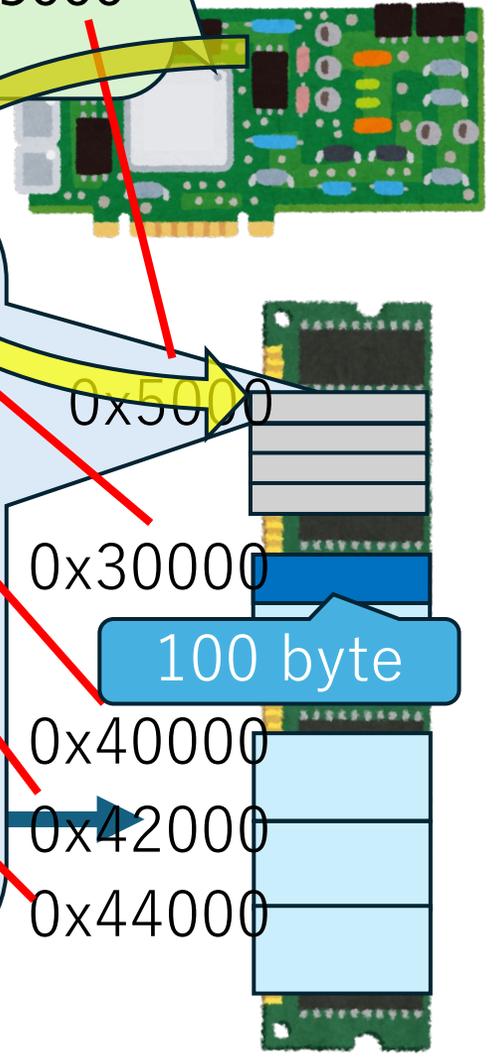
簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```



DRAM 上の配列

[0]	address: 0x30000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length



プログラム

NIC は head の値を進める

NIC のレジスタ

```

ring_head: 1
ring_tail: 0
ring_address: 0x5000
ring_size: 4

```

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

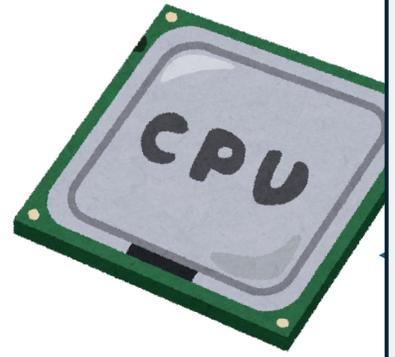
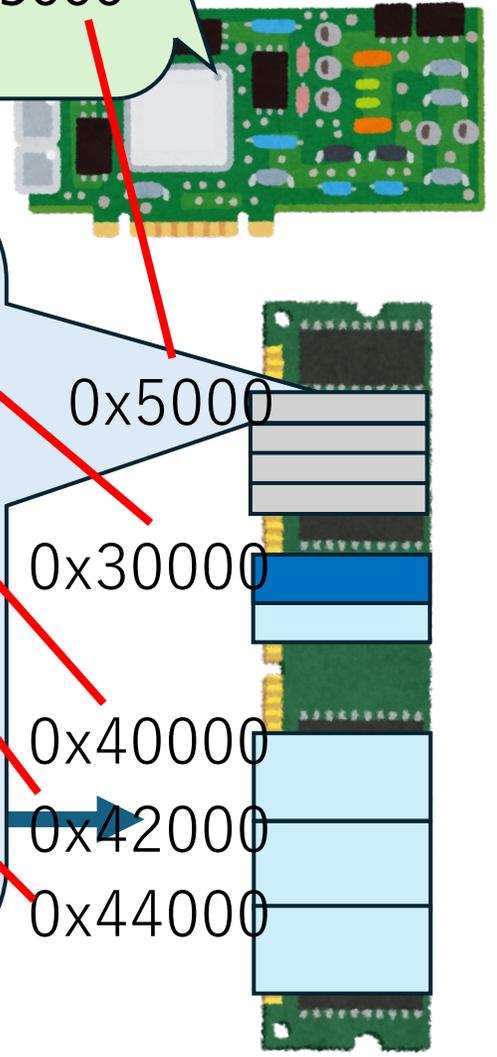
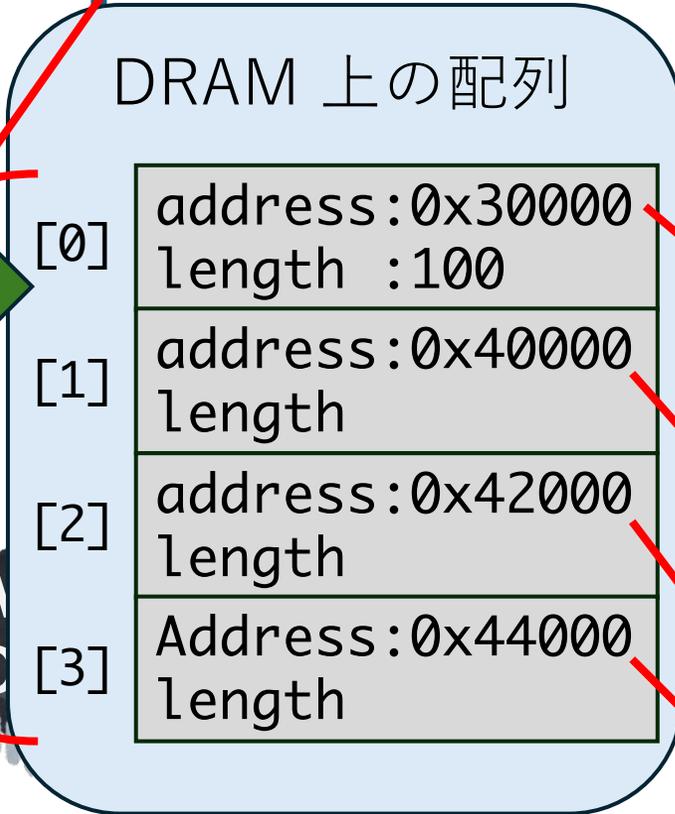
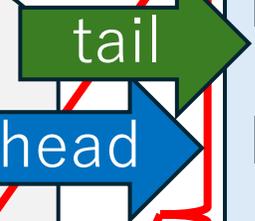
プログラム

簡単な実装

```

int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];

```



プログラムによる NIC の

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

tail

head

プログラムは配列の head と tail の区間で参照される DRAM 上の領域に新しくパケットが受信されていると解釈する

NIC のレジスタ

```
ring_head: 1
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

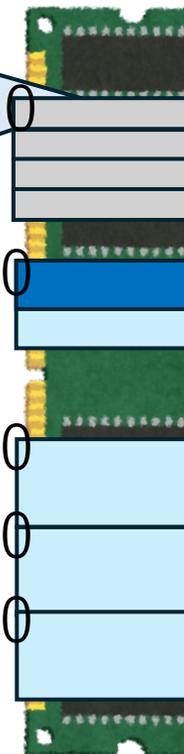
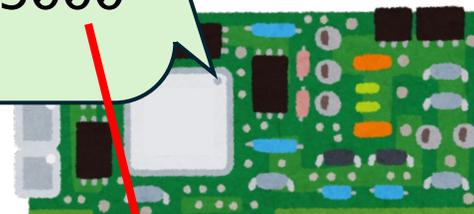
0x5000

0x30000

0x40000

0x42000

0x44000



プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC が受信したデータを配置するための DRAM 領域を確保

```
} slot[NUM_SLOT];
```

tail

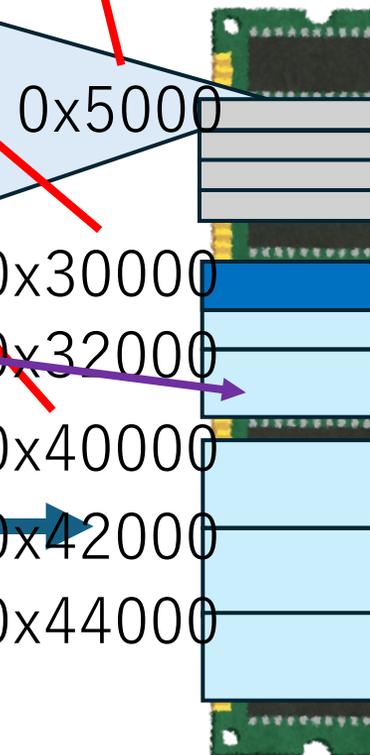
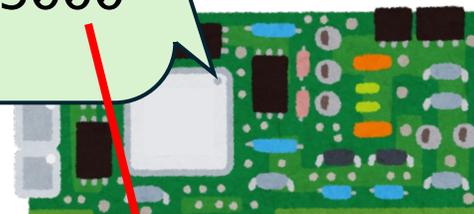
head

DRAM 上の配列

[0]	address: 0x30000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

NIC のレジスタ

```
ring_head: 1
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

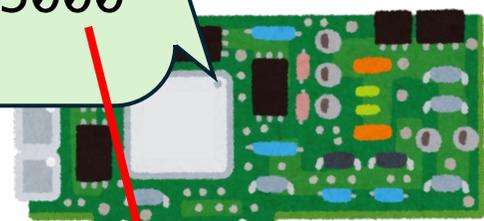
このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

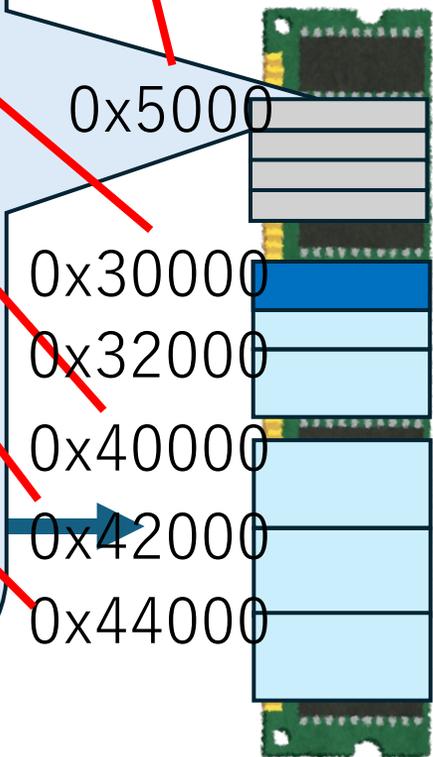
NIC のレジスタ

- ring_head: 1
- ring_tail: 0
- ring_address: 0x5000
- ring_size: 4



DRAM 上の配列

[0]	address: 0x30000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length



プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

Diagram showing 'tail' and 'head' pointers. A green arrow labeled 'tail' points to the start of the first DRAM entry (address 0x30000). A blue arrow labeled 'head' points to the start of the second DRAM entry (address 0x40000).

プログラムは DRAM への書き込みを通じて配列が参照する DRAM 領域を受信データが配置されている DRAM 領域から新しく確保した DRAM 領域へ置き換える

プログラムによる NIC

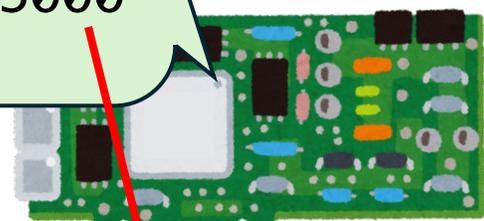
このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

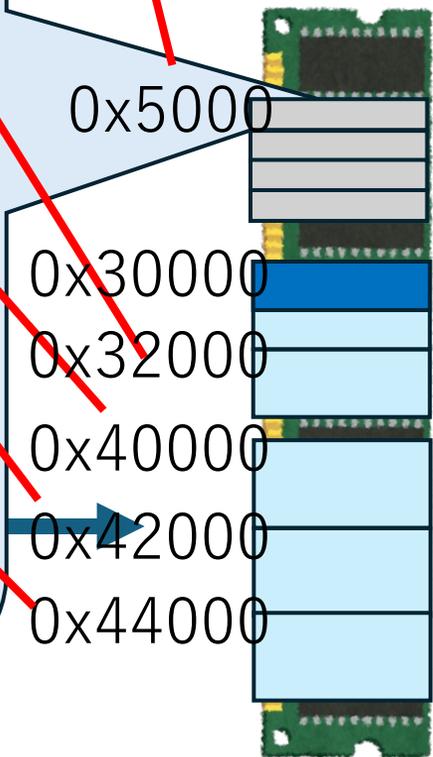
NIC のレジスタ

- ring_head: 1
- ring_tail: 0
- ring_address: 0x5000
- ring_size: 4



DRAM 上の配列

[0]	address: 0x32000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length



簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

tail

head

プログラムは DRAM への書き込みを通じて配列が参照する DRAM 領域を受信データが配置されている DRAM 領域から新しく確保した DRAM 領域へ置き換える

プログラム

プログラムによる NIC の

この
Q. 送
A. キ

これで物理アドレス 0x30000 からの
DRAM 領域は NIC の紐付けが解除され
新たに物理アドレス 0x32000 からの領域が
NIC と紐付けられる

NIC のレジスタ
ring_head: 1
ring_tail: 0
ring_address: 0x5000
ring_size: 4

データで構成されるリングバッファ

DRAM 上の配列

プログラム

簡単な実装

```
int head;  
int tail;  
#define NUM_SLOT 4
```

tail

head

[0]	address: 0x32000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

0x5000

0x30000

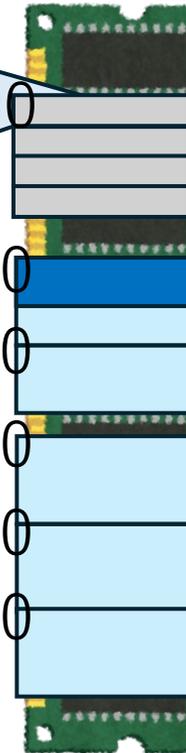
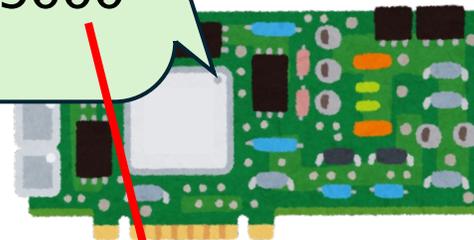
0x32000

0x40000

0x42000

0x44000

プログラムは DRAM への書き込みを通じて
配列が参照する DRAM 領域を受信データが
配置されている DRAM 領域から
新しく確保した DRAM 領域へ置き換える



プログラムによる NIC の

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ書き込むことで配列で参照されるデータを適切に受け取ったことを NIC へ通知する

tail

head

DRAM 上の配列

[0]	address: 0x32000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

NIC のレジスタ

```
ring_head: 1
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

0x5000

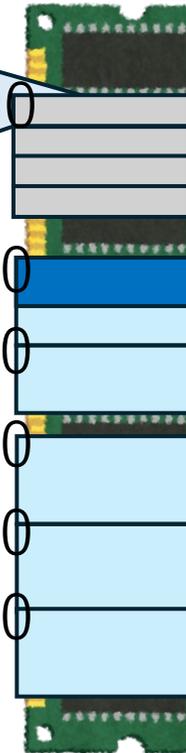
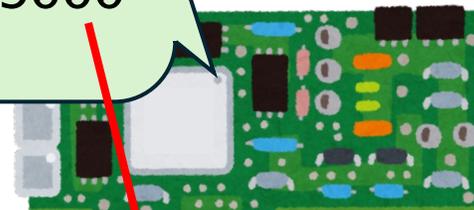
0x30000

0x32000

0x40000

0x42000

0x44000



プログラムによる NIC の

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ書き込むことで配列で参照されるデータを適切に受け取ったことを NIC へ通知する

NIC のレジスタ

```
ring_head: 1
ring_tail: 1
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x32000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

head
tail

0x5000

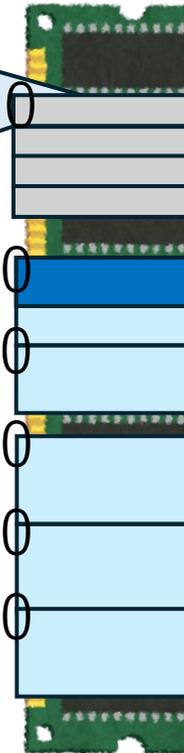
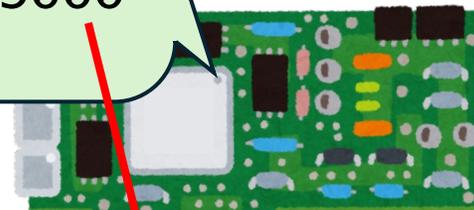
0x30000

0x32000

0x40000

0x42000

0x44000



このキュー

これで NIC は受信したパケットを配列の [0] が参照する DRAM 領域に書き込んで良いと判断できる

Q. 送受信キュー

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

どのように構成されるか？

レジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ書き込むことで配列で参照されるデータを適切に受け取ったことを NIC へ通知する

NIC のレジスタ

```
ring_head: 1
ring_tail: 1
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x32000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

0x5000

0x30000

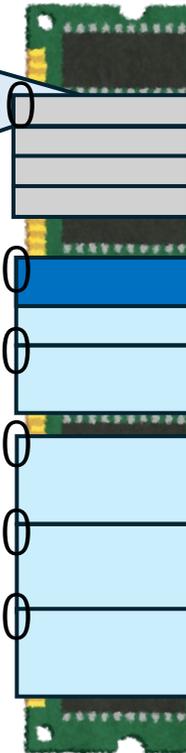
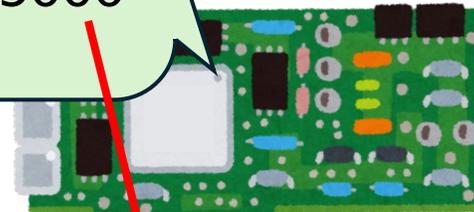
0x32000

0x40000

0x42000

0x44000

head
tail



このキュー

これで NIC は受信したパケットを配列の [0] が参照する DRAM 領域に書き込んで良いと判断できる

NIC のレジスタ
ring_head: 1
ring_tail: 1
ring_address: 0x5000
ring_size: 4

Q. 送受信キュー

どのように構成されるか?

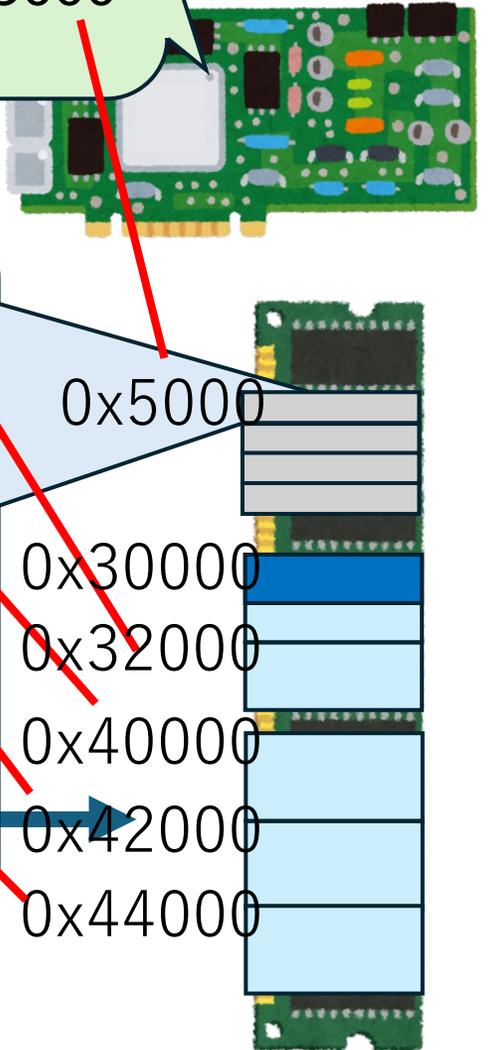
NIC により新しいデータが [0] が参照する物理アドレス 0x32000 に書き込まれても物理アドレス 0x30000 のデータは紐付けが解消されているため上書きされない

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ書き込むことで配列で参照されるデータを適切に受け取ったことを NIC へ通知する

DRAM 上の配列

[0]	address: 0x32000 length: 100
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length



プログラムによる NIC

このキューが受信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

プログラムは DRAM 上の受信データを処理する (例: TCP/IP スタックに渡す)

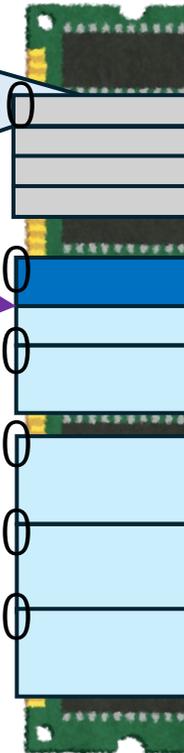
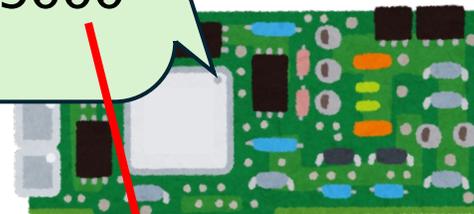
```
void *address;  
int length;  
} slot[NUM_SLOT];
```

NIC のレジスタ

```
ring_head: 1  
ring_tail: 1  
ring_address: 0x5000  
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x32000 length: 100	0x5000
[1]	address: 0x40000 length	0x30000
[2]	address: 0x42000 length	0x32000
		0x40000
[3]	Address: 0x44000 length	0x42000
		0x44000



カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

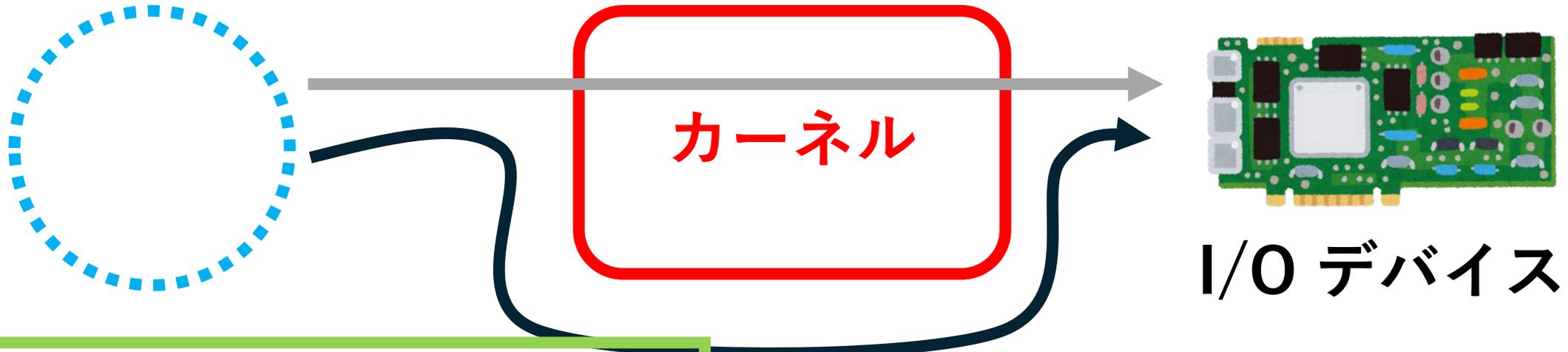
A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常では デバイス操作のメモリ領域へのアクセスが許可されていないから

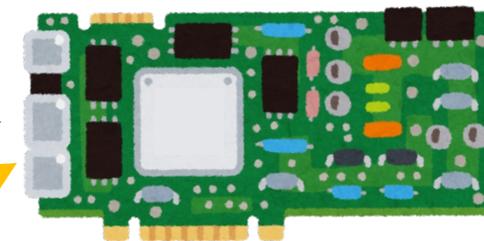
Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	

NIC へのアクセス

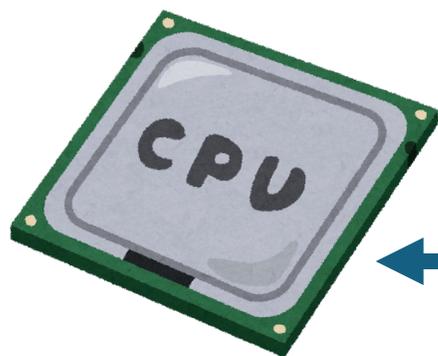


アドレスごまかすという
0xff000000 は
NIC
信号の出入りを振り分ける

物理アドレス
0xff000000 へ
アクセス

Address
Decoder

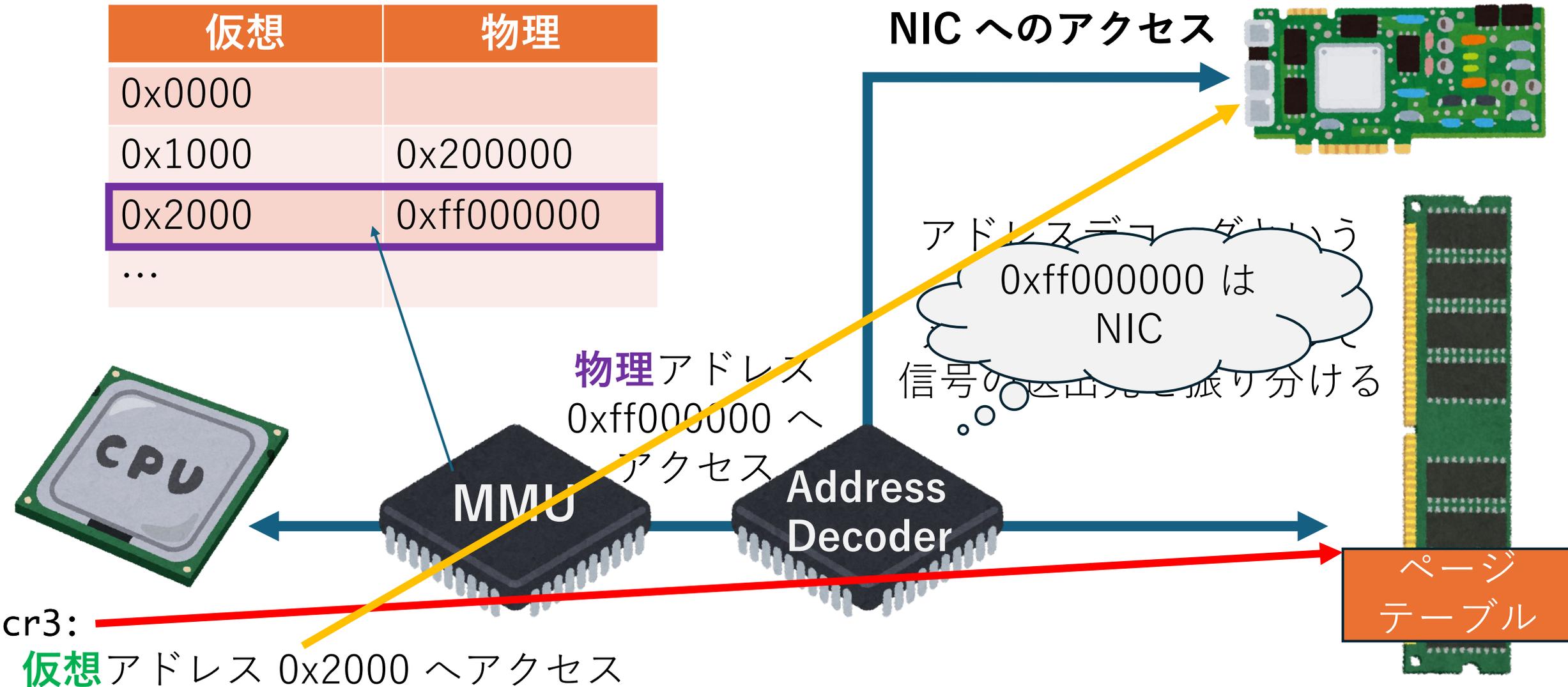
MMU



ページ
テーブル

cr3:

仮想アドレス 0x2000 へアクセス



メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	

基本的にカーネルは一般的なプロセスのページテーブルにデバイスアクセス用の物理アドレスを記載しない
もしくは、記載しても非特権モード時のアクセスを許可しない設定を適用



cr3: **仮想** アドレス 0x2000 へアクセス

ページ
テーブル

メモリの読み書きを通じたデバイス操作

仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	

基本的にカーネルは一般的なプロセスのページテーブルにデバイスアクセス用の物理アドレスを記載しない
もしくは、記載しても非特権モード時のアクセスを許可しない設定を適用

カーネル実行時に参照されるページテーブルにのみデバイスアクセス用の物理アドレスを設定する

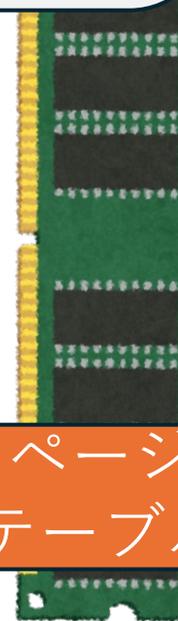
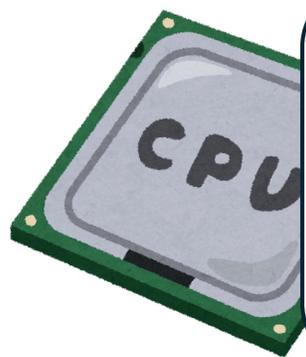
0xff000000はNIC
信号の出入りを振り分ける

Address Decoder

ページ
テーブル

cr3:

仮想アドレス 0x2000 へアクセス



メモリの読み書きを通じたデバイス操作

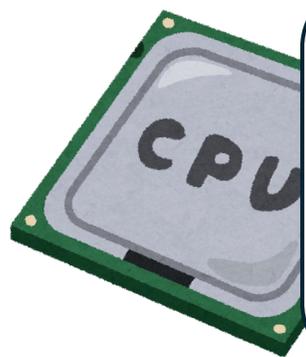
仮想	物理
0x0000	
0x1000	0x200000
0x2000	0xff000000
...	

基本的にカーネルは一般的なプロセスのページテーブルにデバイスアクセス用の物理アドレスを記載しない
もしくは、記載しても非特権モード時のアクセスを許可しない設定を適用

カーネル実行時に参照されるページテーブルにのみデバイスアクセス用の物理アドレスを設定する

一般的なプロセスはシステムコールを通じてカーネルにリクエストすることで代わりにデバイスにアクセスしてもらう

cr3: **仮想** アドレス 0x2000 へアクセス



Addr Dec

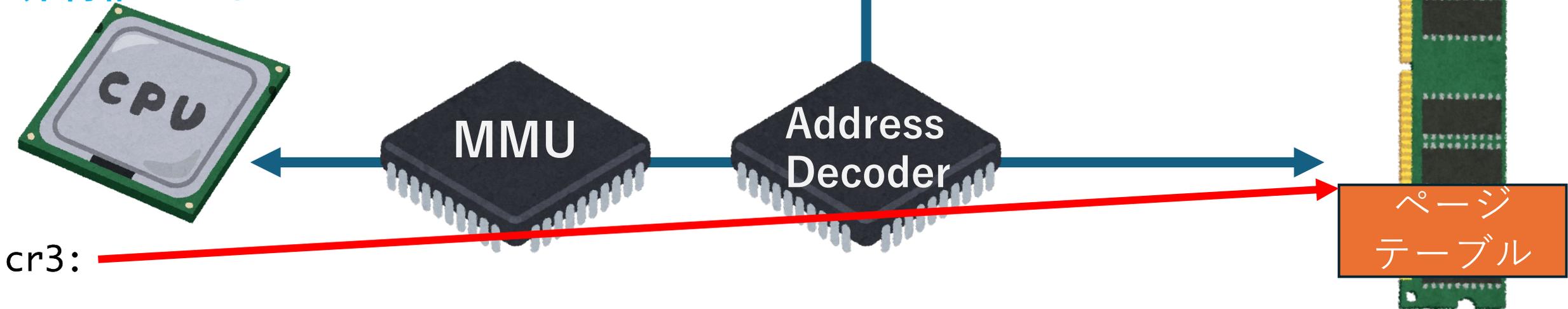
0xff000000
NIC

シ
ブル

メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

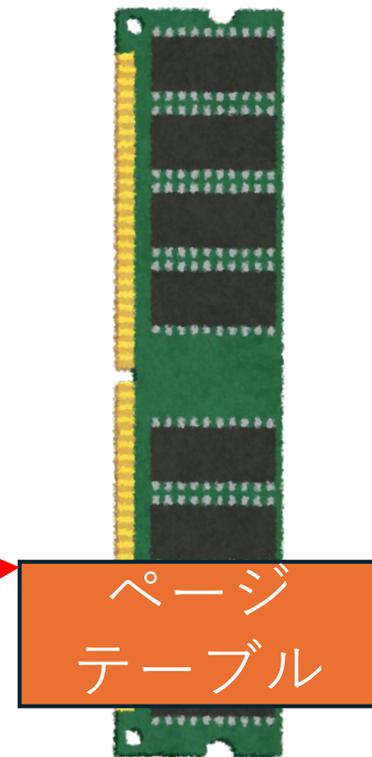
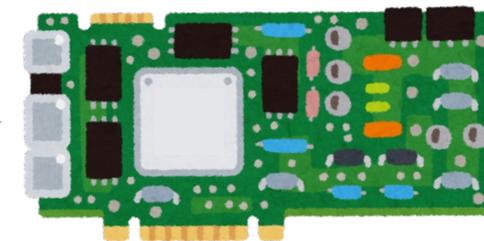
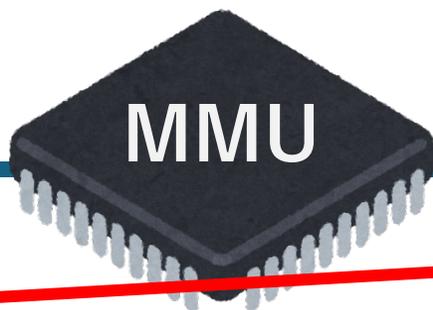
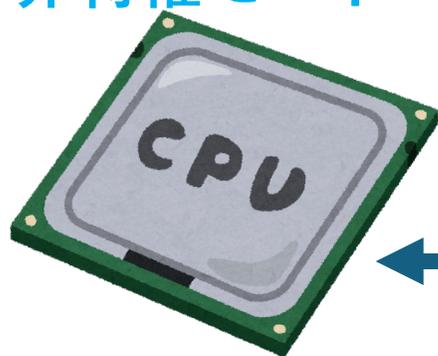
非特権モード



メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

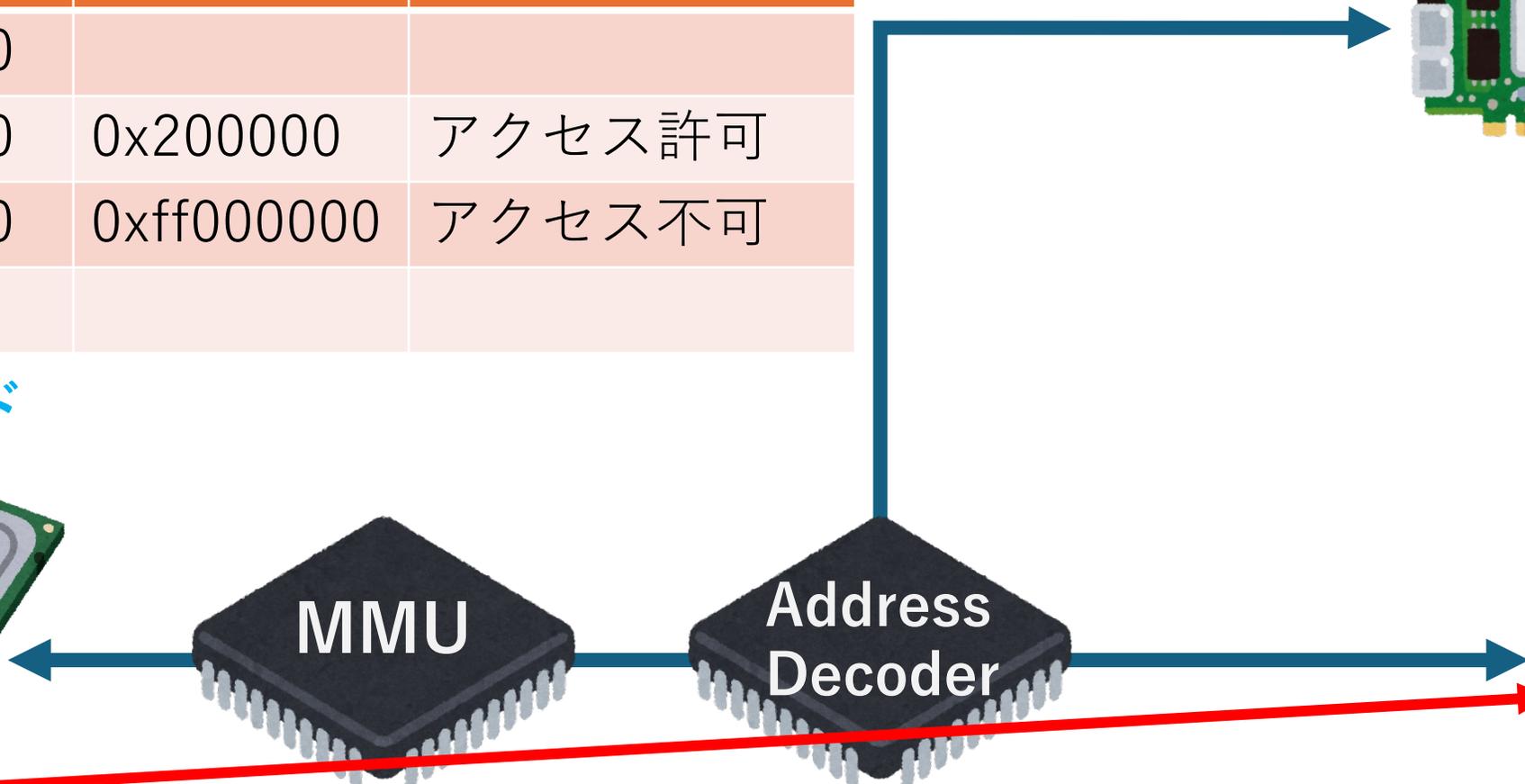
非特権モード



ページ
テーブル

cr3:

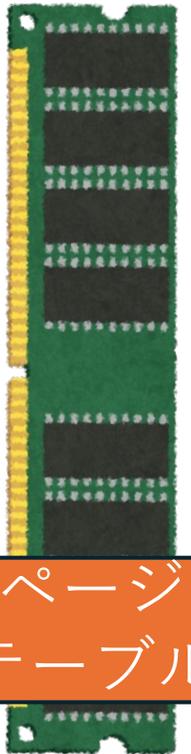
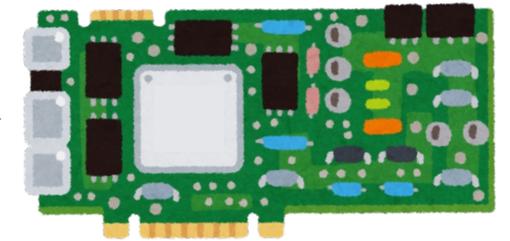
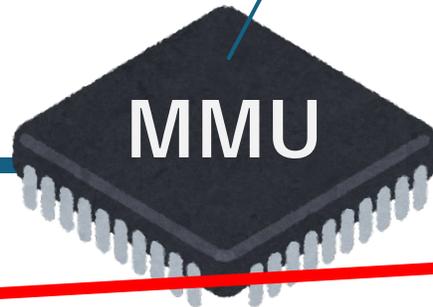
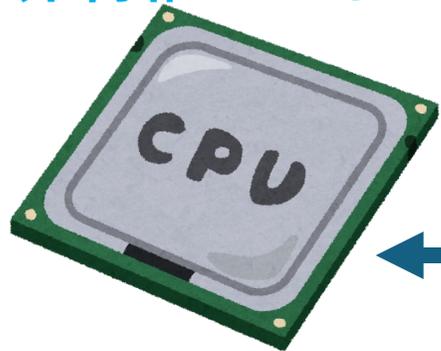
仮想アドレス 0x2000 へアクセス



メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

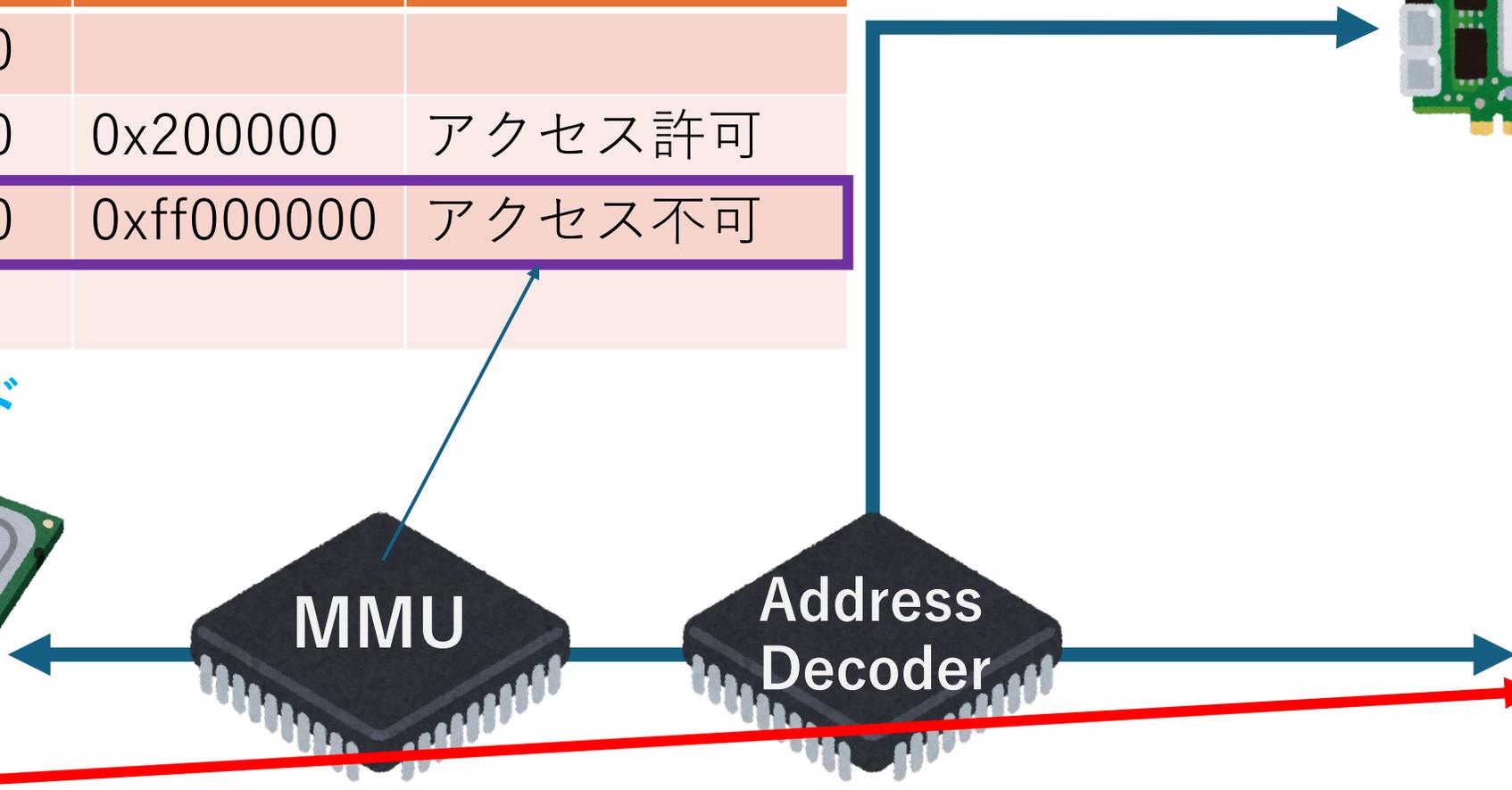
非特権モード



ページ
テーブル

cr3:

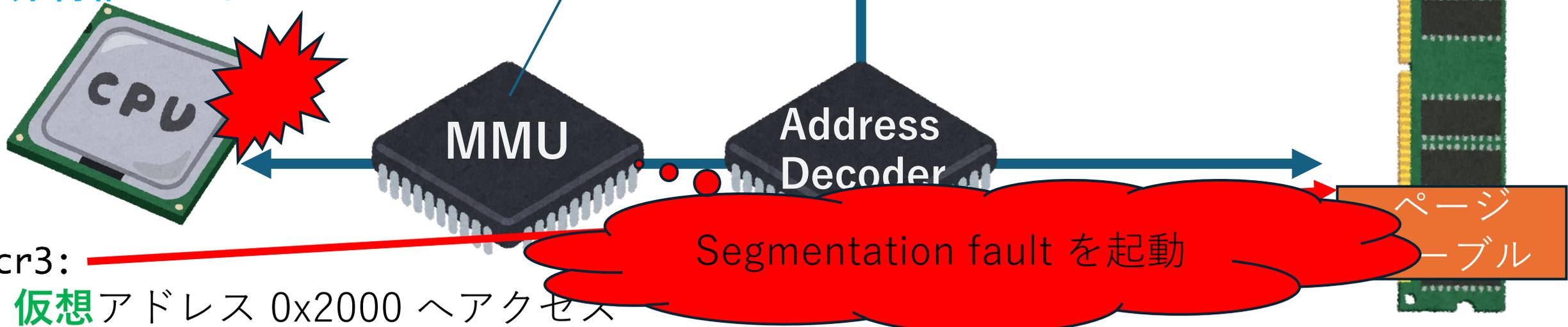
仮想アドレス 0x2000 へアクセス



メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

非特権モード



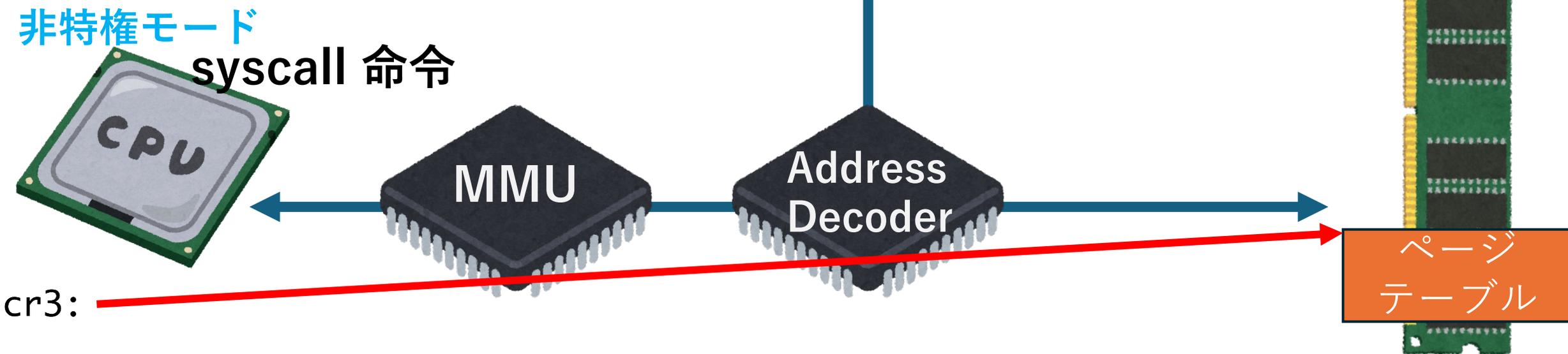
cr3: 仮想アドレス 0x2000 へアクセス

Segmentation fault を起動

ページ
テーブル

メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

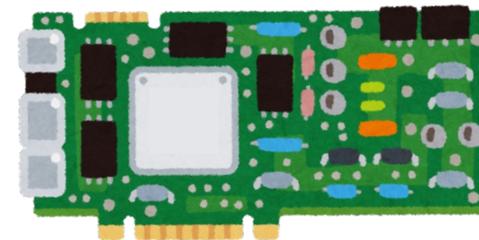
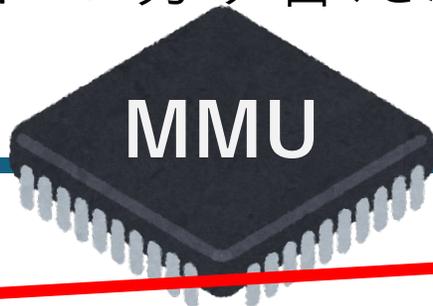
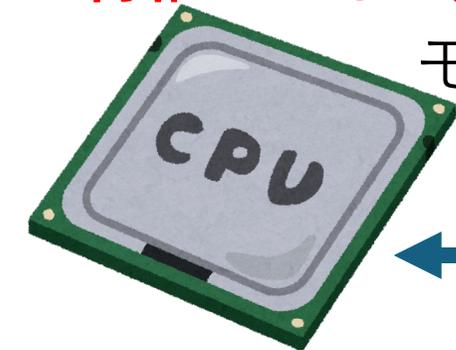


メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

特権モード (カーネル)

モードの切り替えが完了



ページ
テーブル

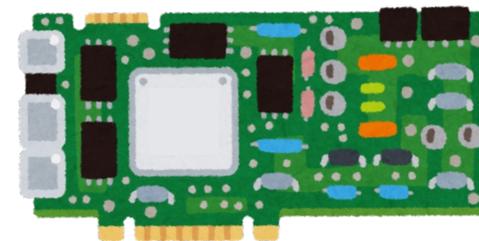
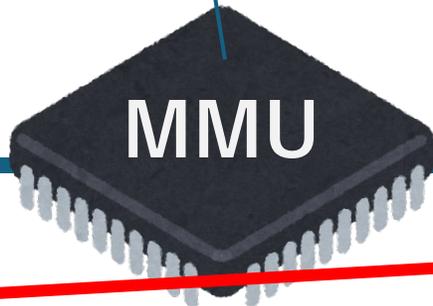
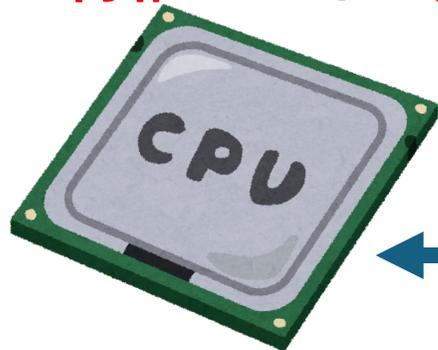
cr3:



メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

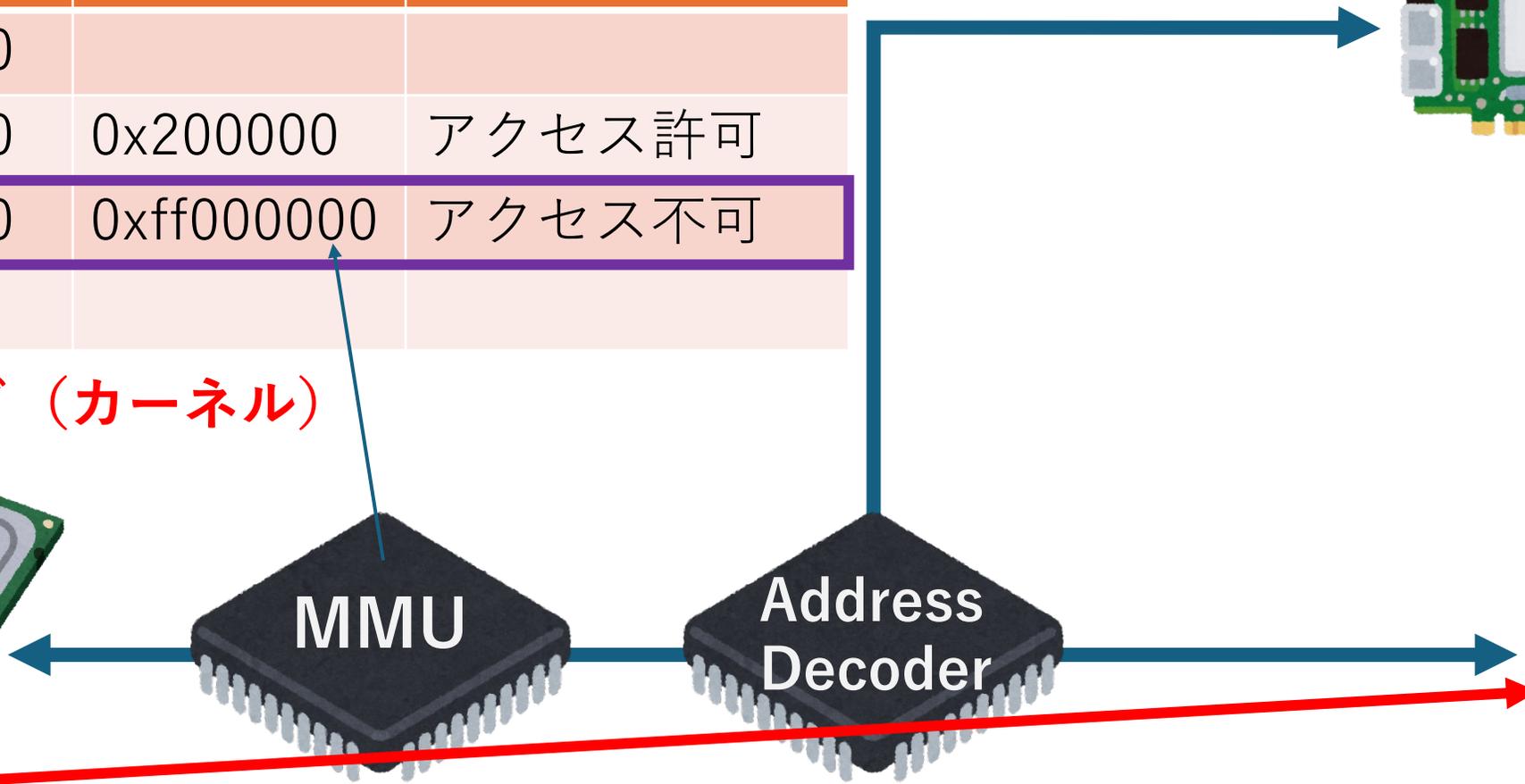
特権モード (カーネル)



ページ
テーブル

cr3:

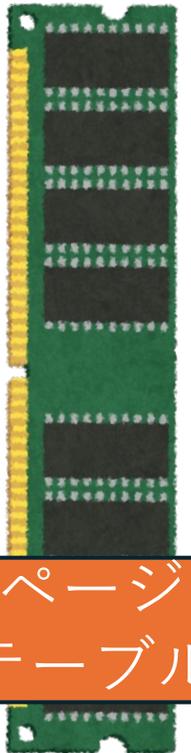
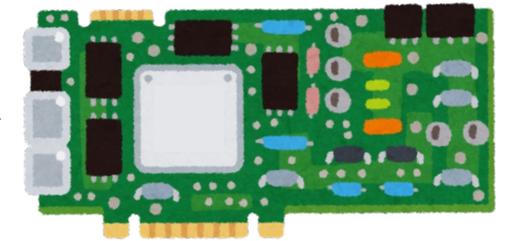
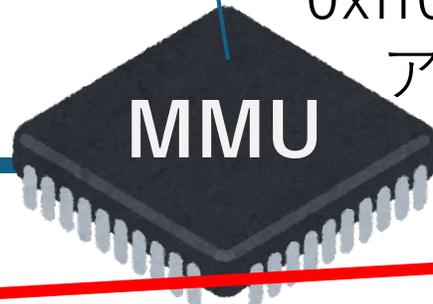
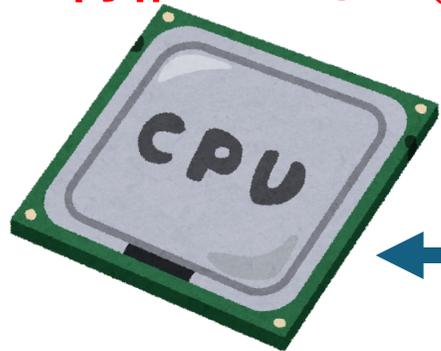
仮想アドレス 0x2000 へアクセス



メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

特権モード (カーネル)



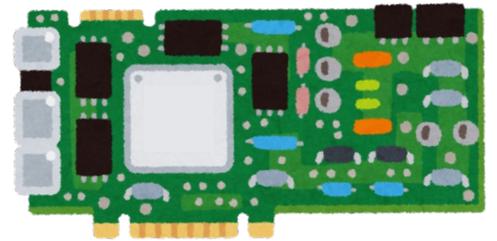
ページ
テーブル

物理アドレス
0xff000000 へ
アクセス

cr3: **仮想** アドレス 0x2000 へアクセス

メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

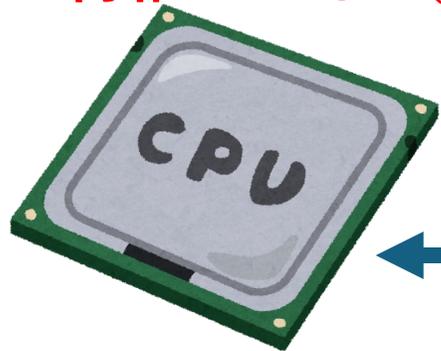


特権モードが適用されているのでアクセス可能

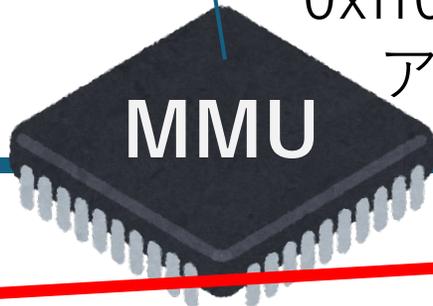


ページ
テーブル

特権モード (カーネル)



物理アドレス
0xff000000 へ
アクセス

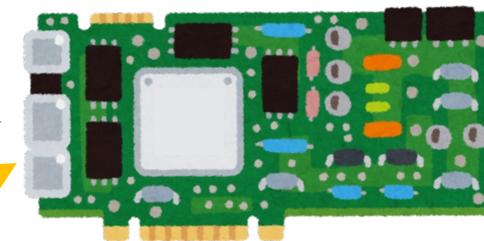


cr3: **仮想** アドレス 0x2000 へアクセス

メモリの読み書きを通じたデバイス操作

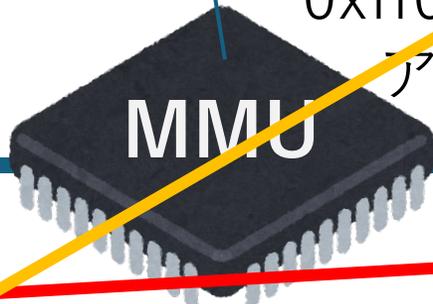
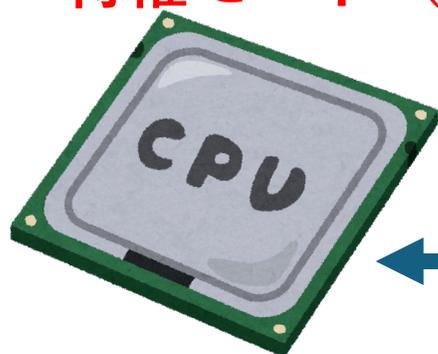
仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

NIC へのアクセス



0xff000000 は NIC

特権モード (カーネル)

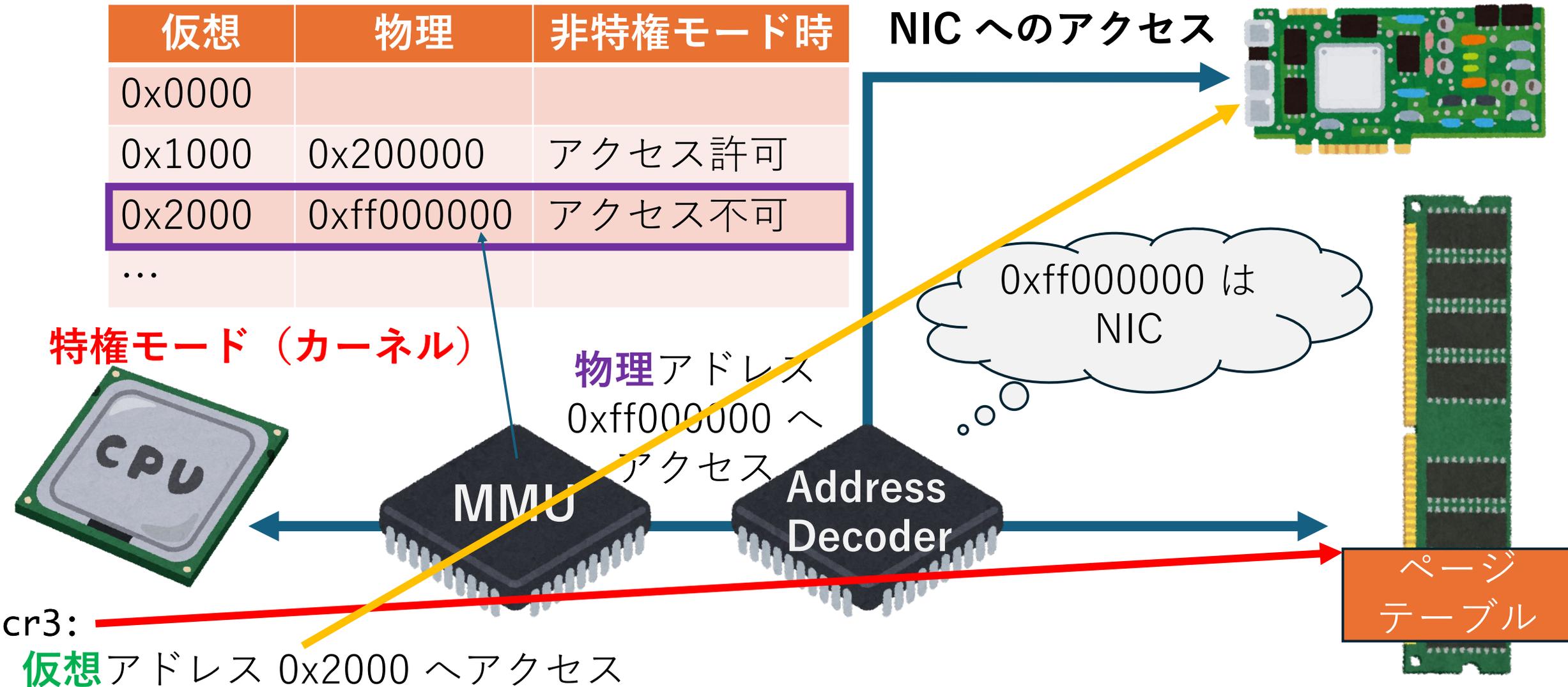


物理アドレス
0xff000000 へ
アクセス

cr3:

仮想アドレス 0x2000 へアクセス

ページ
テーブル



カーネルバイパスとは？

Q. これは何？

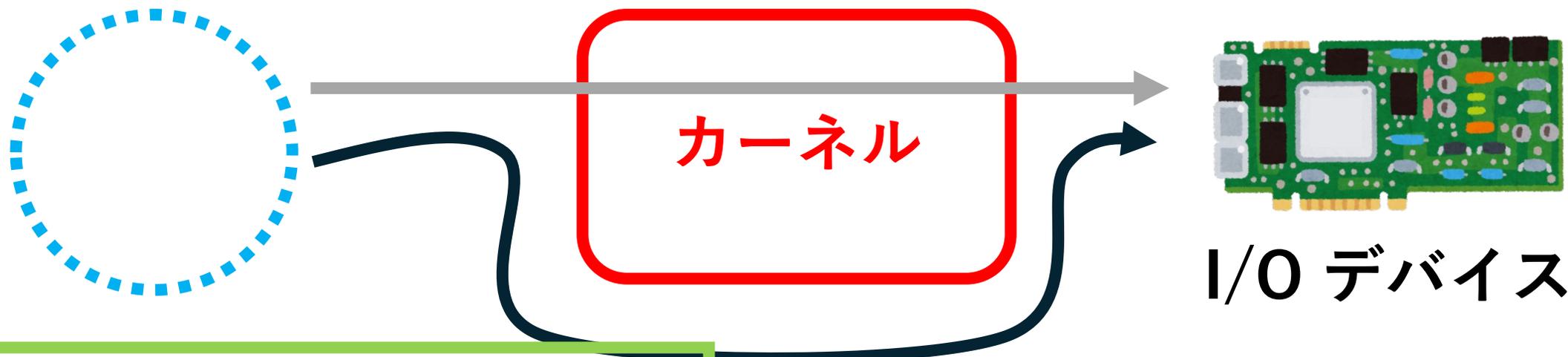
A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常では デバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作メモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

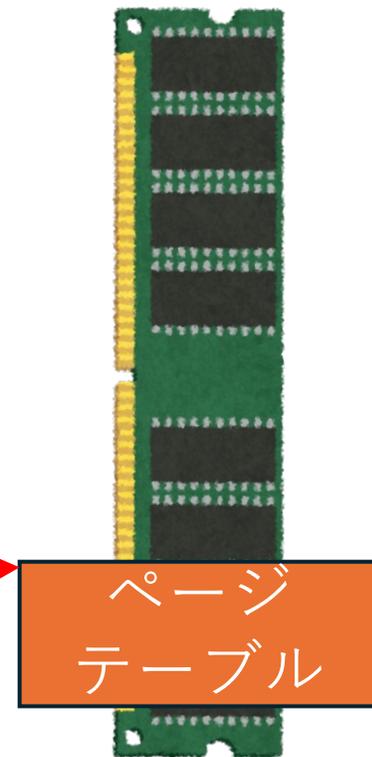
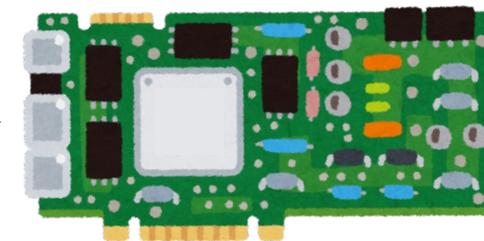
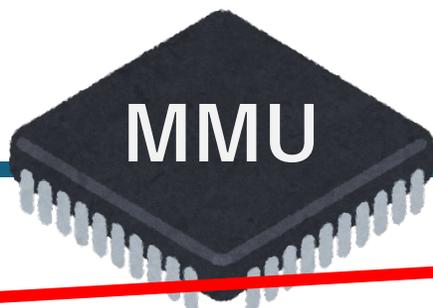
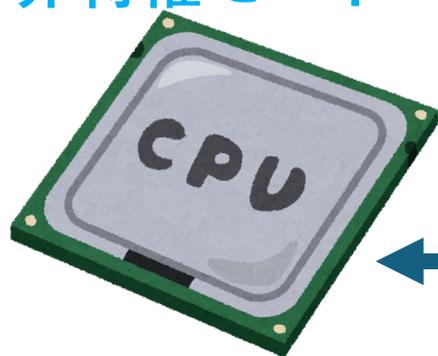
Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス不可
...		

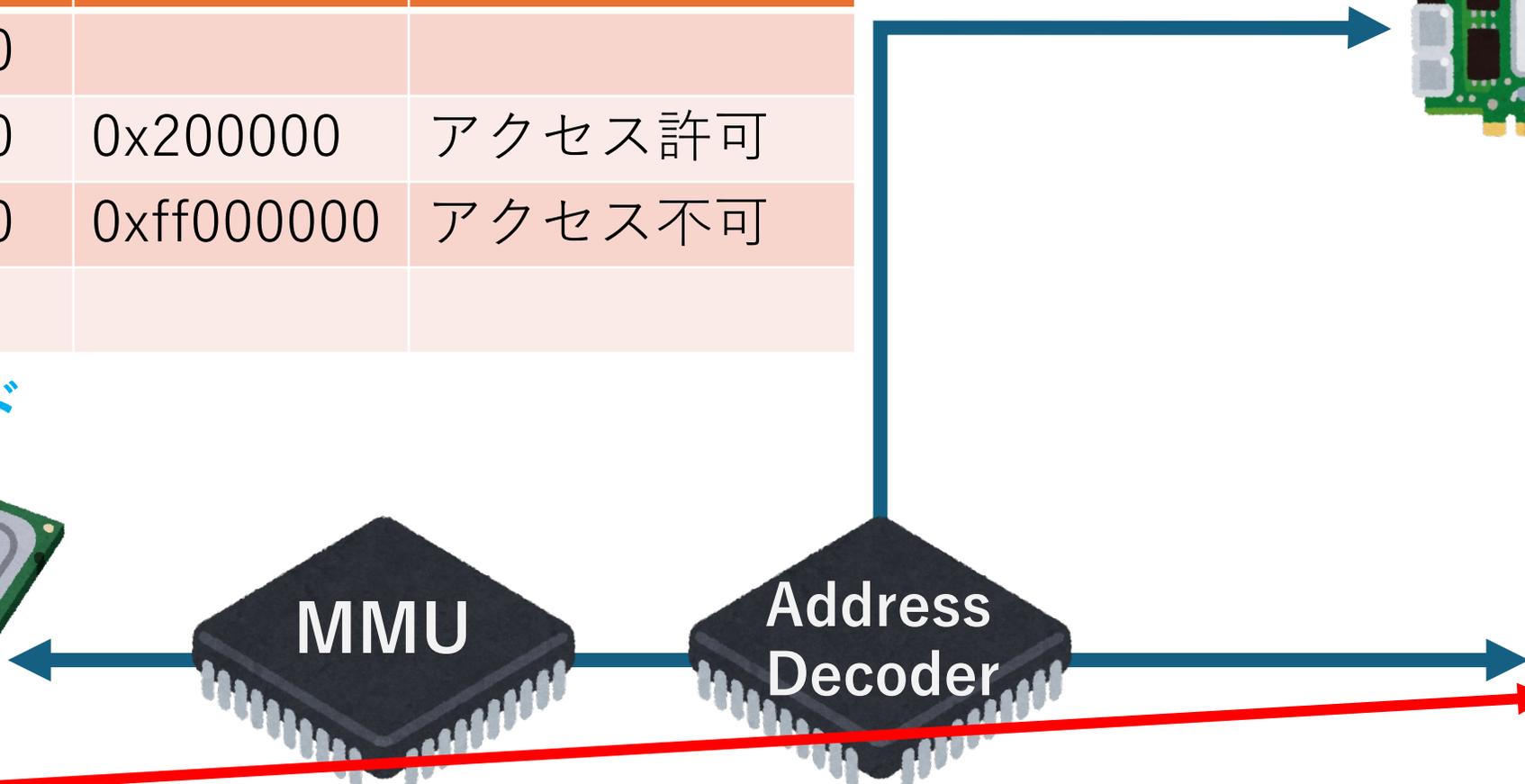
非特権モード



ページ
テーブル

cr3:

仮想アドレス 0x2000 へアクセス

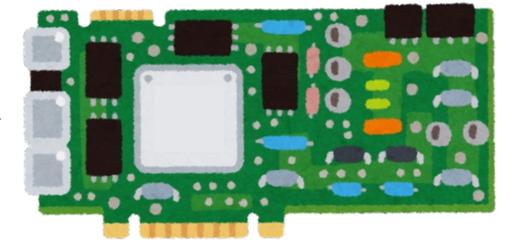
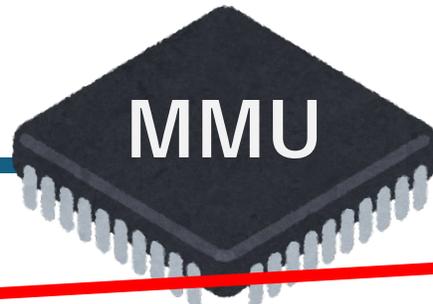
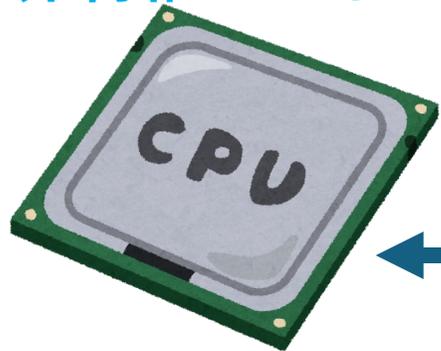


メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス 許可
...		

ページテーブル設定を通して
デバイス操作用メモリ領域へ
非特権モードで
アクセス可能にする

非特権モード



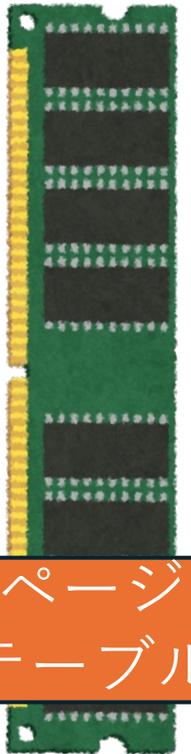
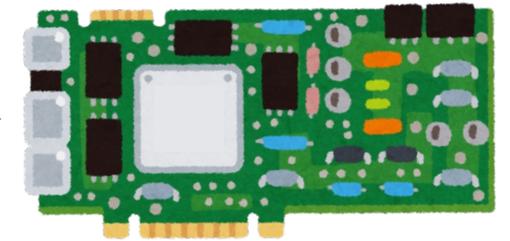
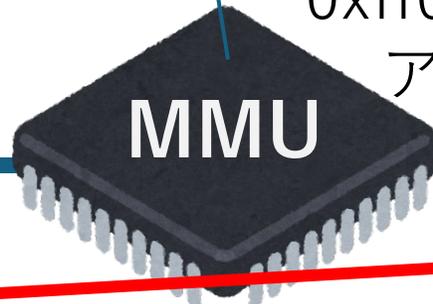
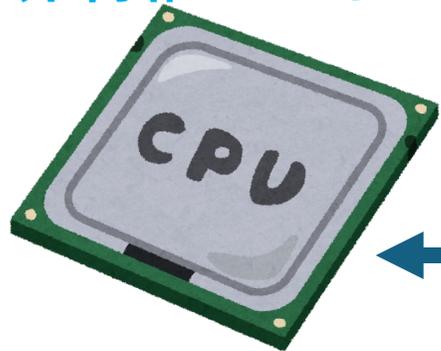
ページ
テーブル

cr3: **仮想** アドレス 0x2000 へアクセス

メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス許可
...		

非特権モード



ページ
テーブル

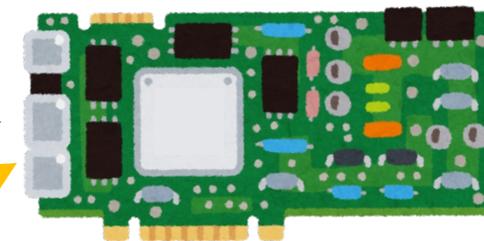
物理アドレス
0xff000000 へ
アクセス

cr3: **仮想**アドレス 0x2000 へアクセス

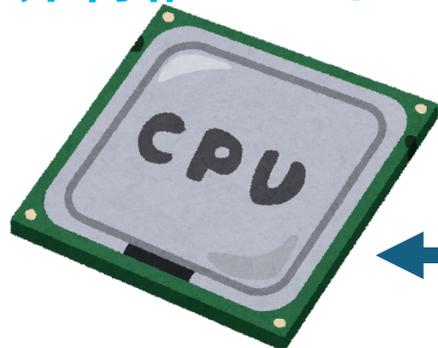
メモリの読み書きを通じたデバイス操作

仮想	物理	非特権モード時
0x0000		
0x1000	0x200000	アクセス許可
0x2000	0xff000000	アクセス許可
...		

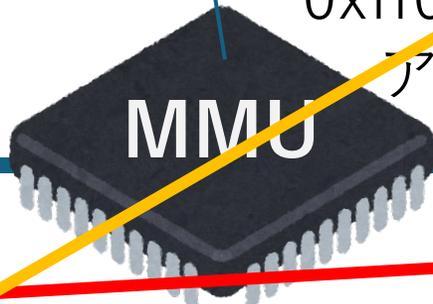
NIC へのアクセス



非特権モード



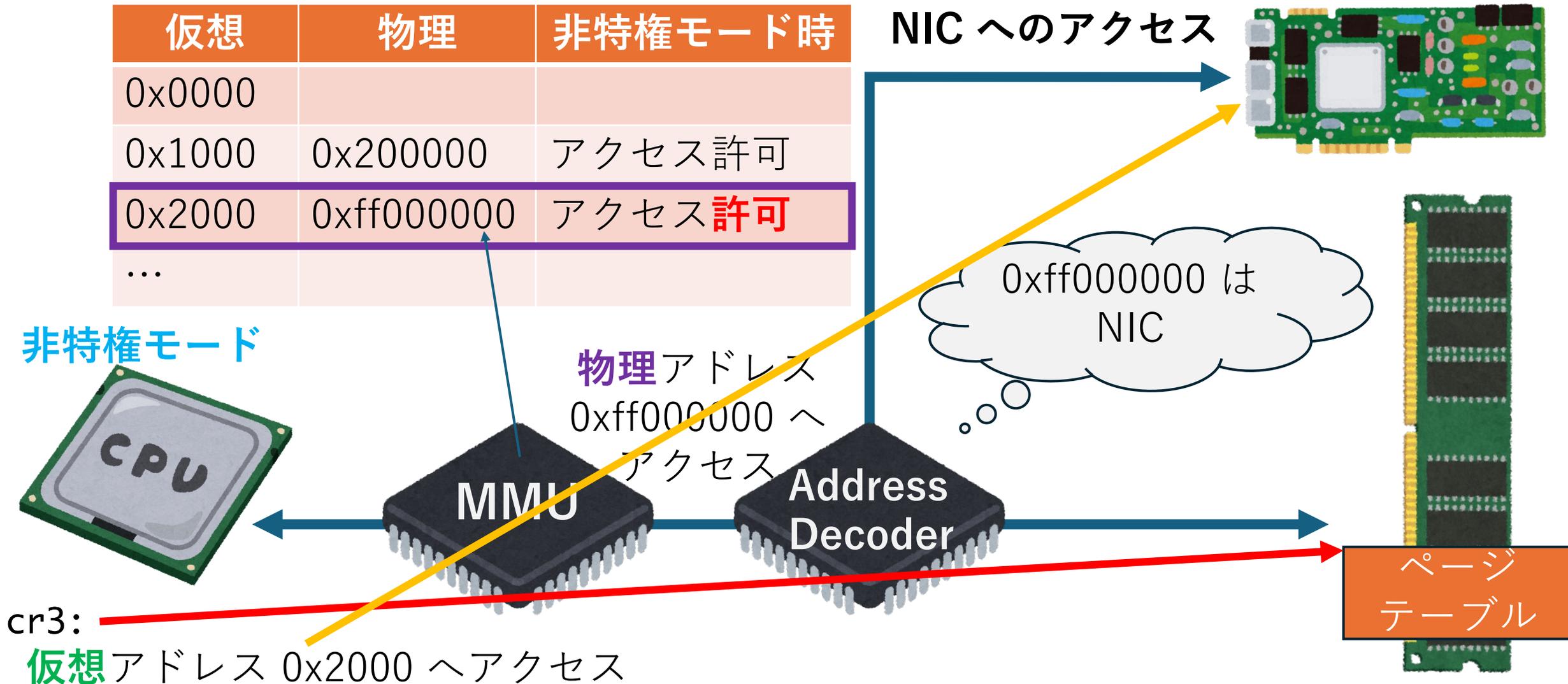
物理アドレス
0xff000000 へ
アクセス



cr3:

仮想アドレス 0x2000 へアクセス

ページ
テーブル



カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. デバイスアクセスのための処理と I/O デバイスの I/O 対象データ

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

カーネルバイパスとは？

Q. これは何？

A. ユーザー空間で動作するプログラム

Q. 何がカーネルを中継しなくなる？

A. 実装次第

Q. デバイスへのアクセスとは具体的に何？

A. 基本的にはメモリの読み書き



I/O デバイス

Q. 何故通常はカーネルを経由する？

A. ユーザー空間プログラムは通常ではデバイス操作のメモリ領域へのアクセスが許可されていないから

Q. どうやってバイパスする？

A. ユーザー空間プログラムへデバイス操作のメモリ領域へのアクセスを許可する

プログラムによる NIC

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

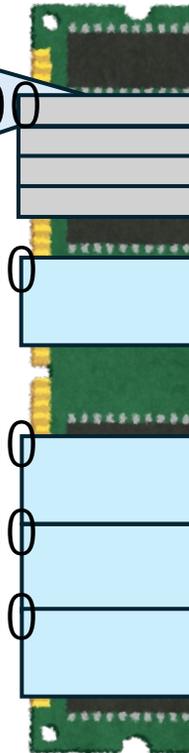
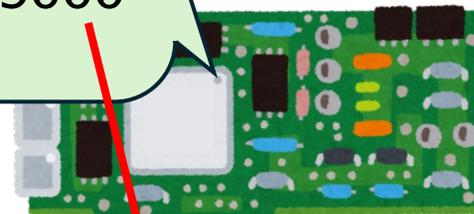
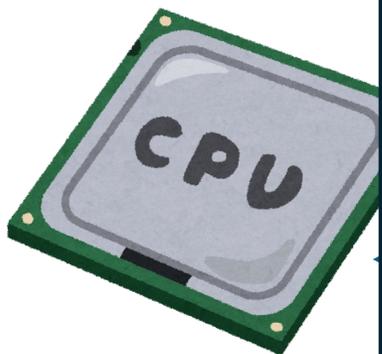
```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head
tail

DRAM 上の配列

[0]	address: 0x30000 length	0x5000
[1]	address: 0x40000 length	0x30000
[2]	address: 0x42000 length	0x40000
[3]	Address: 0x44000 length	0x42000 0x44000

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4



プログラムによる NIC

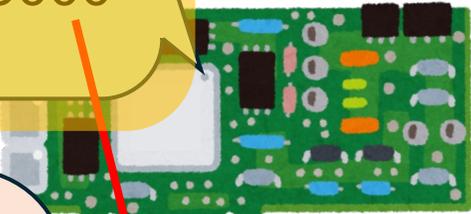
Q. 送受信キューはどのように構成されるか？

A. **メモリアクセス設定 2 パターン**

パターン 1 : NIC のレジスタへのアクセスを許可する

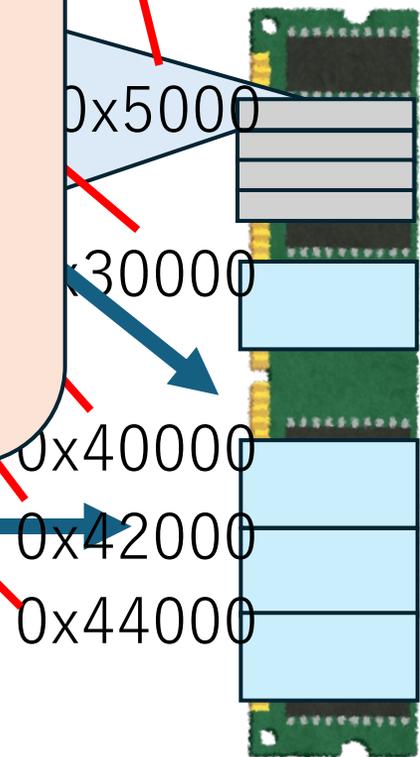
NIC のレジスタ

- ring_head: 0
- ring_tail: 0
- ring_address: 0x5000
- ring_size: 4



```
struct slot {  
    void *address;  
    int length;  
} slot[NUM_SLOT];
```

[3] Address: 0x44000
length



プログラムによる NIC

Q. 送受信キューはどのように構成されるか？

A. **メモリアクセス設定 2 パターン**

パターン 1 : NIC のレジスタへのアクセスを許可する

ユーザー空間のプログラムでもNICのレジスタへアクセスできればリングバッファ用の配列と送受信に利用する DRAM 領域を mmap のような一般的なメモリ確保機能を使って確保してNICと紐づけることができます (NIC と紐づける DRAM 領域の物理アドレスを取得できる必要があります)

```
NIC のレジスタ  
ring_head: 0  
ring_tail: 0  
ring_address: 0x5000  
ring_size: 4
```

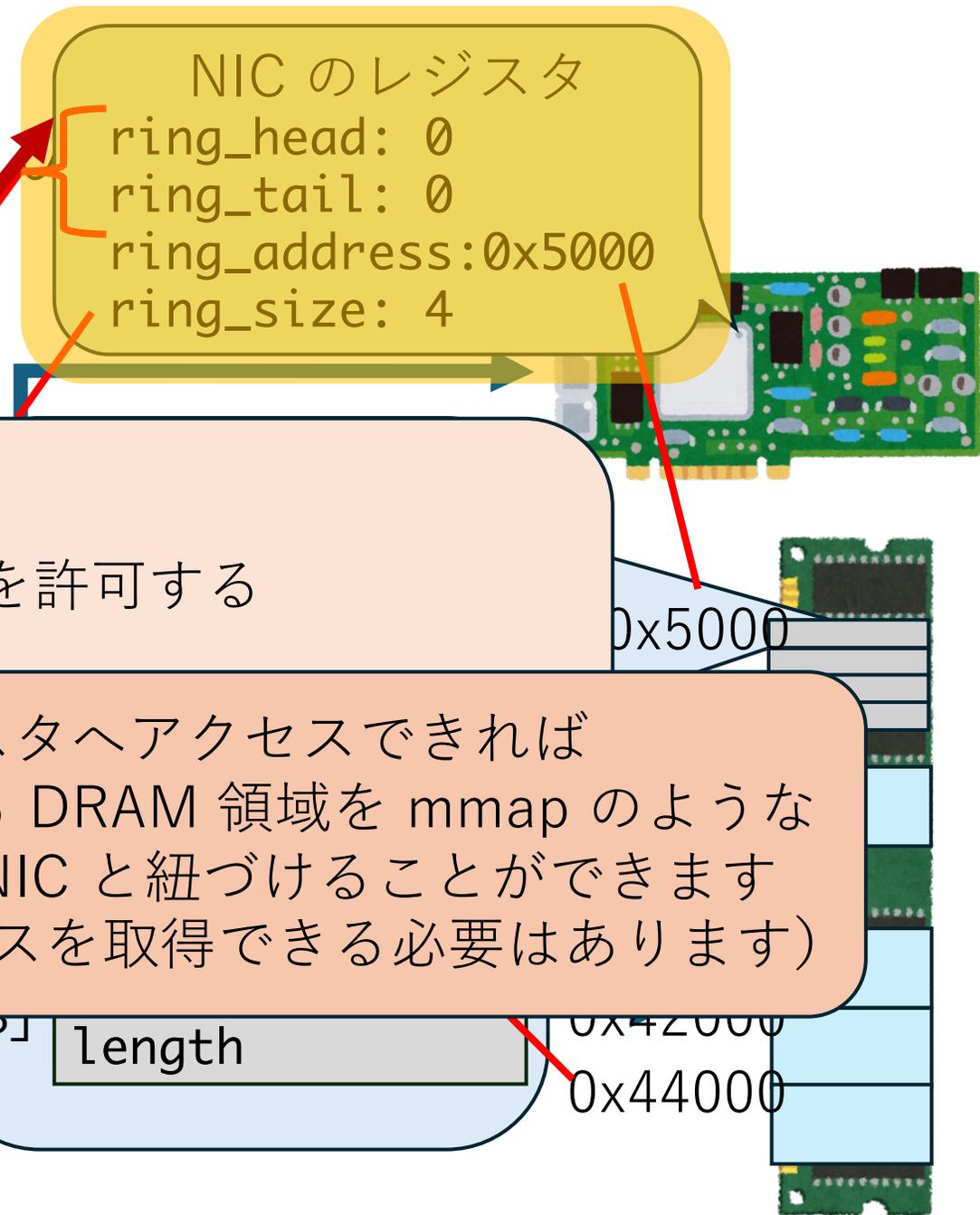
```
int length;  
} slot[ NUM_SLOT ];
```

```
length
```

0x5000

0x42000

0x44000



プログラムによる NIC

Q. 送受信キューはどのように構成されるか？

A

メモリアクセス設定 2 パターン

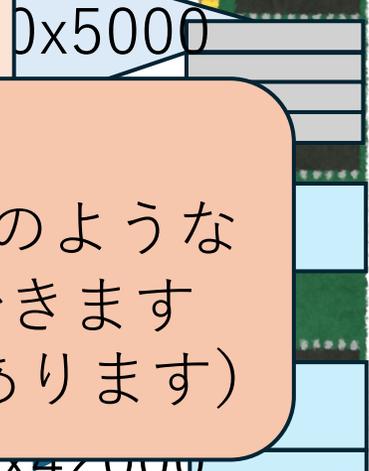
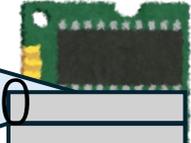
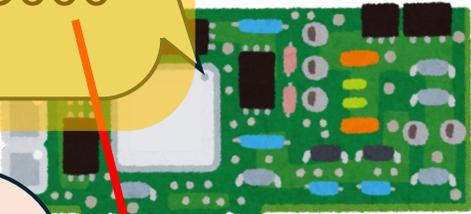
パターン 1 : NIC のレジスタへのアクセスを許可する

ユーザー空間のプログラムでもNICのレジスタへアクセスできればリングバッファ用の配列と送受信に利用する DRAM 領域を mmap のような一般的なメモリ確保機能を使って確保してNICと紐づけることができます (NIC と紐づける DRAM 領域の物理アドレスを取得できる必要があります)

この場合はユーザー空間で各 NIC のハードウェア仕様に対応した操作が必要となるため、ユーザー空間でデバイスドライバを動作させる必要があります

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```



プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A

メモリアクセス設定 2 パターン

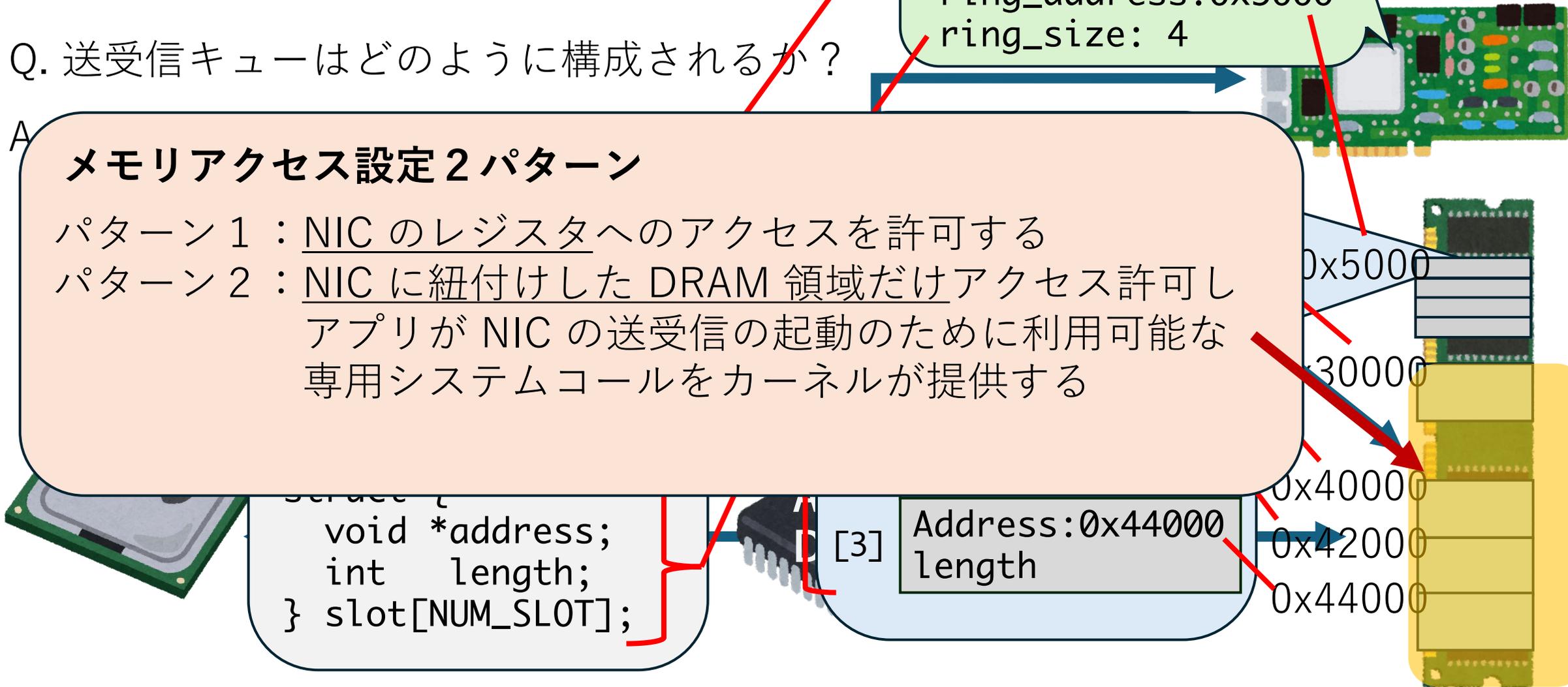
パターン 1 : NIC のレジスタへのアクセスを許可する

パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

```
struct {  
    void *address;  
    int length;  
} slot[ NUM_SLOT ];
```

[3] Address: 0x44000
length

0x5000
0x30000
0x40000
0x42000
0x44000



プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

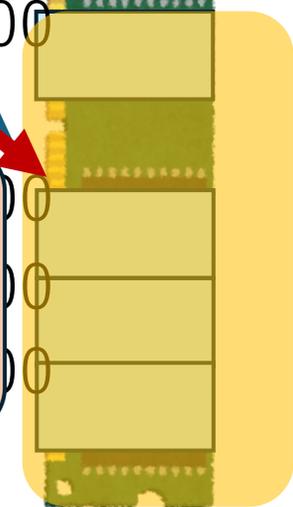
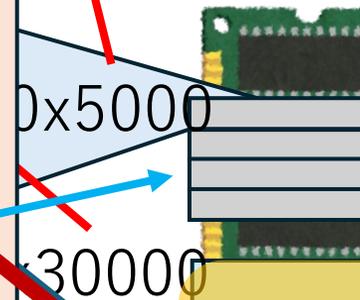
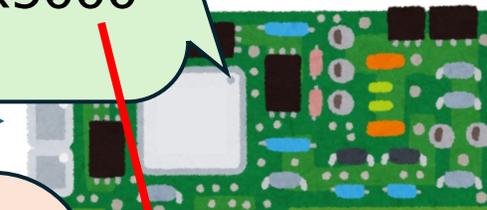
A

メモリアクセス設定 2 パターン

パターン 1 : NIC のレジスタへのアクセスを許可する

パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

NIC のレジスタとリングバッファ用配列はカーネルからのみ
アクセス可能なように設定されており、アプリがシステムコールを
通してリクエストした時に、それらへの読み書きが行われます



プログラムによる NIC

```
NIC のレジスタ  
ring_head: 0  
ring_tail: 0  
ring_address: 0x5000  
ring_size: 4
```

Q. 送受信キューはどのように構成されるか？

A. メモリアクセス設定 2 パターン

パターン 1: NIC のレジスタへのアクセスを許可する

この場合は、NIC ごとのハードウェア仕様の差異はカーネル内で対応します：
カーネルに実装された各 NIC 用のドライバが利用されるため
ユーザー空間でデバイスドライバを実行する必要はありません

NIC のレジスタとリングバッファ用配列はカーネルからのみ
アクセス可能なように設定されており、アプリがシステムコールを
通してリクエストした時に、それらへの読み書きが行われます

プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A

メモリアクセス設定 2 パターン

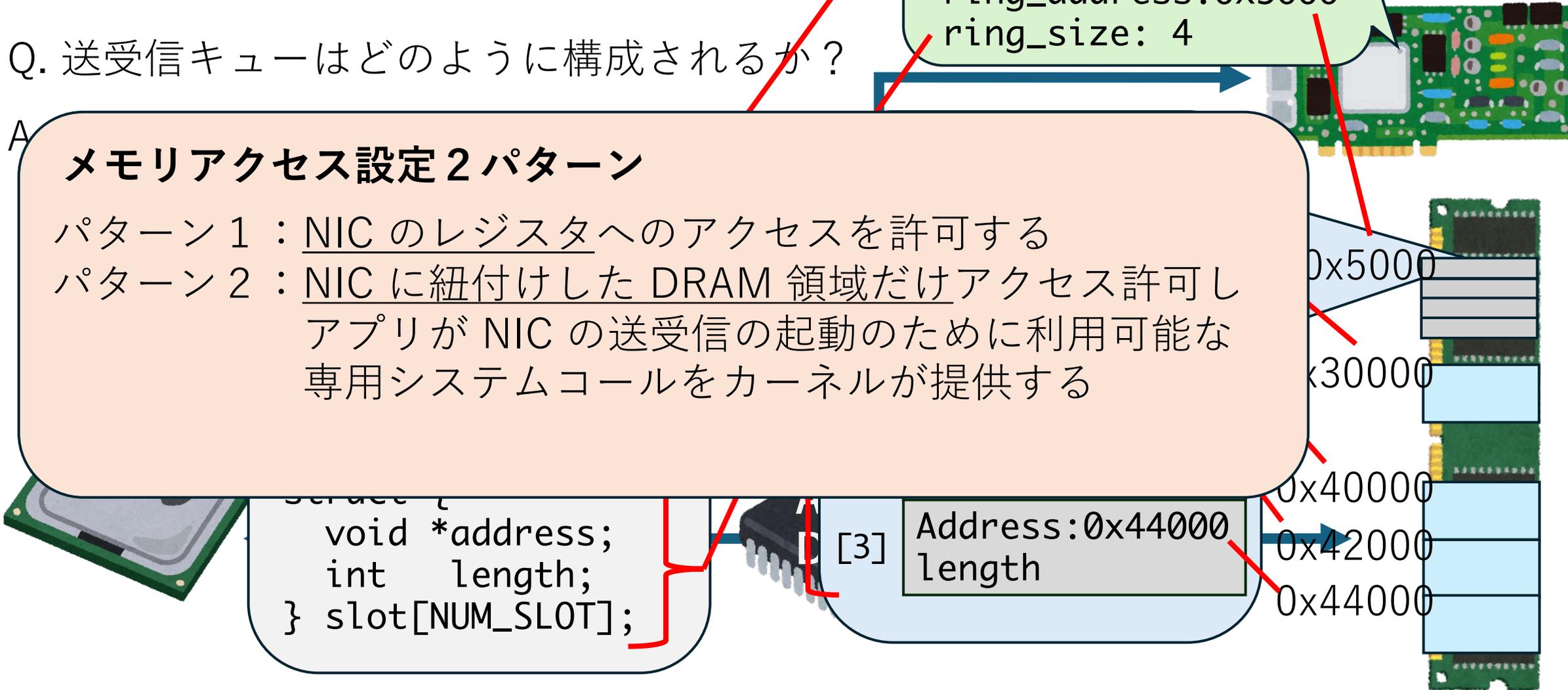
パターン 1 : NIC のレジスタへのアクセスを許可する

パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

```
struct {  
    void *address;  
    int length;  
} slot[ NUM_SLOT ];
```

[3] Address: 0x44000
length

0x5000
k30000
0x40000
0x42000
0x44000



プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A

メモリアクセス設定 2 パターン

パターン 1 : NIC のレジスタへのアクセスを許可する

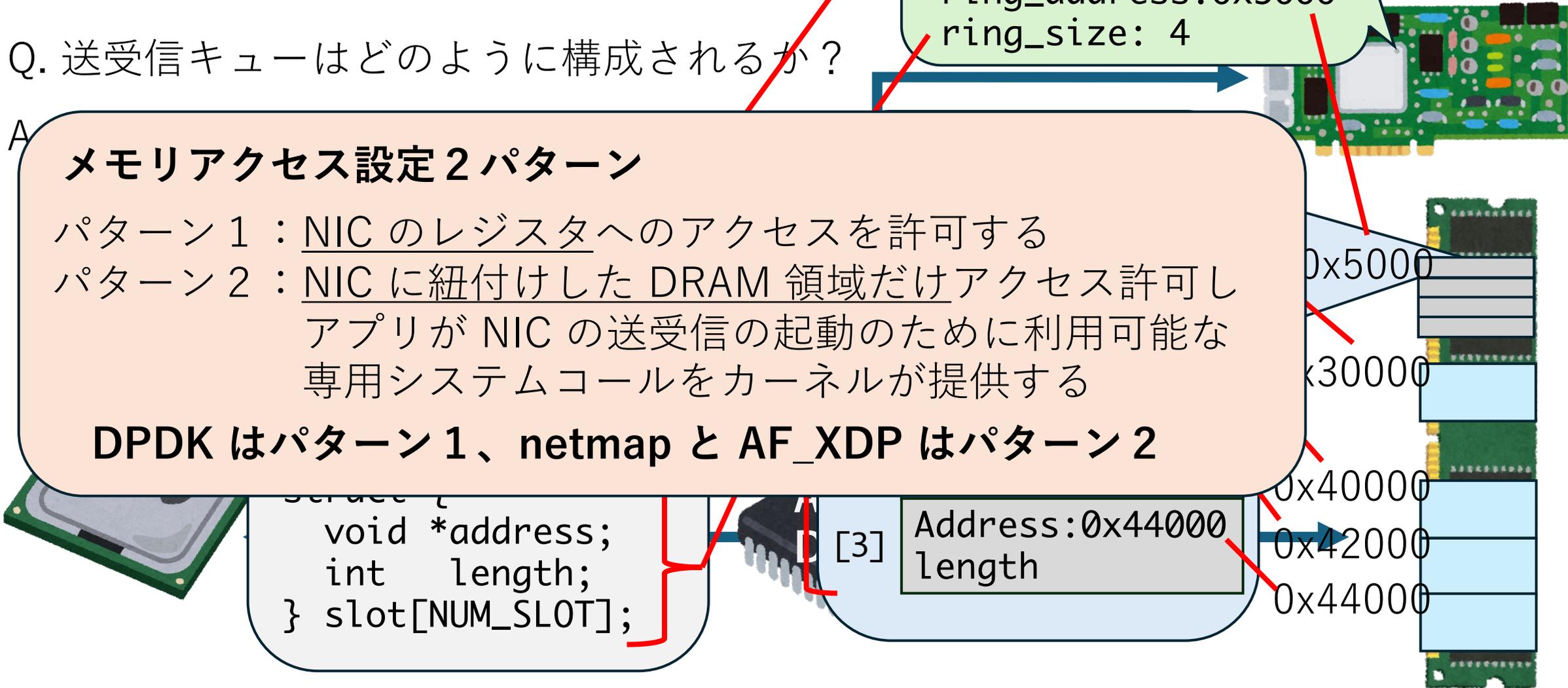
パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

DPDK はパターン 1、netmap と AF_XDP はパターン 2

```
struct {  
    void *address;  
    int length;  
} slot[ NUM_SLOT ];
```

[3] Address: 0x44000
length

0x5000
k30000
0x40000
0x42000
0x44000



簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

プログラムによるメモリアクセス

例：メモリアドレス 0x100000000 に 0x12345678 を書き込むプログラム

```
int main(void)
{
  *((int *) 0x100000000) = 0x12345678;
}
```

gcc -O0 program.c

コンパイル

a.out

ディスアセンブル

objdump -d ./a.out

00000000000001129 <main>:

rax レジスタに 0x100000000 を設定する

```
1131: 48 b8 00 00 00 00 01
1138: 00 00 00
113b: c7 00 78 56 34 12
```

rax レジスタの値により参照される
メモリアドレス (0x100000000) へ
0x12345678 を書き込む

```
endbr64
push   %rbp
mov    %rsp,%rbp
movabs $0x100000000,%rax
movl   $0x12345678,(%rax)
mov    $0x0,%eax
pop    %rbp
retq
nopl   0x0(%rax,%rax,1)
```

簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

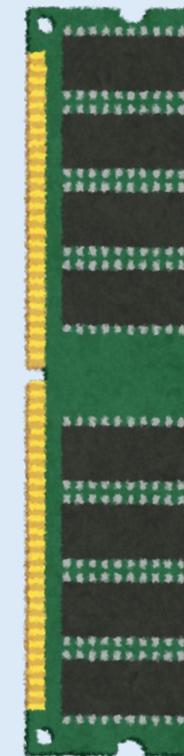
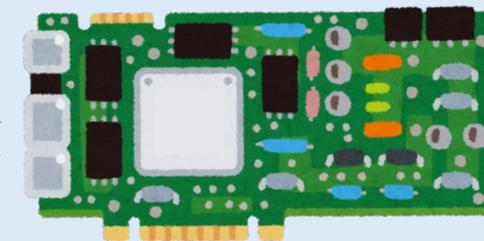
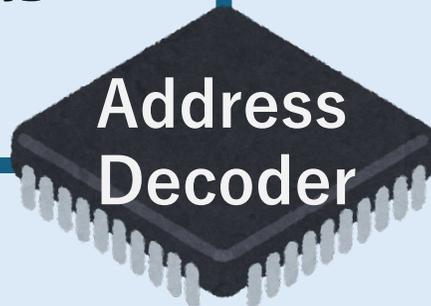
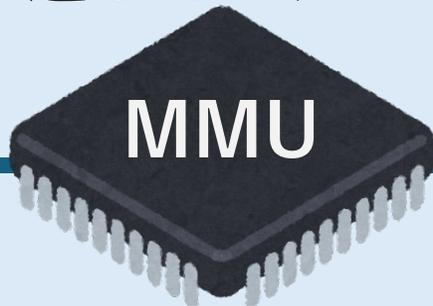
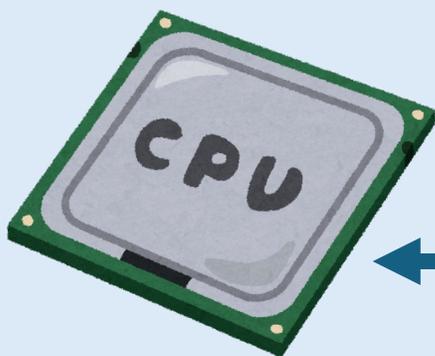
プログラムによる NIC の取り扱い

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラムから
NIC のレジスタと DRAM は
MMU / Address Decoder を
通じてアクセス可能

アドレスデコーダという
ハードウェアが
メモリアドレスに応じて
信号の送出先を振り分ける



(どのアドレスが NIC のどのレジスタに対応するかは NIC の仕様に依存)

簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

プログラムによるメモリアクセス

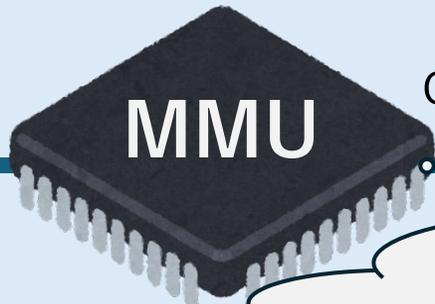
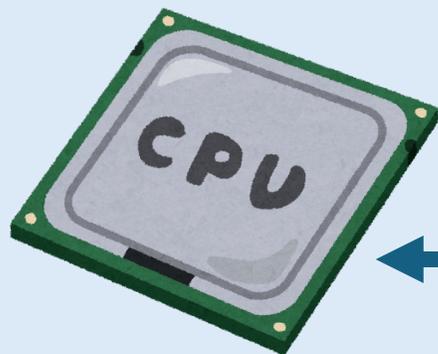
対応する**物理アドレス**として
0x2000 が設定されていた場合

仮想アドレス 0x100000000 へ
0x12345678 を書き込み

```
movl $0x12345678, (%rax)
```

仮想	物理
0x0000	
0x1000	
0x100000000	0x2000
...	

物理アドレス 0x2000 へ
0x12345678 を書き込み



MMU はアドレス変換に際して
cr3 で示されるページテーブルを参照

ページ
テーブル

rax:0x100000000
cr3: ページテーブルの物理アドレス

仮想アドレス 0x100000000 は
物理アドレスの 0x2000

簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

簡単にまとめると

- CPU 命令にはメモリの読み書き機能を提供するものがある
- DRAM と NIC は CPU のメモリ読み書き命令を通して操作可能
- カーネルはユーザー空間プログラムがアクセス可能な物理アドレスを制限でき、同じ仕組みで NIC へのアクセスを制限できる
- カーネルバイパス構成は、ユーザー空間プログラムへ NIC のレジスタへのアクセス、もしくは NIC と紐付けられた送受信データ配置用の DRAM 領域へのアクセスを許可する

プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A

メモリアクセス設定 2 パターン

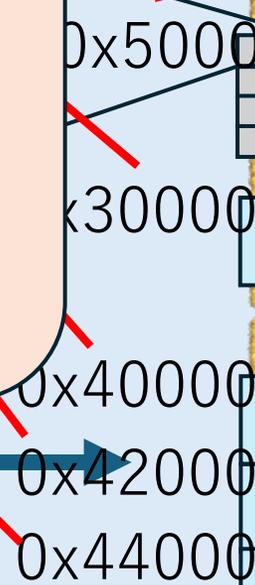
パターン 1 : NIC のレジスタへのアクセスを許可する

パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

DPDK はパターン 1、netmap と AF_XDP はパターン 2

```
struct {  
    void *address;  
    int length;  
} slot[ NUM_SLOT ];
```

[3] Address: 0x44000
length



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

- 一般的な（おそらく暗黙の）想定
 - ハードウェアに脆弱性および悪意はない
 - カーネルに脆弱性はなく悪意のあるコードは実行されない
 - ユーザー空間で動作するプログラムの悪意や脆弱性の可能性は想定するものの、操作性の観点から防御策は現実的な範囲に限定する

CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります
- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されている場合

アプリのクラッシュの
影響はパソコン全体には
波及しない

インストールしたアプリが
バグでクラッシュした！けど
パソコンはフリーズしない



CPU のモードと基本的な運用方針

- CPU のモードは大まかに分けて二種類
 1. **特権モード**：大体の CPU 機能を使える（システムの破壊が容易）
 2. **非特権モード**：できることが制限されている（破壊操作が難しい）

特権モードはカーネルモード、非特権モードはユーザーモードと呼ばれることもあります

- 基本的な運用方針：**なるべく破壊操作が難しいことを目指す**

この方針が適用されている場合

アプリのクラッシュの
影響はパソコン全体には
波及しない

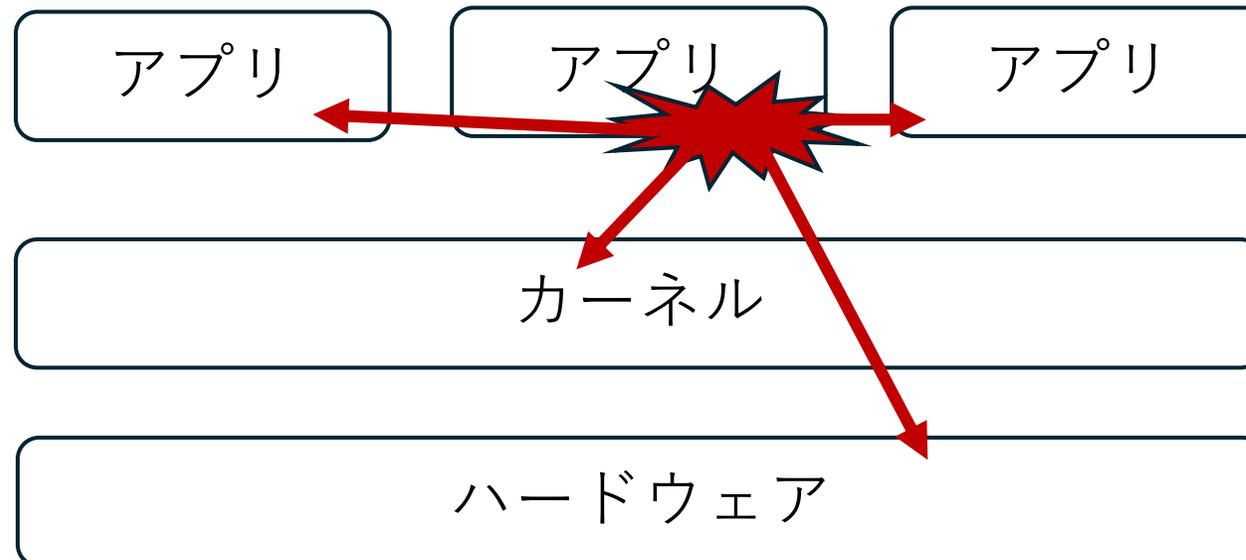
インストールしたアプリが
バグでクラッシュした！けど
パソコンはフリーズしない



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

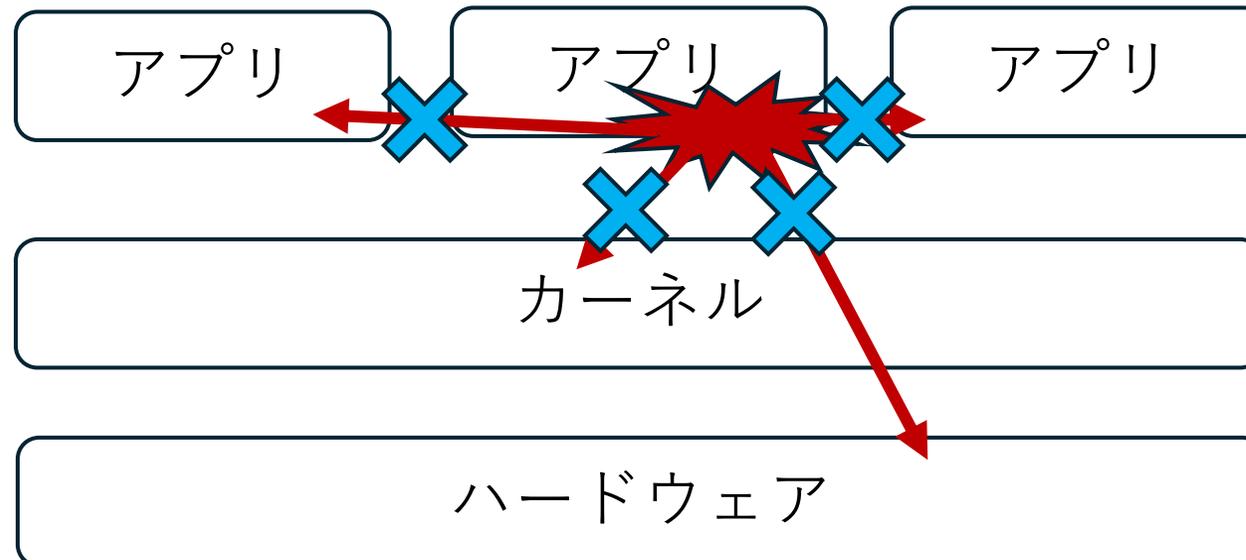
- 基本的な目標：一つのアプリのバグや脆弱性の影響がハードウェア、カーネル、別のアプリへ波及しにくくする
- カーネルバイパス構成を適用すると何が変わる？



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

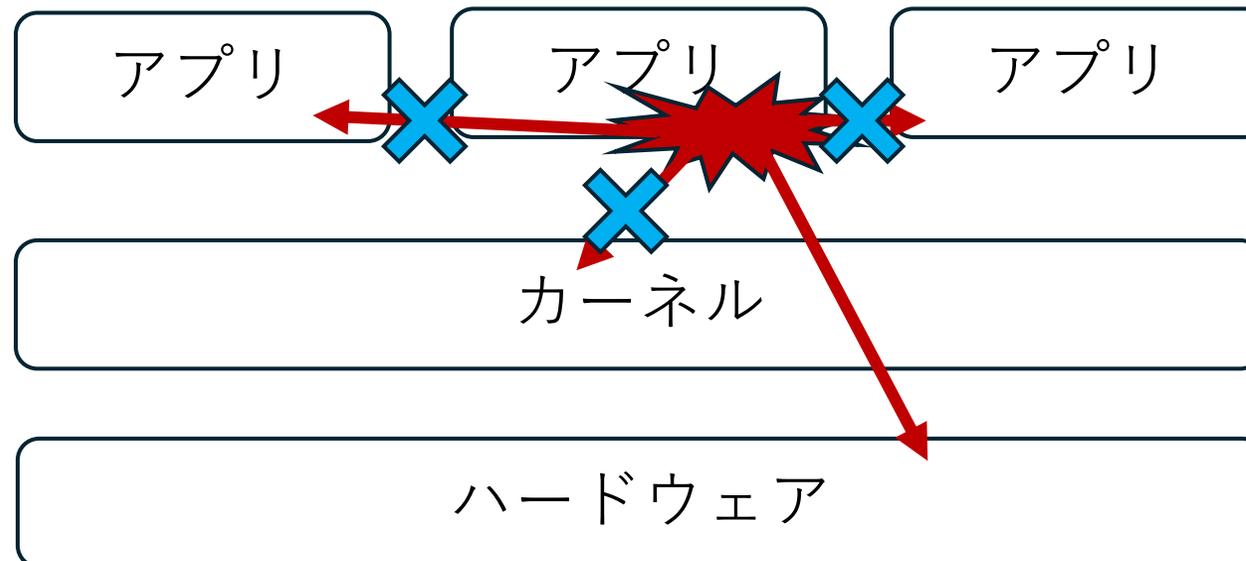
- 基本的な目標：一つのアプリのバグや脆弱性の影響がハードウェア、カーネル、別のアプリへ波及しにくくする



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

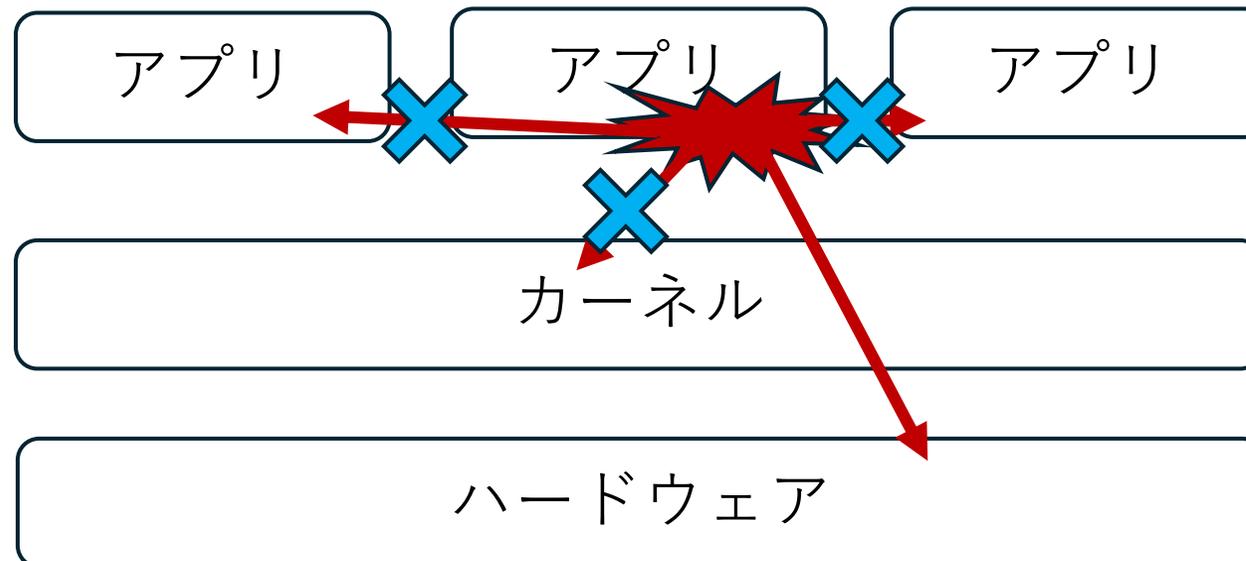
- 基本的な目標：一つのアプリのバグや脆弱性の影響がハードウェア、カーネル、別のアプリへ波及しにくくする
- カーネルバイパス構成を適用すると何が変わる？



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

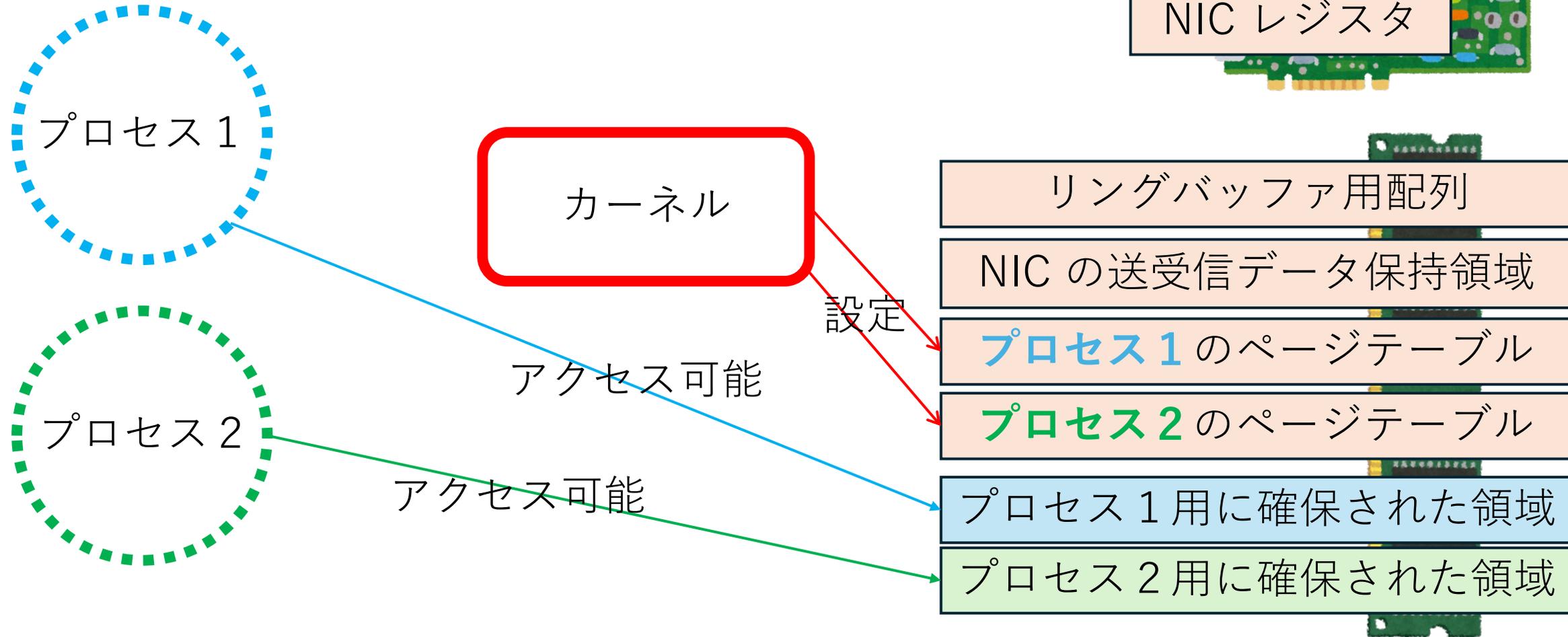
- 基本的な目標：一つのアプリのバグや脆弱性の影響がハードウェア、カーネル、別のアプリへ波及しにくくする
- **影響の波及を防ぐポイント：そもそも複数の要素が同一のソースへアクセスしないようにする**



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

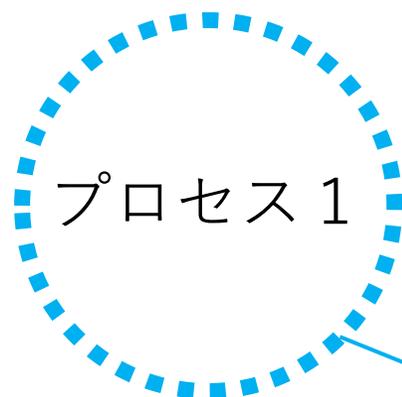
複数の要素が同一のリソースへアクセスしないようにする



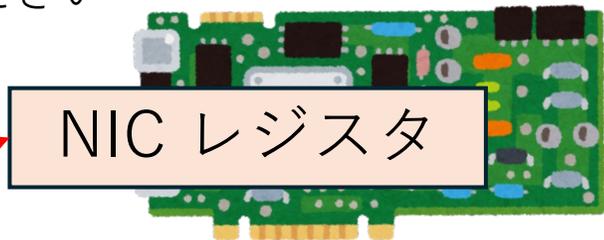
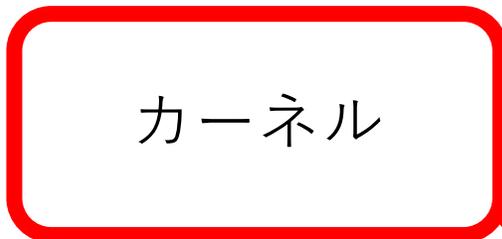
セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

複数の要素が同一のリソースへアクセスしないようにする



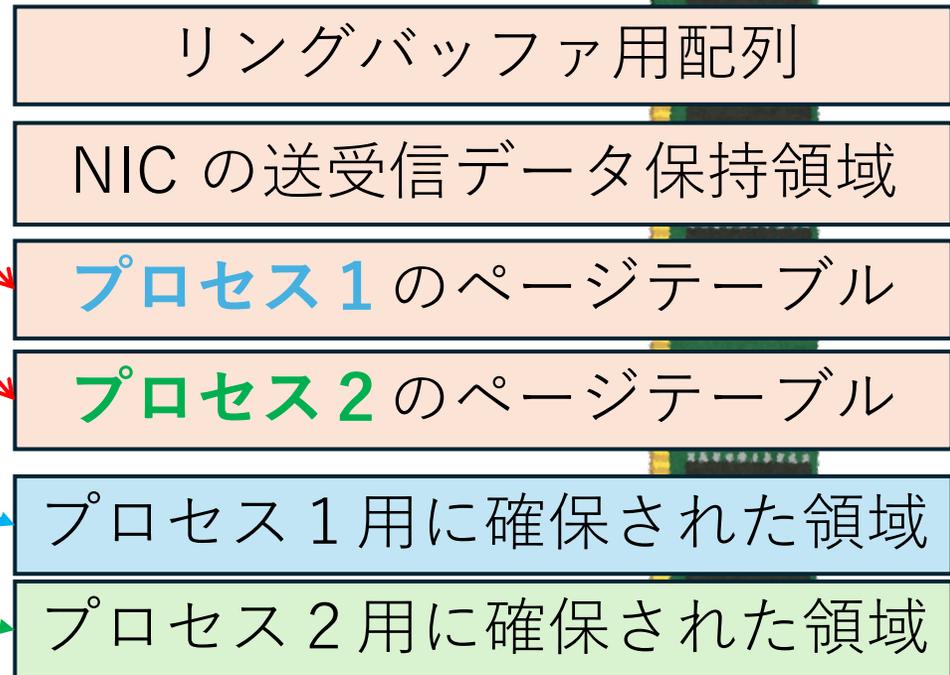
一般的な構成ではカーネルしか
NIC のレジスタにアクセスしない



設定

アクセス可能

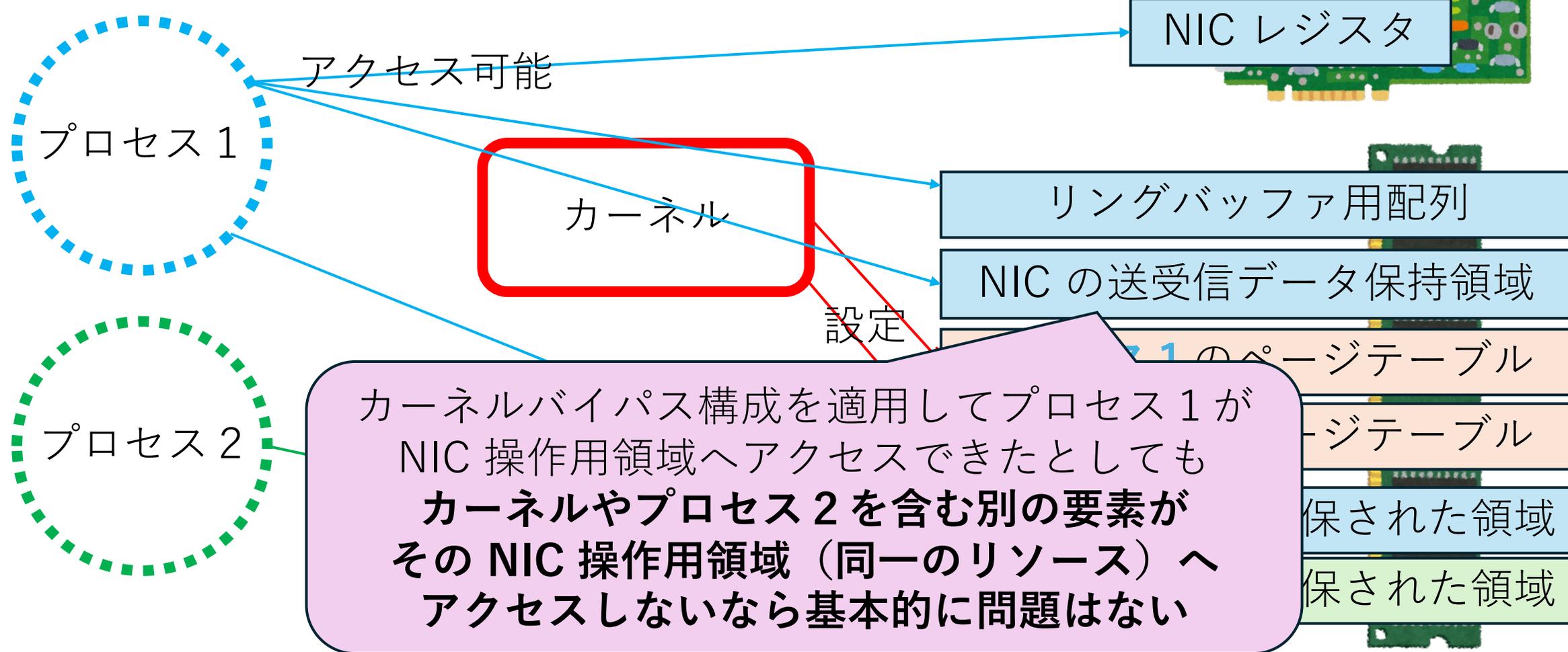
アクセス可能



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

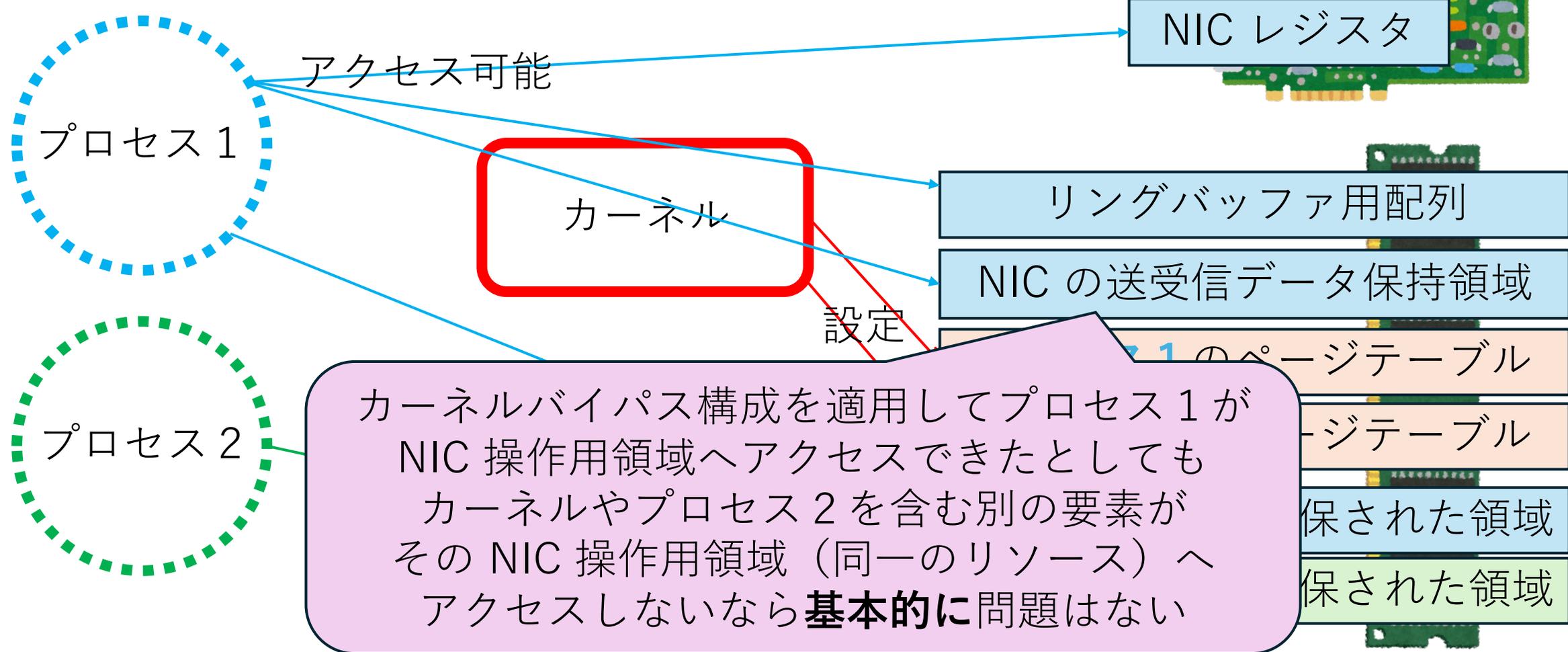
複数の要素が同一のリソースへアクセスしないようにする



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

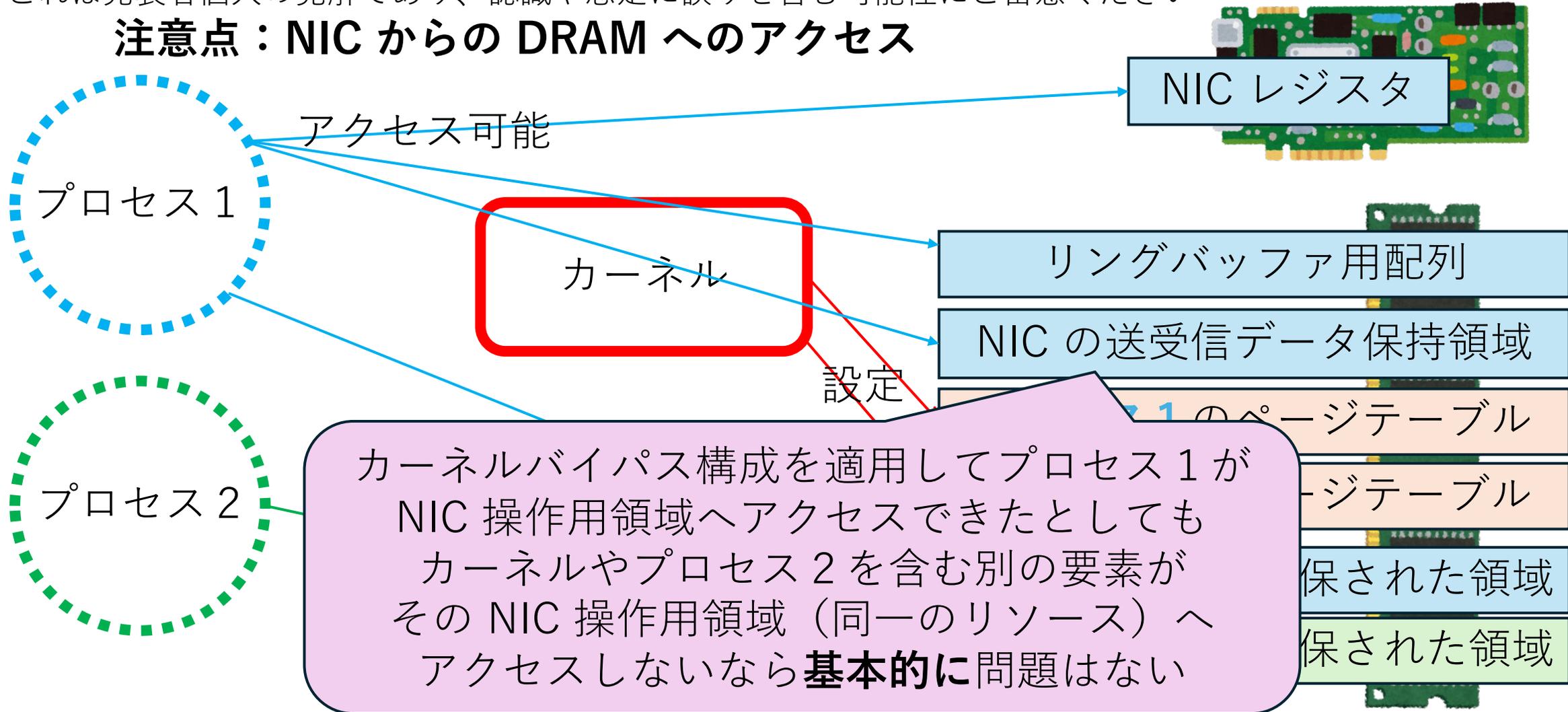
複数の要素が同一のリソースへアクセスしないようにする



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス



プログラマー
このキュー

NIC は配列中の head と tail の間の区間が参照する DRAM 上のデータを外部へ送信する

Q. 送受信

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラマー

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラマーは NIC のレジスタへ値を書き込むことでパケットの送信開始をリクエストする

NIC のレジスタ

ring_head:	0
ring_tail:	2
ring_address:	0x5000
ring_size:	4

DRAM 上の配列

[0]	address: 0x30000	length: 500
[1]	address: 0x40000	length: 800
[2]	address	length
[3]	address	length

head

tail

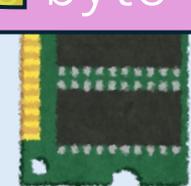
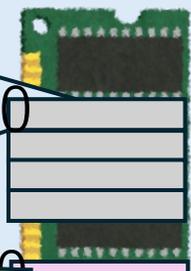
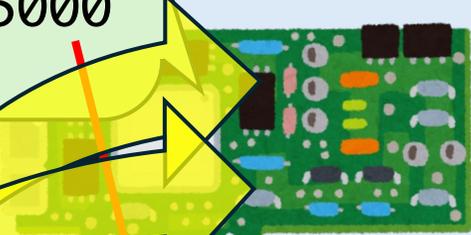
0x5000

0x30000

0x40000

500 byte

800 byte



このキュー

Q. 送受信

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは NIC のレジスタへ値を書き込むことでパケットの送信開始をリクエストする

NIC は配列中の head と tail の間の区間が参照する DRAM 上のデータを外部へ送信する

NIC のレジスタ

```
ring_head: 0
ring_tail: 2
ring_address: 0x5000
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address: 0x40000 length: 800
[2]	address length
[3]	address length

0x5000

0x30000

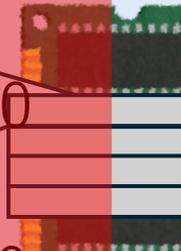
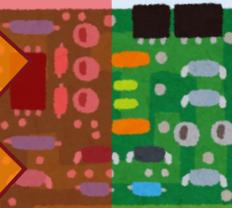
500 byte

0x40000

800 byte

head

tail



プログラムのこのキューが受信

NIC が受信したパケットを配列の head で参照される DRAM 領域へ書き込む

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成される？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head
tail

DRAM 上の配列

[0]	address: 0x30000 length
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length

0x5000

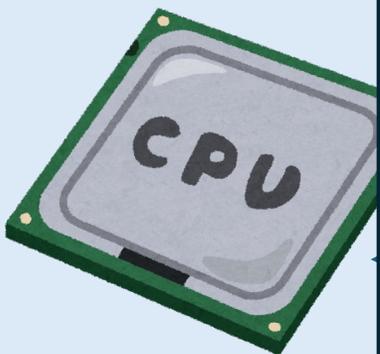
0x30000

100 byte

0x40000

0x42000

0x44000



プログラム

このキューが受信

NIC が受信したパケットを配列の head で参照される DRAM 領域へ書き込む

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成される?

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

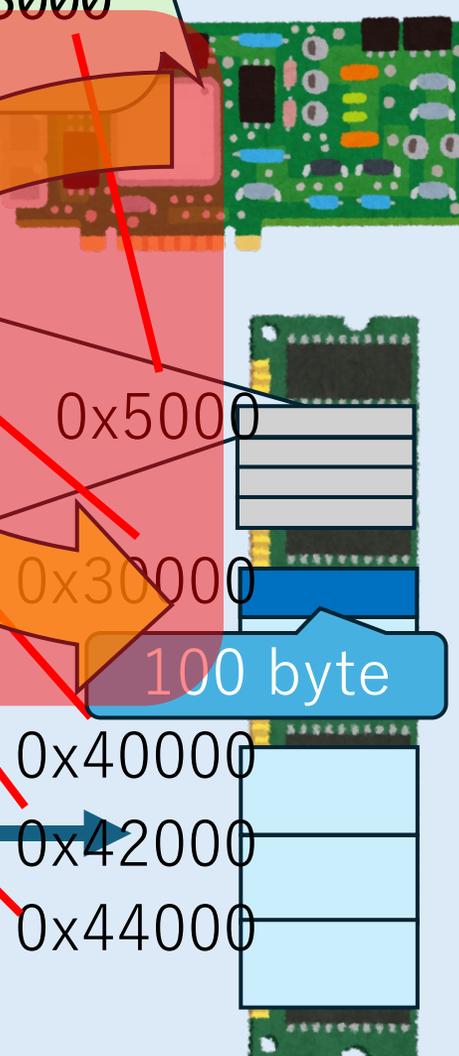
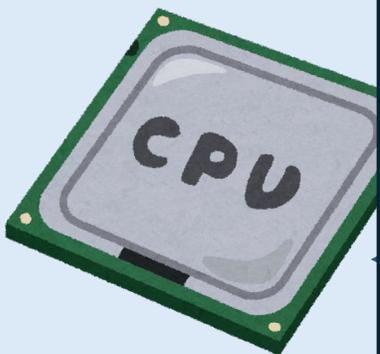
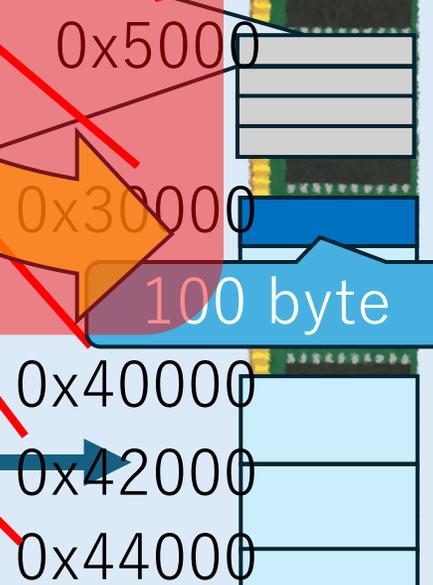
簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
struct {
    void *address;
    int length;
} slot[NUM_SLOT];
```

head
tail

DRAM 上の配列

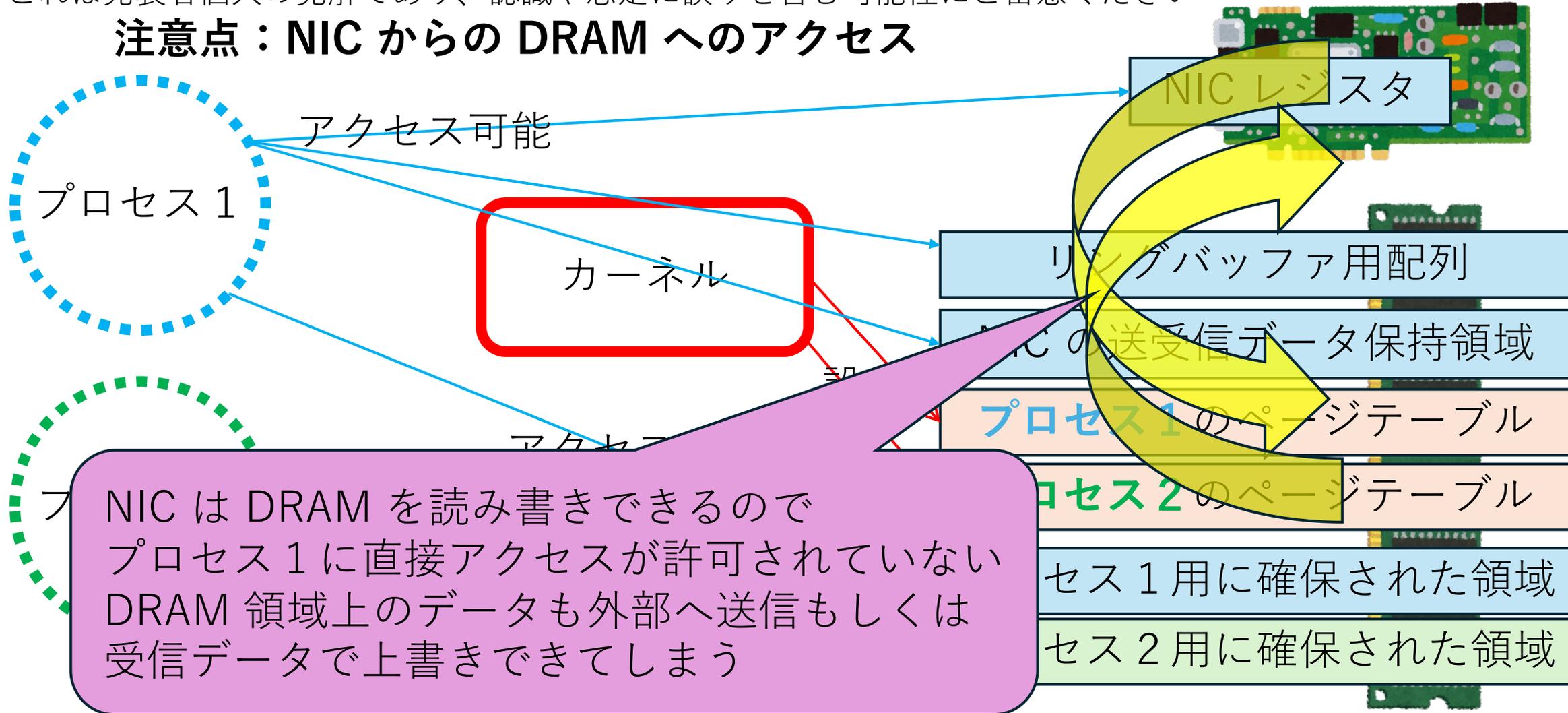
[0]	address: 0x30000 length
[1]	address: 0x40000 length
[2]	address: 0x42000 length
[3]	Address: 0x44000 length



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス



プログラムによる NIC

このキューが送信キューとして使われる場合

Q. 送受信キューはどのように構成されるか？

A. キューは NIC のレジスタとメモリ上のデータで構成されるリングバッファ

プログラム

簡単な実装

```
int head;
int tail;
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを通して配列に送信データへの参照を設定

```
} slot[NUM_SLOT];
```

head

tail

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address length
[2]	address length
[3]	address length

NIC のレジスタ

```
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4
```

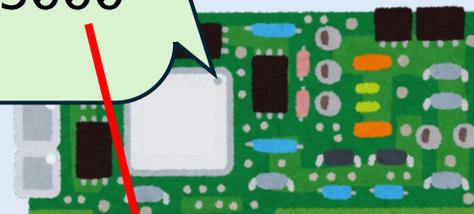
0x5000

0x30000

500 byte

0x40000

800 byte



プログラムによる NIC

このキューが送信キューとして使われる場合

物理アドレス 0x30000 の領域は
プロセス 1 に割り当てられていなくても
プロセス 1 は配列に参照を設定することで
NIC からの送信対象に設定できる

NIC のレジスタ

```
ring_head: 0  
ring_tail: 0  
ring_address: 0x5000  
ring_size: 4
```

DRAM 上の配列

[0]	address: 0x30000 length: 500
[1]	address length
[2]	address length
[3]	address length

プログラム

簡単な実装

```
int head;  
int tail;  
#define NUM_SLOT 4
```

プログラムは DRAM への書き込みを
通して配列に送信データへの参照を設定

```
} slot[NUM_SLOT];
```

head

tail

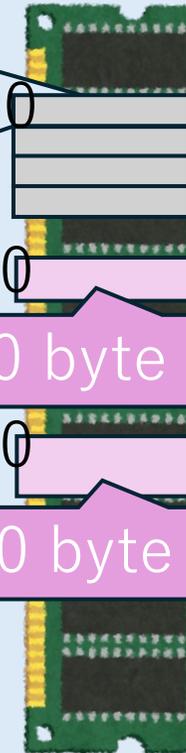
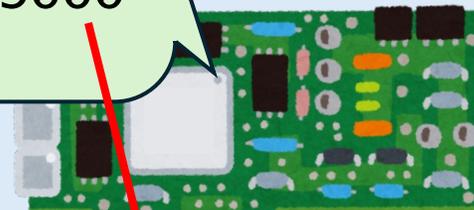
0x5000

0x30000

500 byte

0x40000

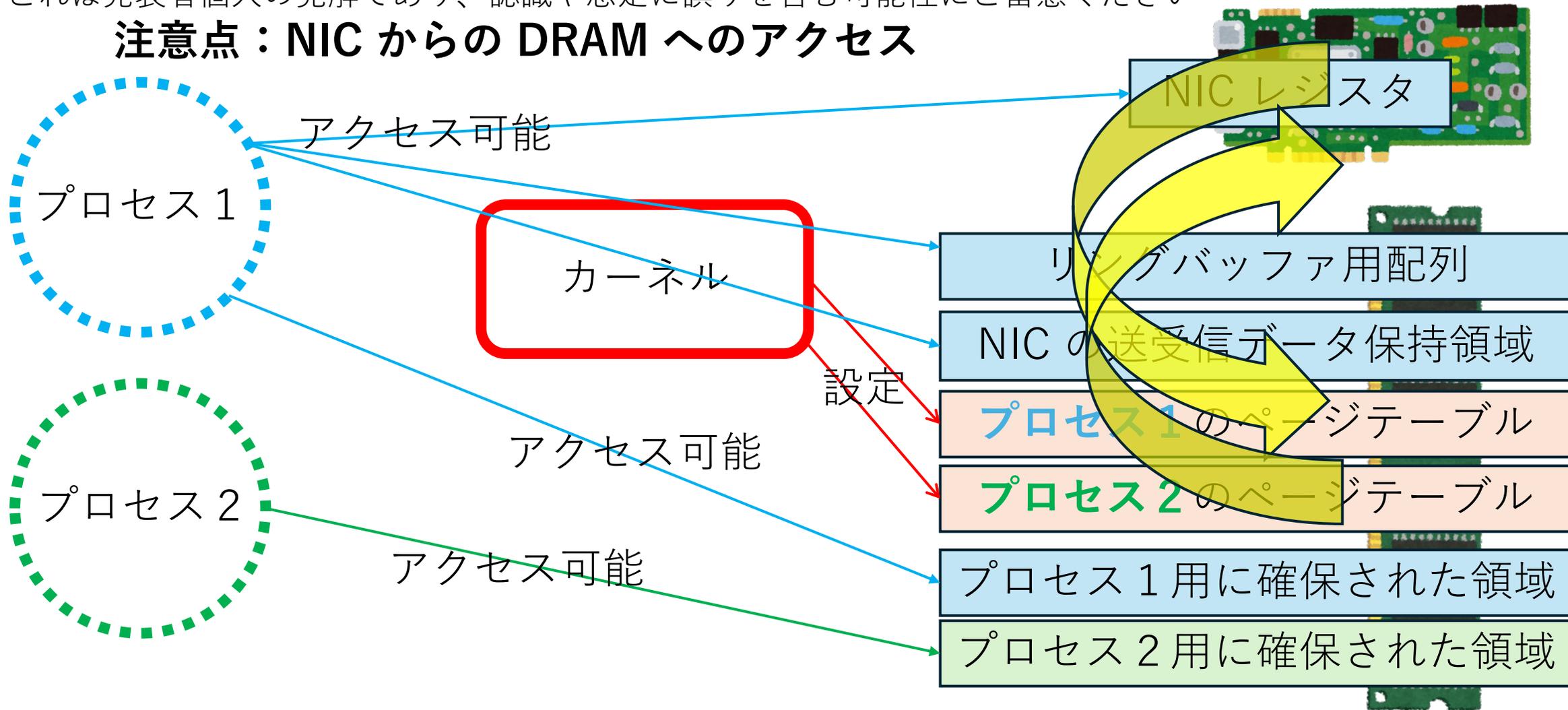
800 byte



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

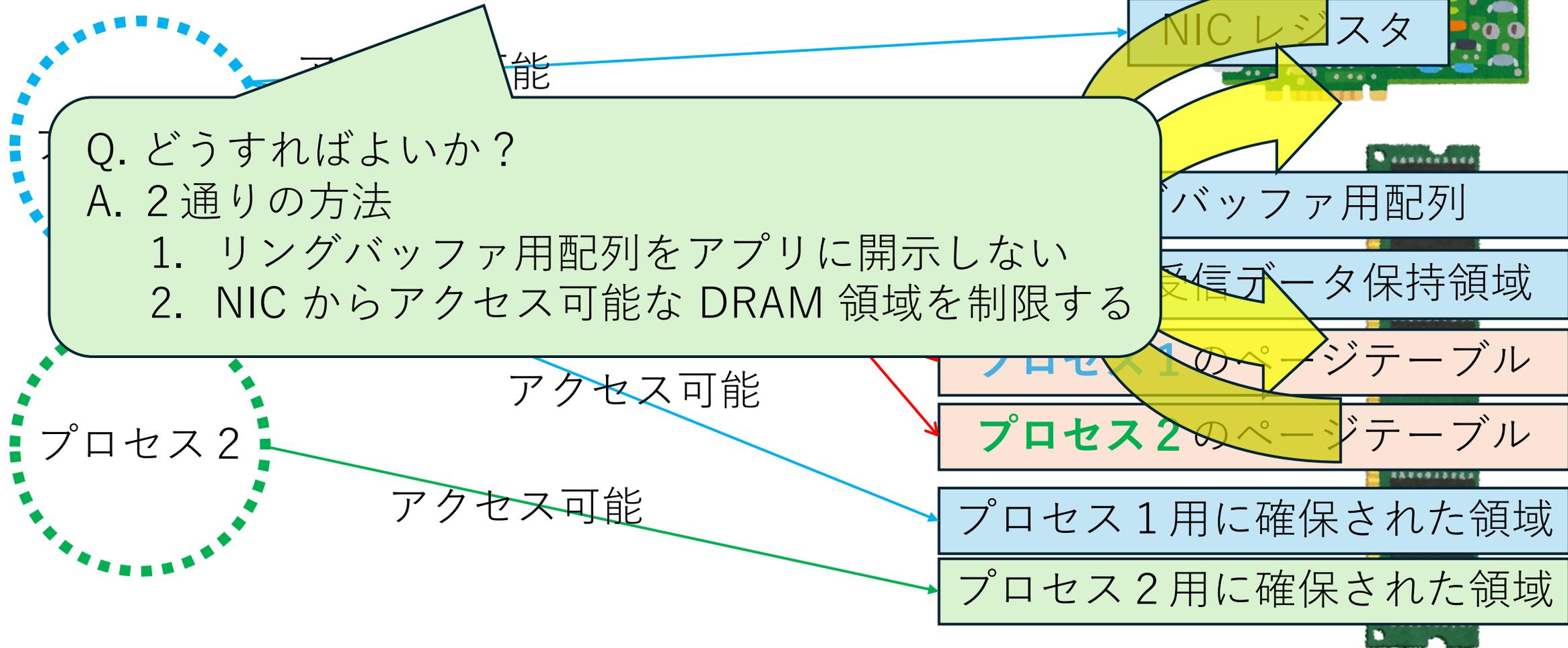
注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

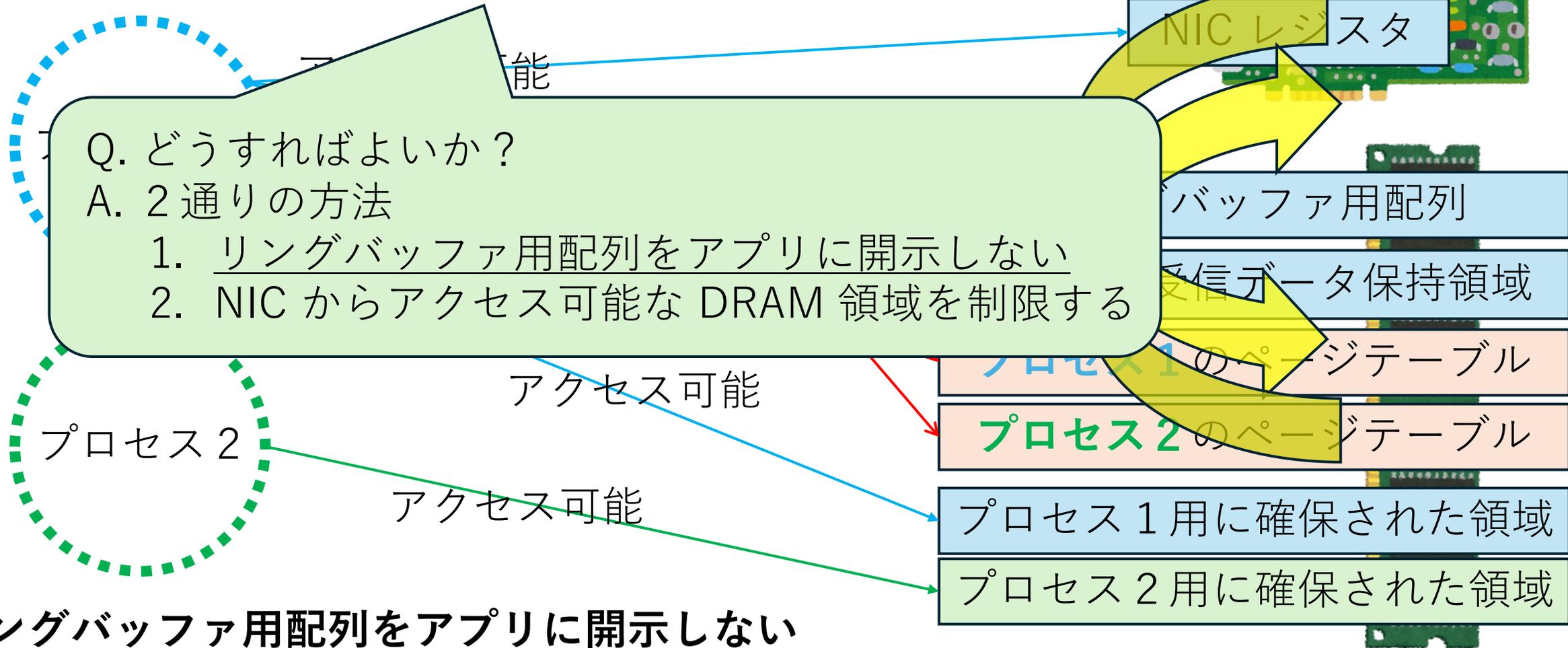
注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

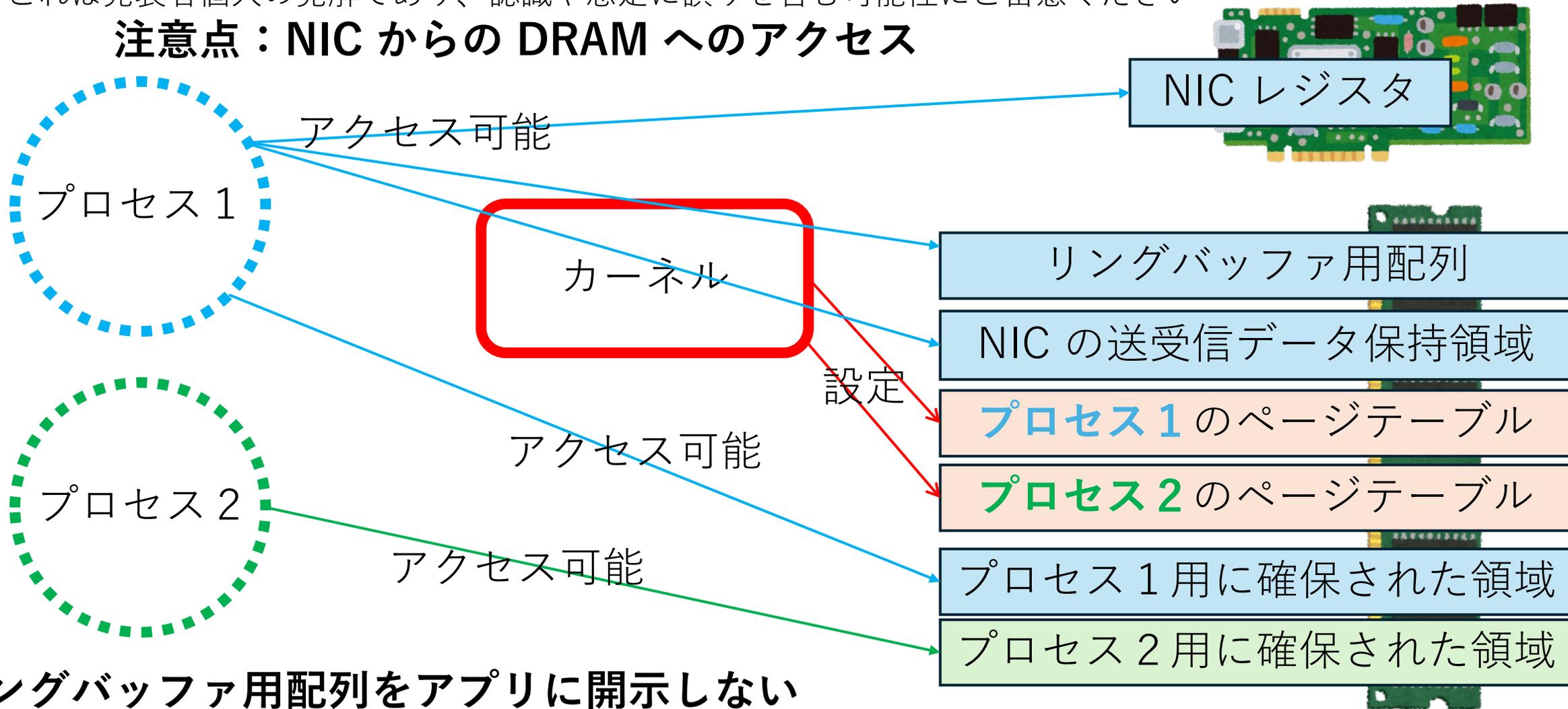
注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

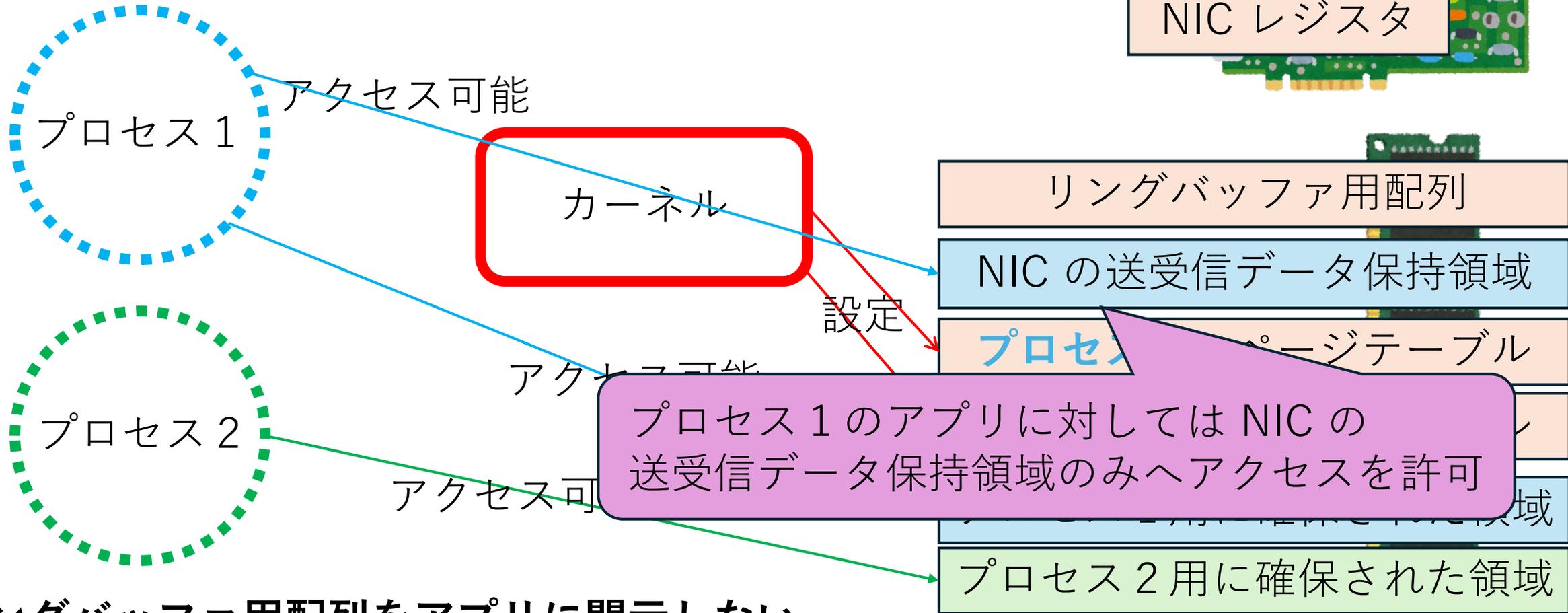
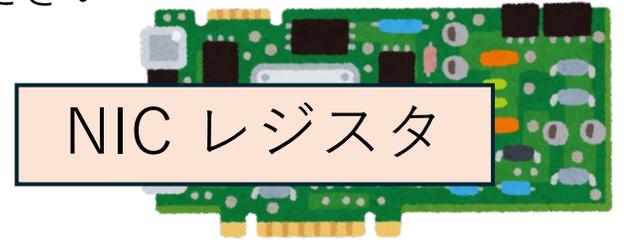


リングバッファ用配列をアプリに開示しない

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス



リングバッファ用配列をアプリに開示しない

セキュリティ

これは発生
注

プロセス1のアプリはカーネルへリクエスト
(専用のシステムコールの呼び出し)を通して
代わりにNICを操作してもらう

ください



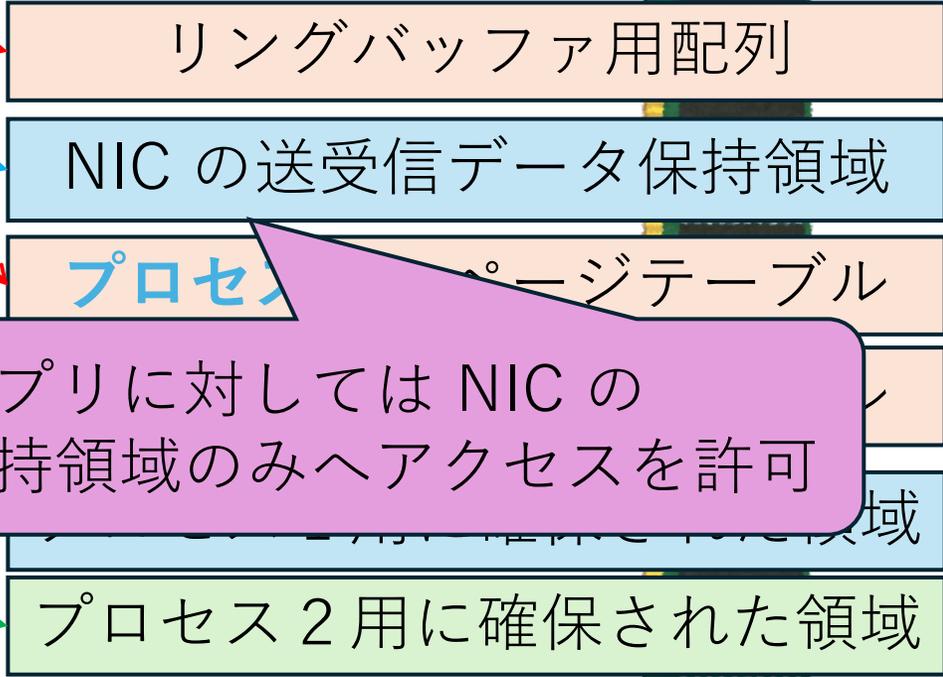
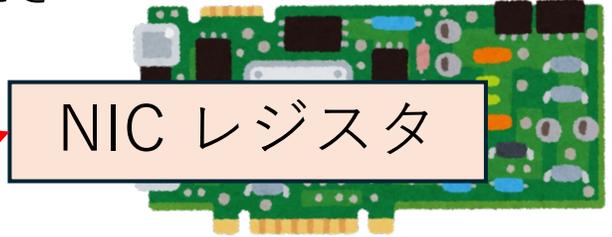
アクセス可能

リクエスト



操作

設定



プロセス1のアプリに対してはNICの
送受信データ保持領域のみへアクセスを許可



アクセス可能

リングバッファ用配列をアプリに開示しない

セキュリティ

これは発生
注

プロセス1のアプリはカーネルへリクエスト
(専用のシステムコールの呼び出し)を通して
代わりにNICを操作してもらう



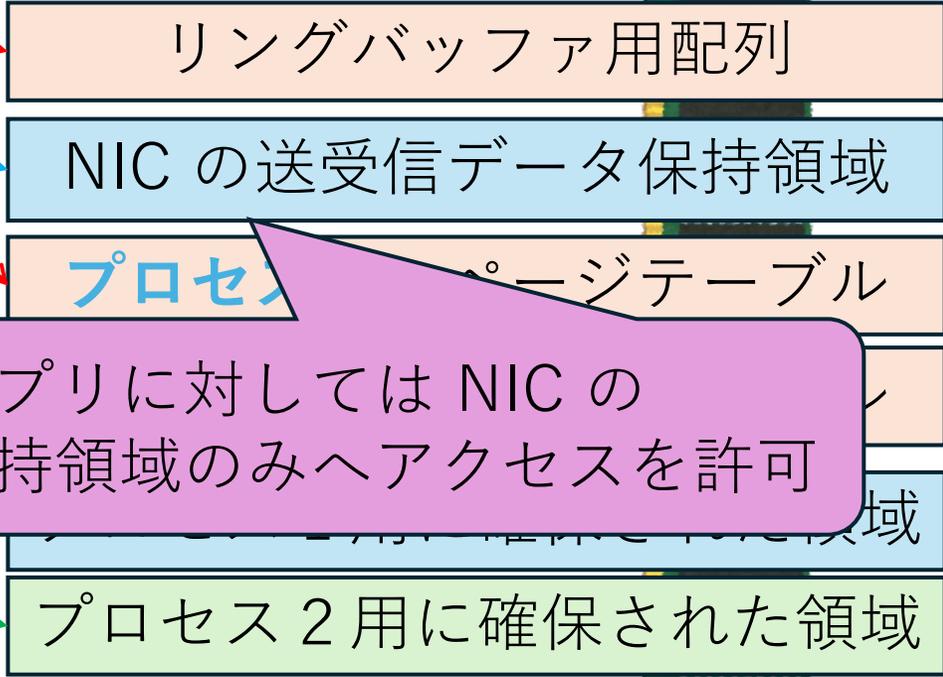
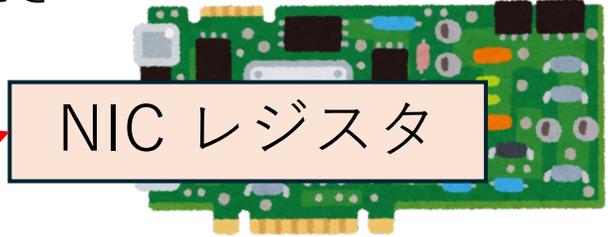
アクセス可能

リクエスト



操作

設定



想定として
アプリは信頼されず
カーネルは信頼されている
というところがポイントです

プロセス1のアプリに対してはNICの
送受信データ保持領域のみへアクセスを許可

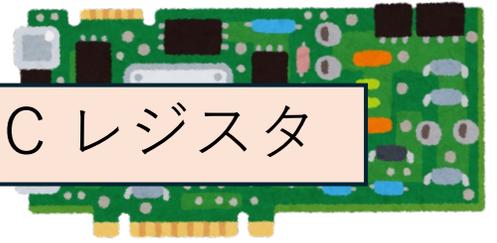
リングバッファ用配列をアプリに開示しない

セキュリティ

これは発生
注

プロセス1のアプリはカーネルヘリクエスト
(専用のシステムコールの呼び出し)を通して
代わりにNICを操作してもらう

ください



NIC レジスタ



アクセス可能

操作

仮にアプリをカーネルと同じく信頼する想定であれば
リングバッファ用配列を直接開示して問題ないです

想定として

アプリは信頼されず
カーネルは信頼されている
というところがポイントです

プロセス

ページテーブル

プロセス1のアプリに対してはNICの
送受信データ保持領域のみへアクセスを許可

プロセス2用に確保された領域

リングバッファ用配列をアプリに開示しない

プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A

メモリアクセス設定 2 パターン

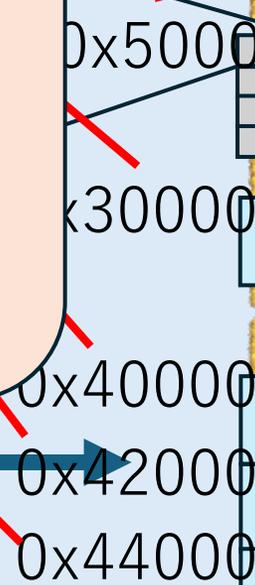
パターン 1 : NIC のレジスタへのアクセスを許可する

パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

DPDK はパターン 1、netmap と AF_XDP はパターン 2

```
struct {  
    void *address;  
    int length;  
} slot[ NUM_SLOT ];
```

[3] Address: 0x44000
length



プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A

メモリアクセス設定 2 パターン

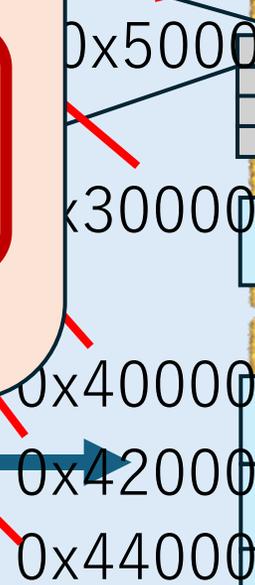
パターン 1 : NIC のレジスタへのアクセスを許可する

パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

DPDK はパターン 1、netmap と AF_XDP はパターン 2

```
struct {  
    void *address;  
    int length;  
} slot[ NUM_SLOT ];
```

[3] Address: 0x44000
length



プログラムによる NIC

NIC のレジスタ
ring_head: 0
ring_tail: 0
ring_address: 0x5000
ring_size: 4

Q. 送受信キューはどのように構成されるか？

A. メモリアクセス設定 2 パターン

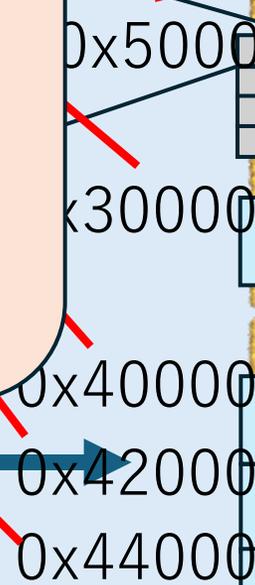
パターン 1 : NIC のレジスタへのアクセスを許可する

パターン 2 : NIC に紐付けした DRAM 領域だけアクセス許可し
アプリが NIC の送受信の起動のために利用可能な
専用システムコールをカーネルが提供する

DPDK はパターン 1、netmap と AF_XDP はパターン 2

```
void *address;  
int length;  
} slot[ NUM_SLOT ];
```

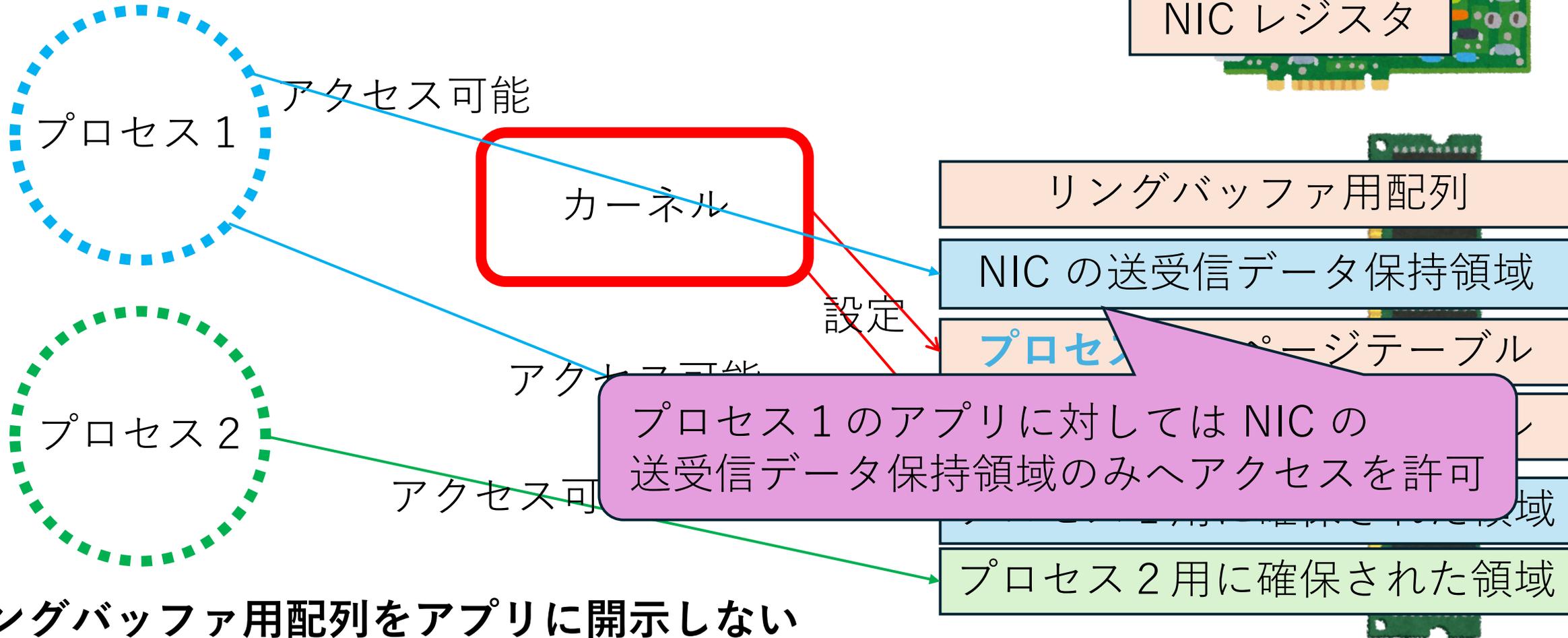
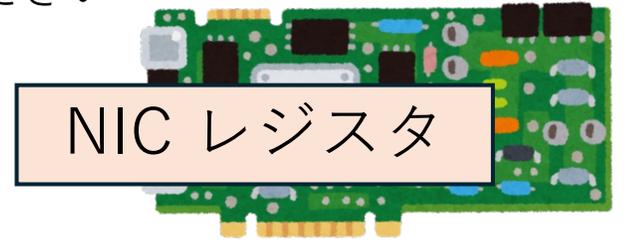
[3] Address: 0x44000
length



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

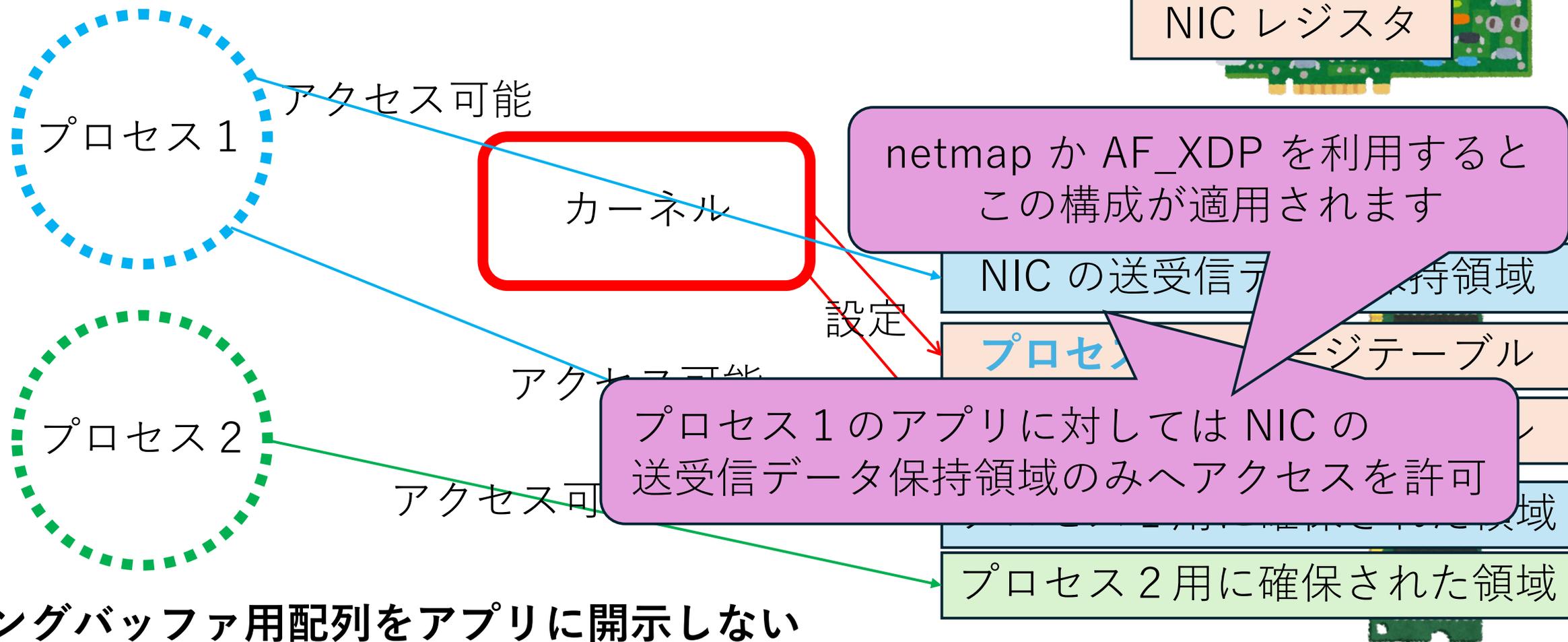
注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

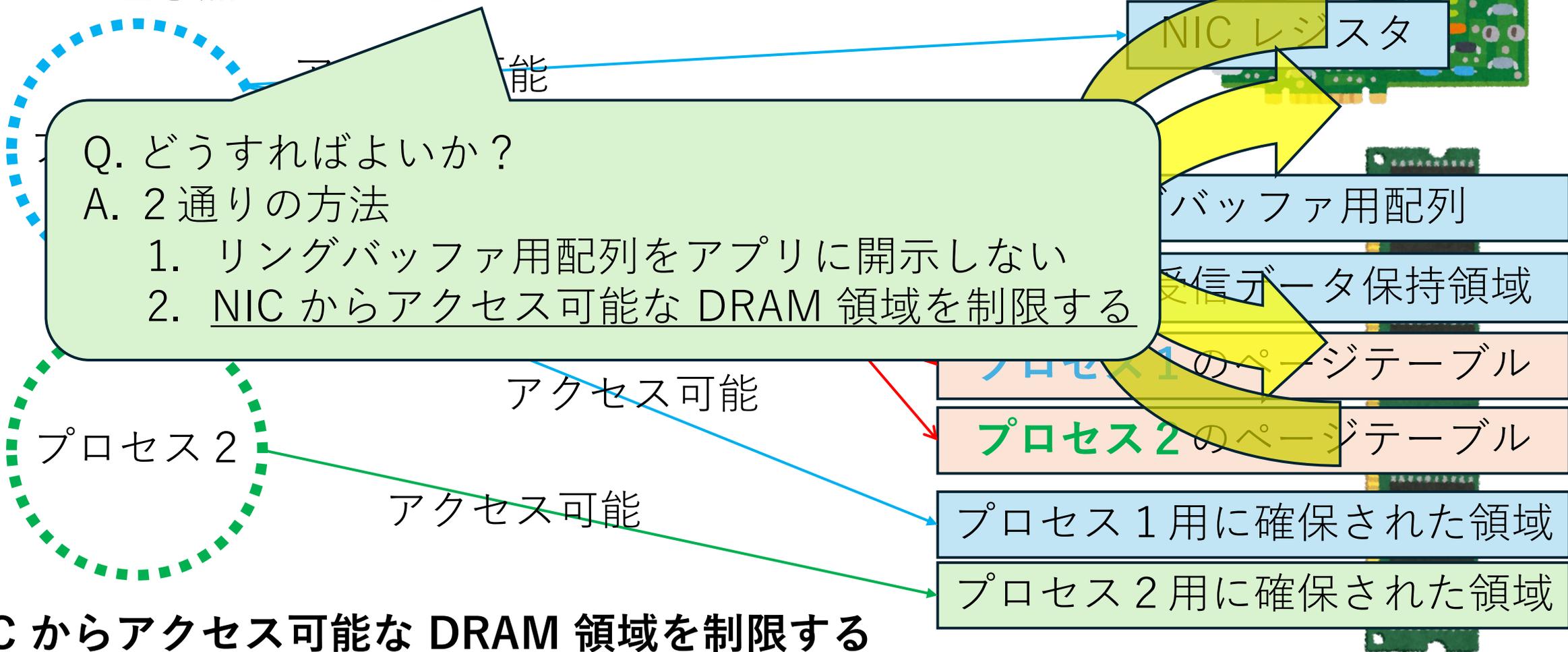
注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

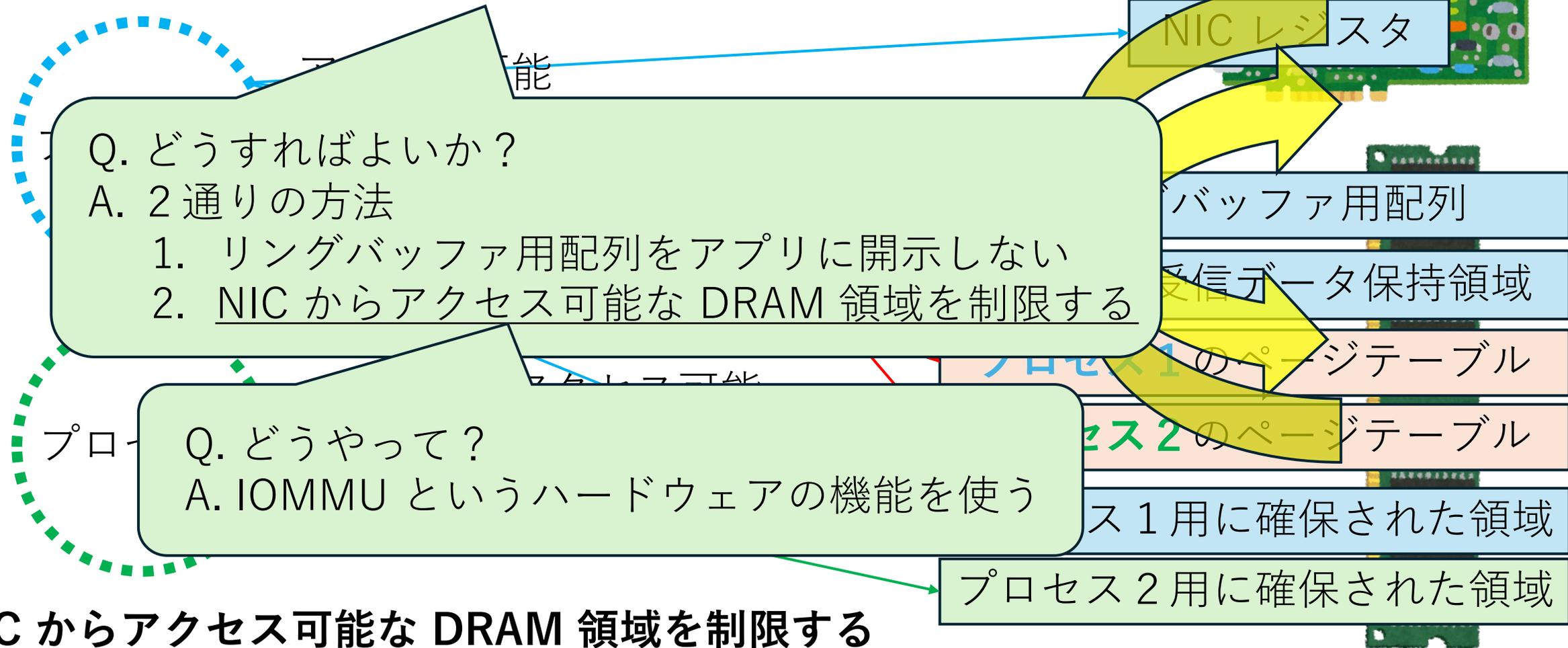


NIC からアクセス可能な DRAM 領域を制限する

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

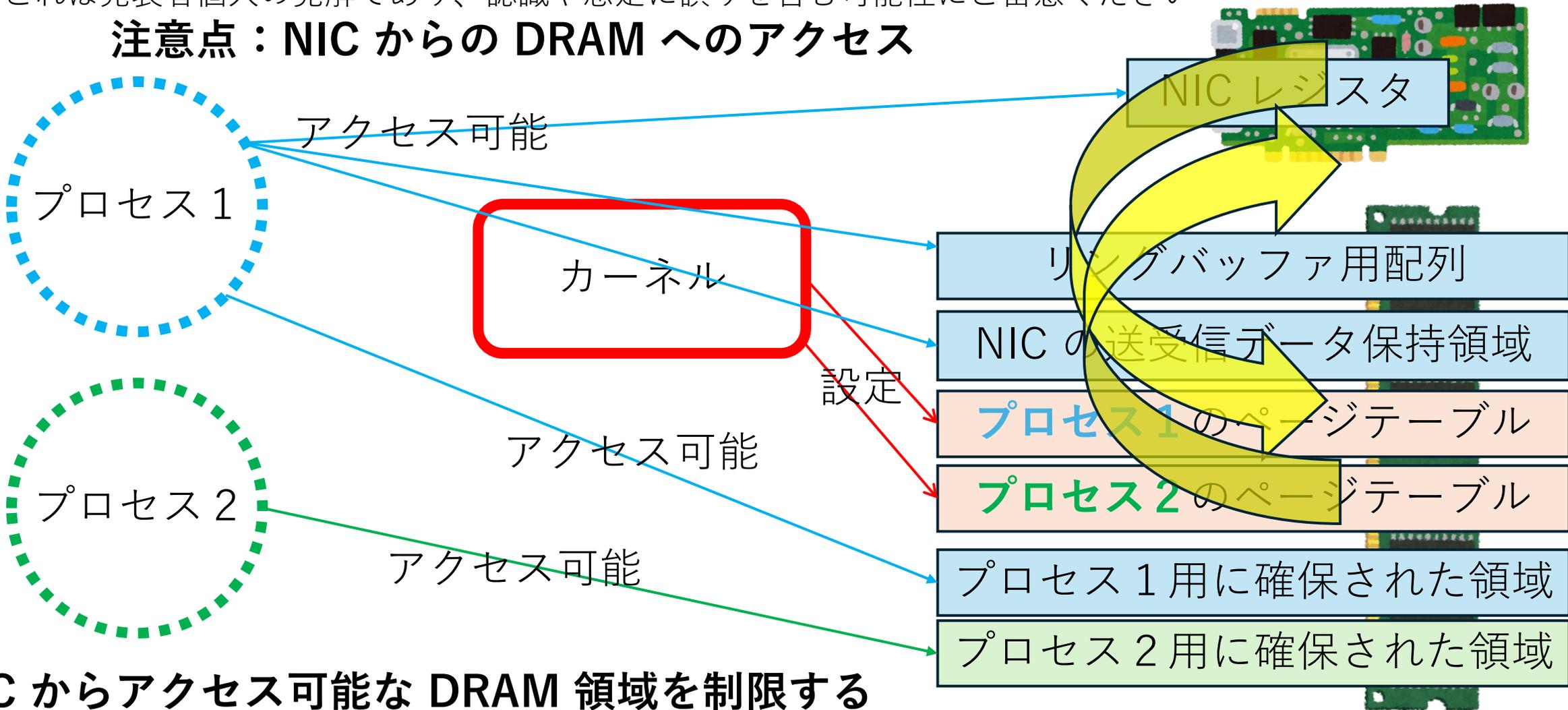


NIC からアクセス可能な DRAM 領域を制限する

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

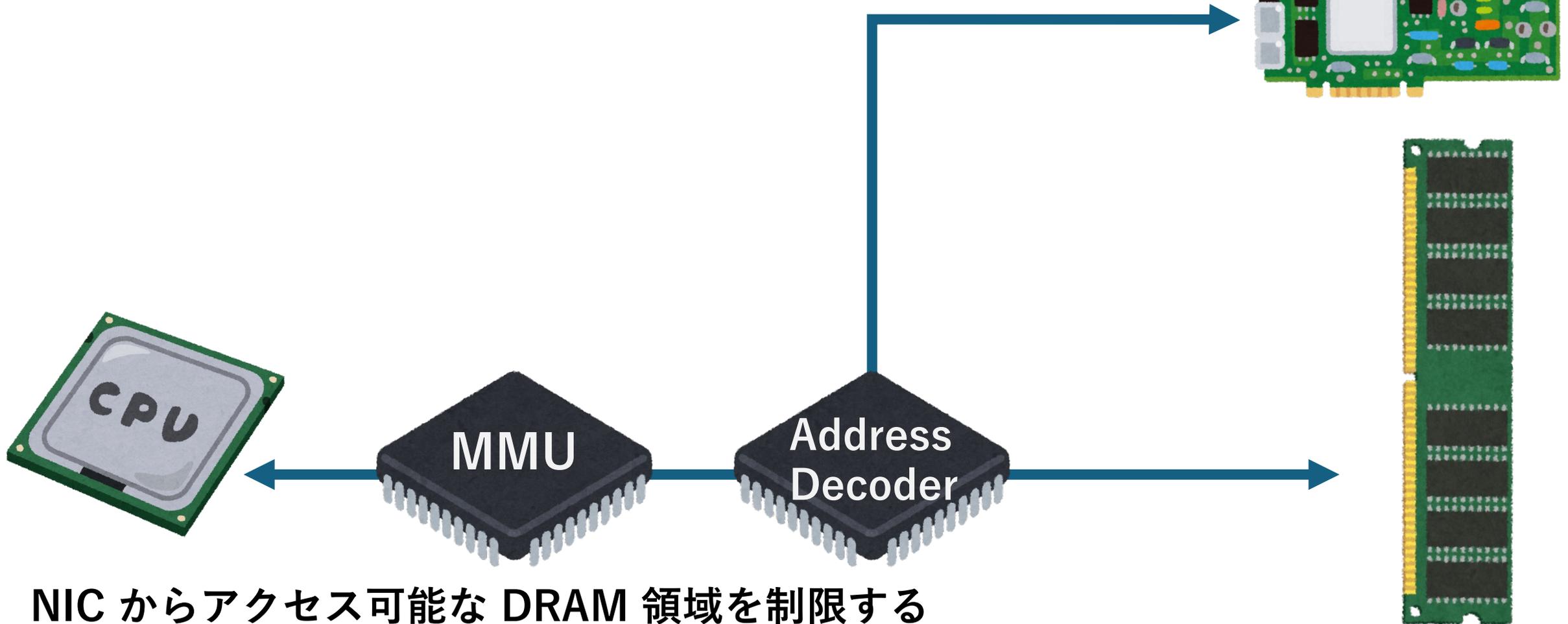


NIC からアクセス可能な DRAM 領域を制限する

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス



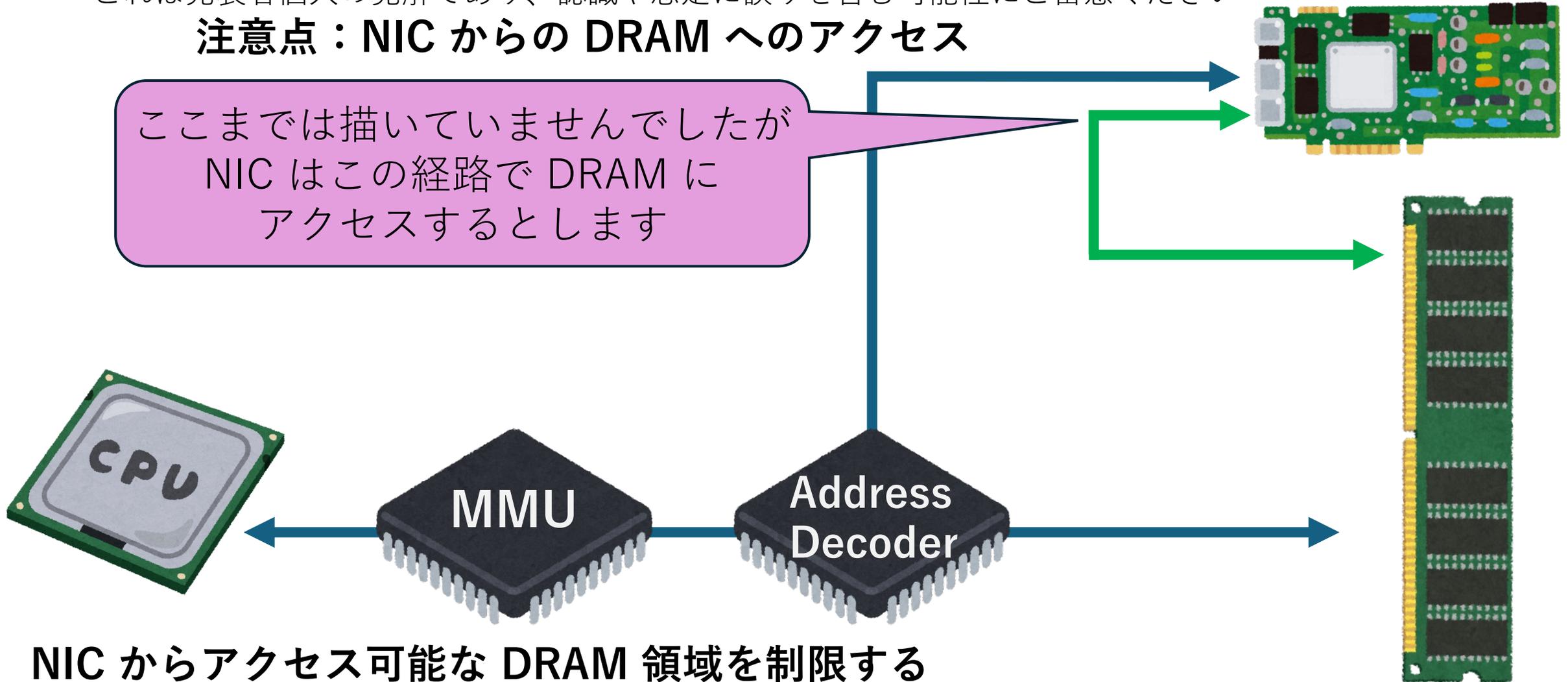
NIC からアクセス可能な DRAM 領域を制限する

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

ここまでは描いていませんでしたが
NICはこの経路で DRAM に
アクセスするとします

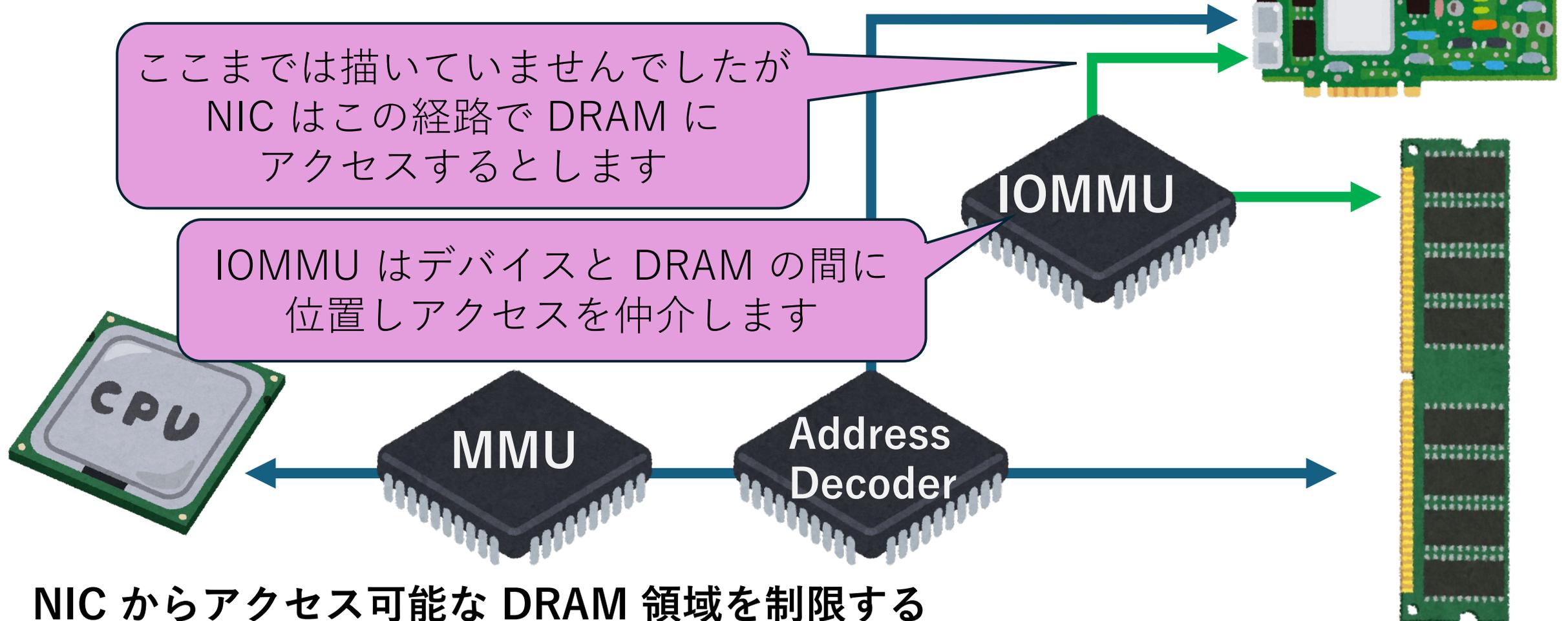


NIC からアクセス可能な DRAM 領域を制限する

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

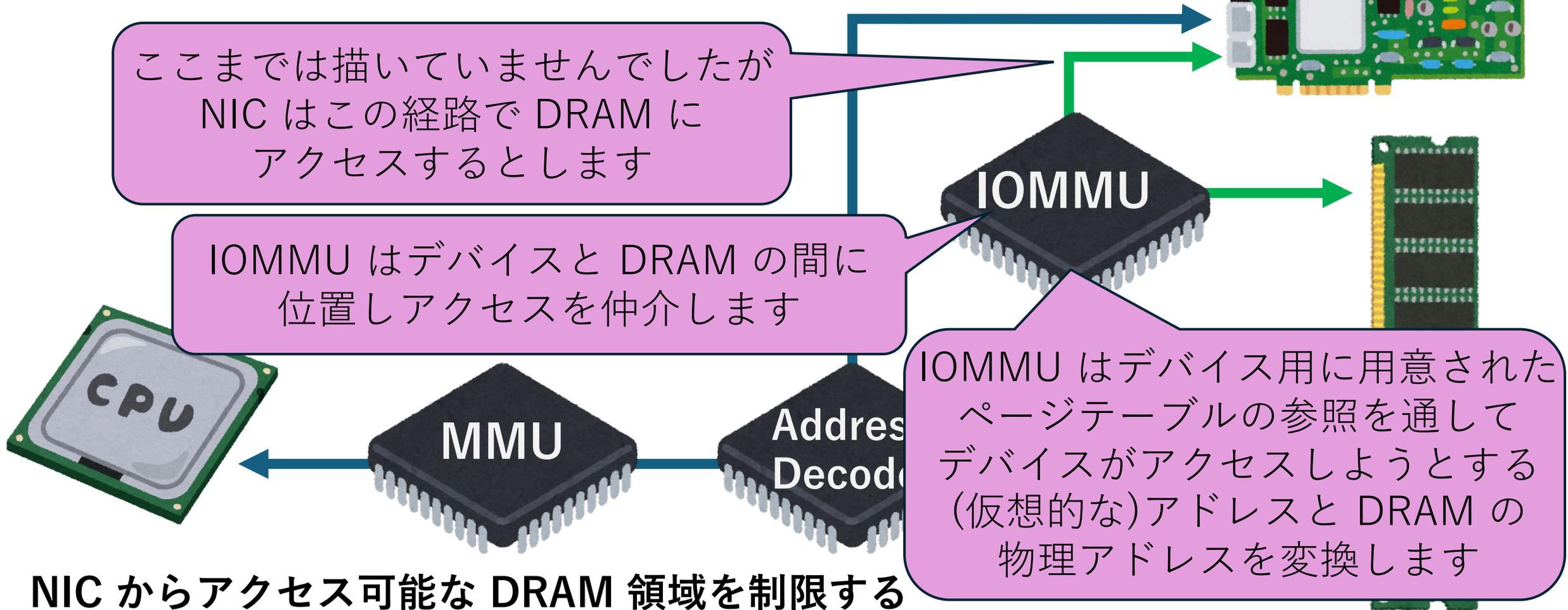


NIC からアクセス可能な DRAM 領域を制限する

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

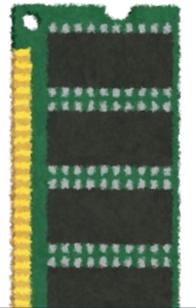
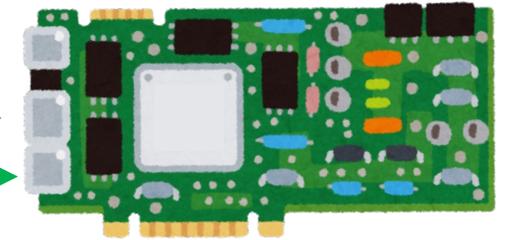
ここまでは描いていませんでしたが
NIC はこの経路で DRAM に
アクセスするとします

IOMMU はデバイスと DRAM の間に
位置しアクセスを仲介します

IOMMU が有効な場合には NIC からは
このページテーブルに記載のある
DRAM の物理アドレスにしか
アクセスできません

NIC からアクセス可能な DRAM 領域を制限する

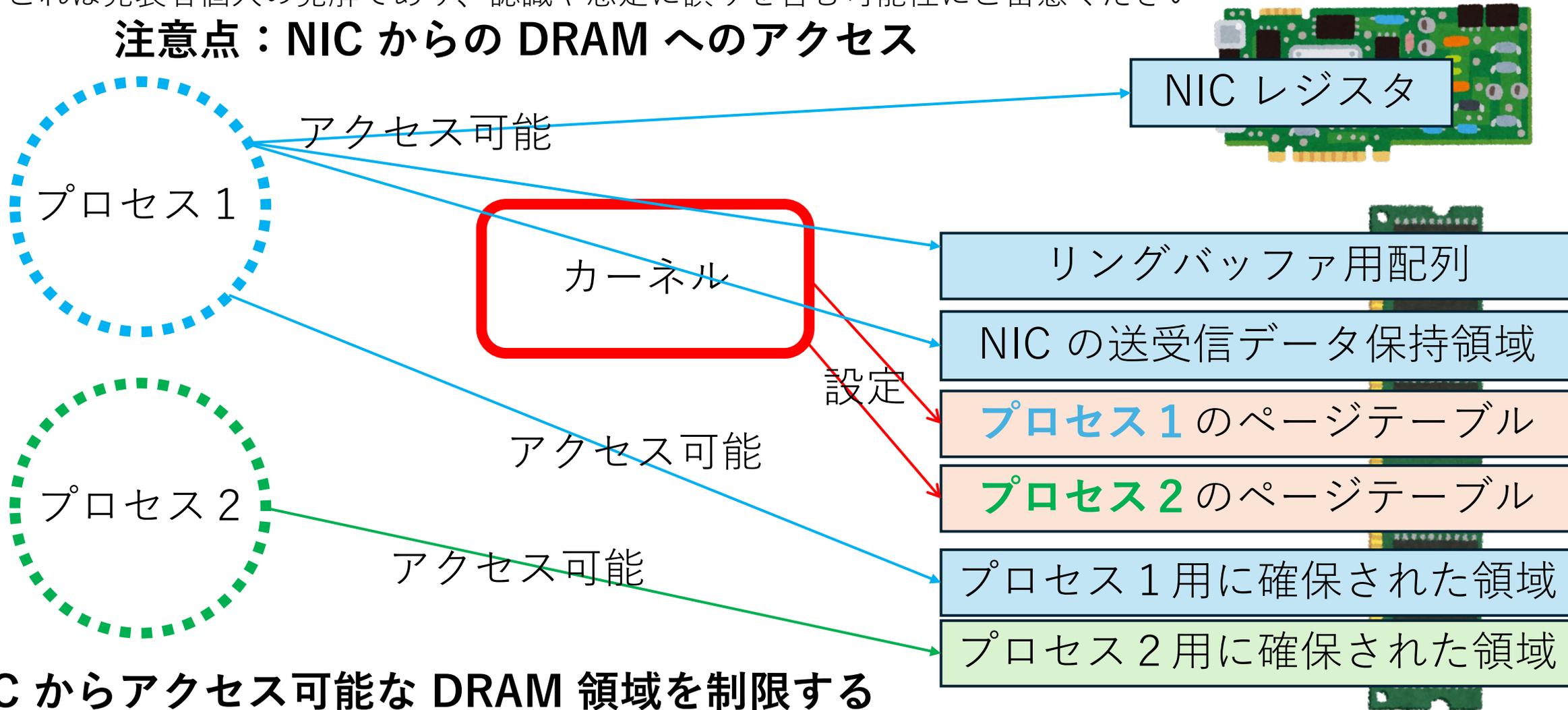
IOMMU はデバイス用に用意された
ページテーブルの参照を通して
デバイスがアクセスしようとする
(仮想的な)アドレスと DRAM の
物理アドレスを変換します



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

注意点：NIC からの DRAM へのアクセス

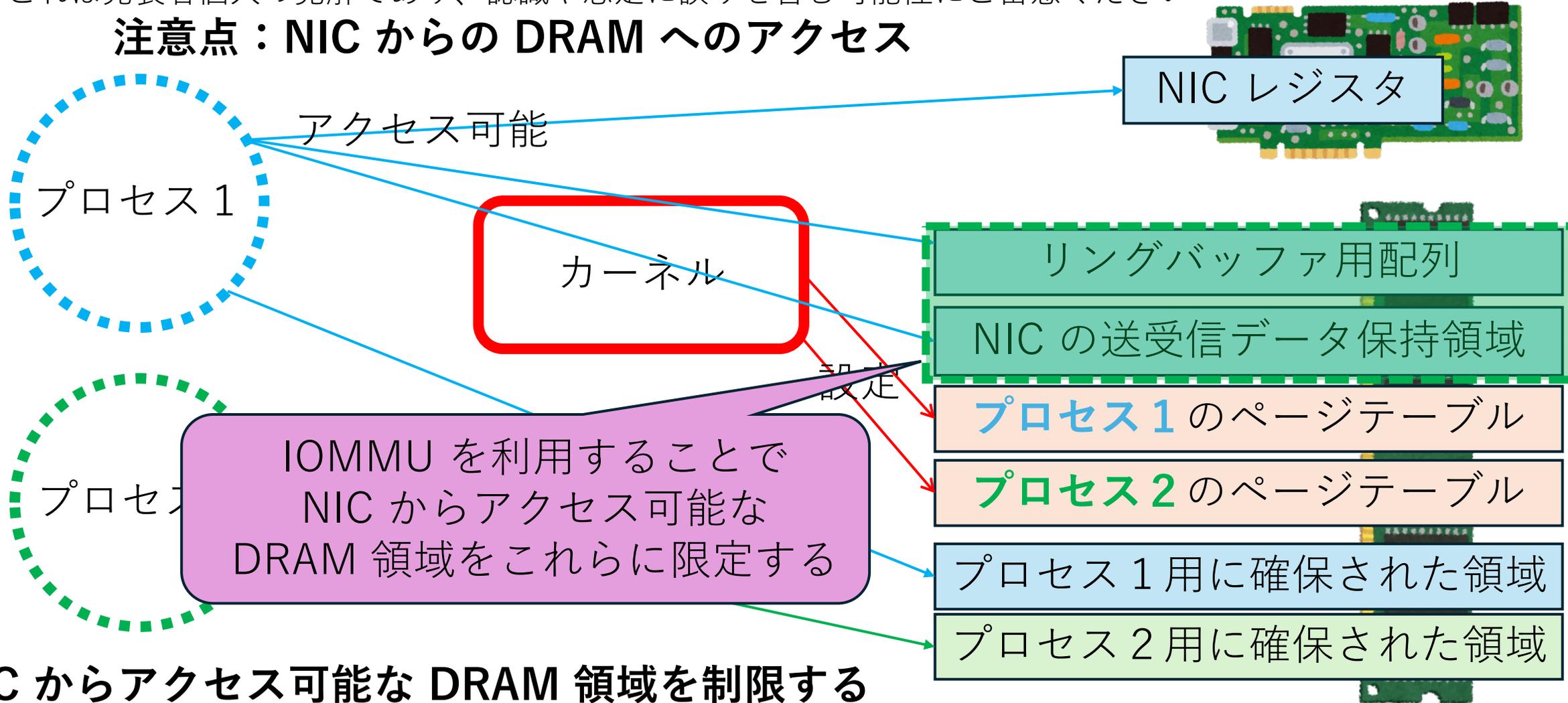


NIC からアクセス可能な DRAM 領域を制限する

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

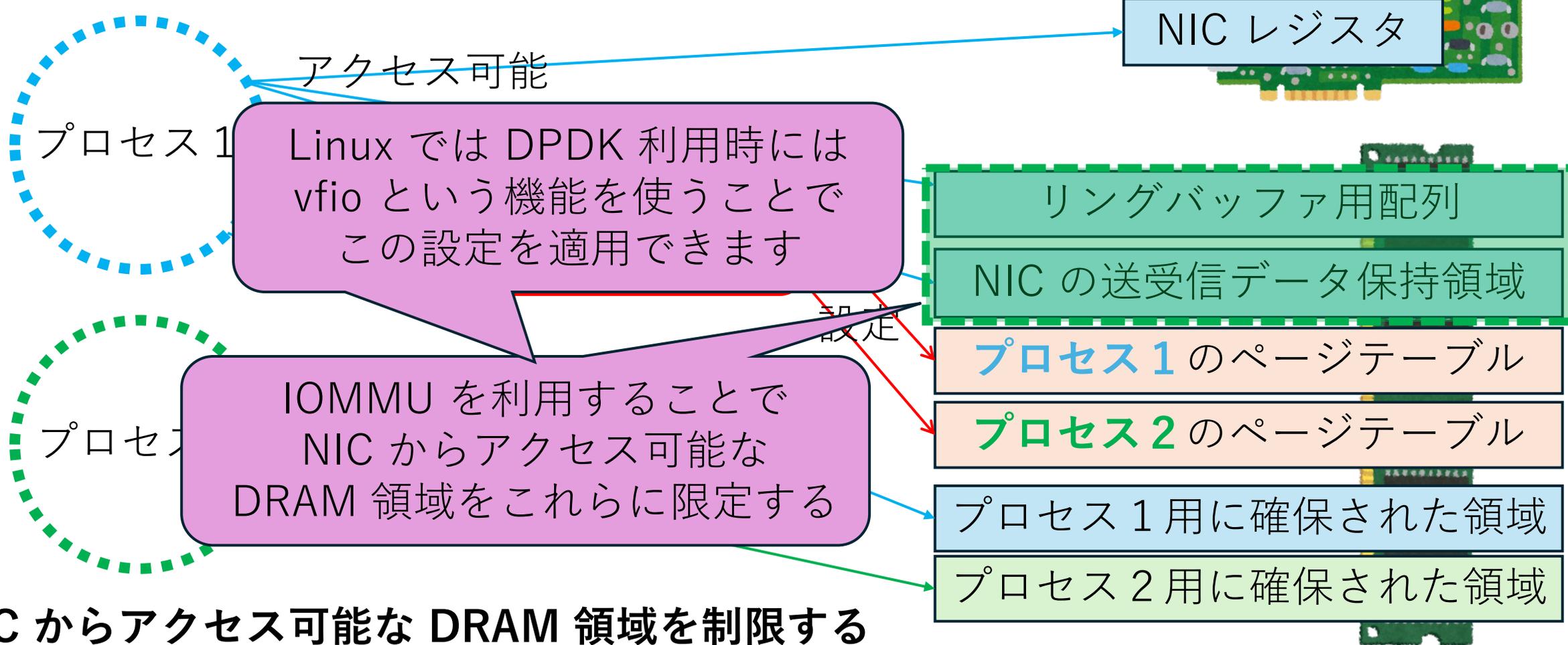
注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

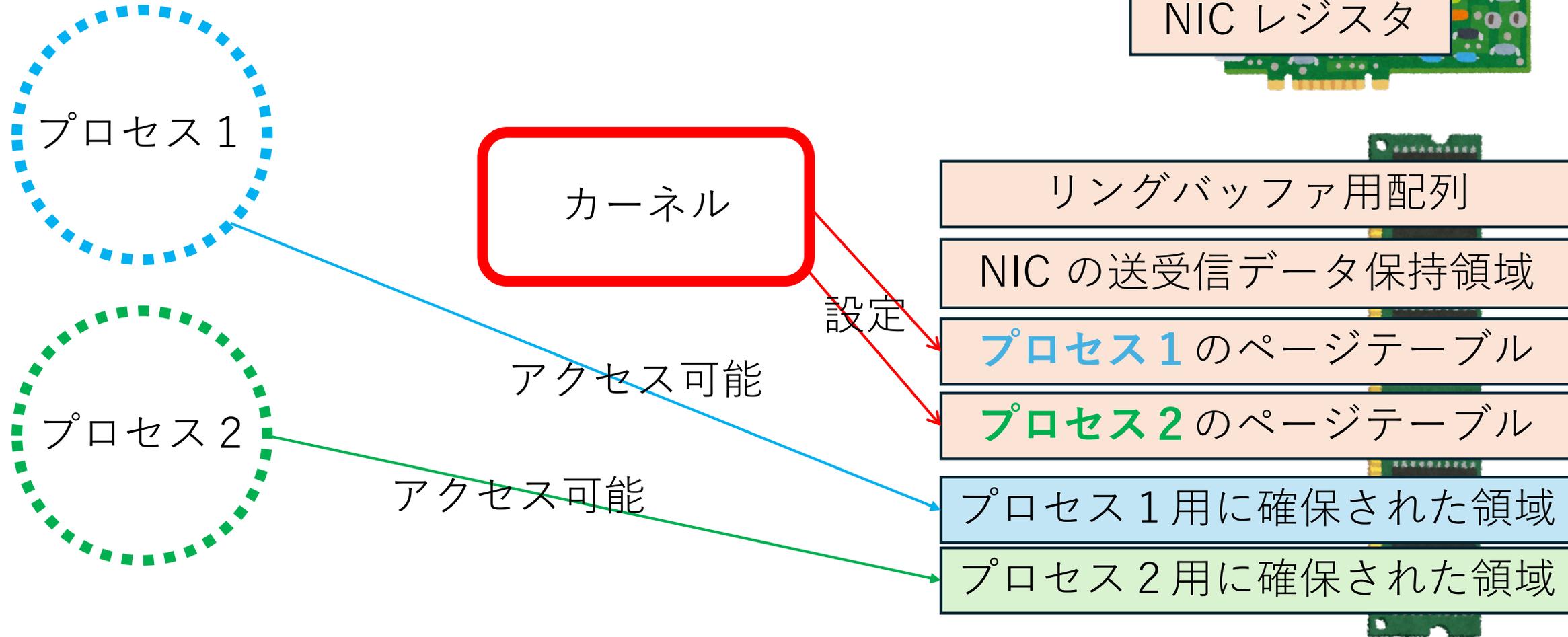
注意点：NIC からの DRAM へのアクセス



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

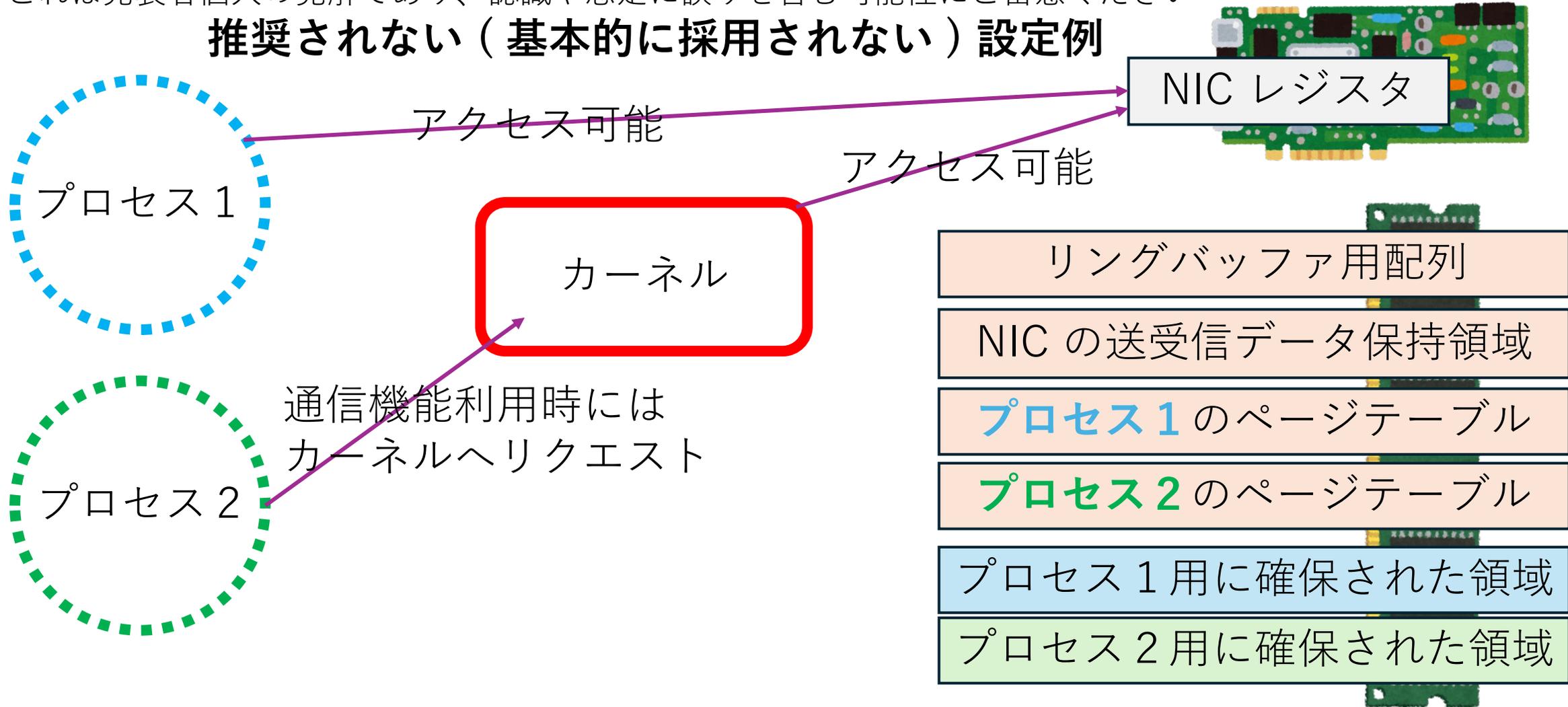
推奨されない（基本的に採用されない）設定例



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

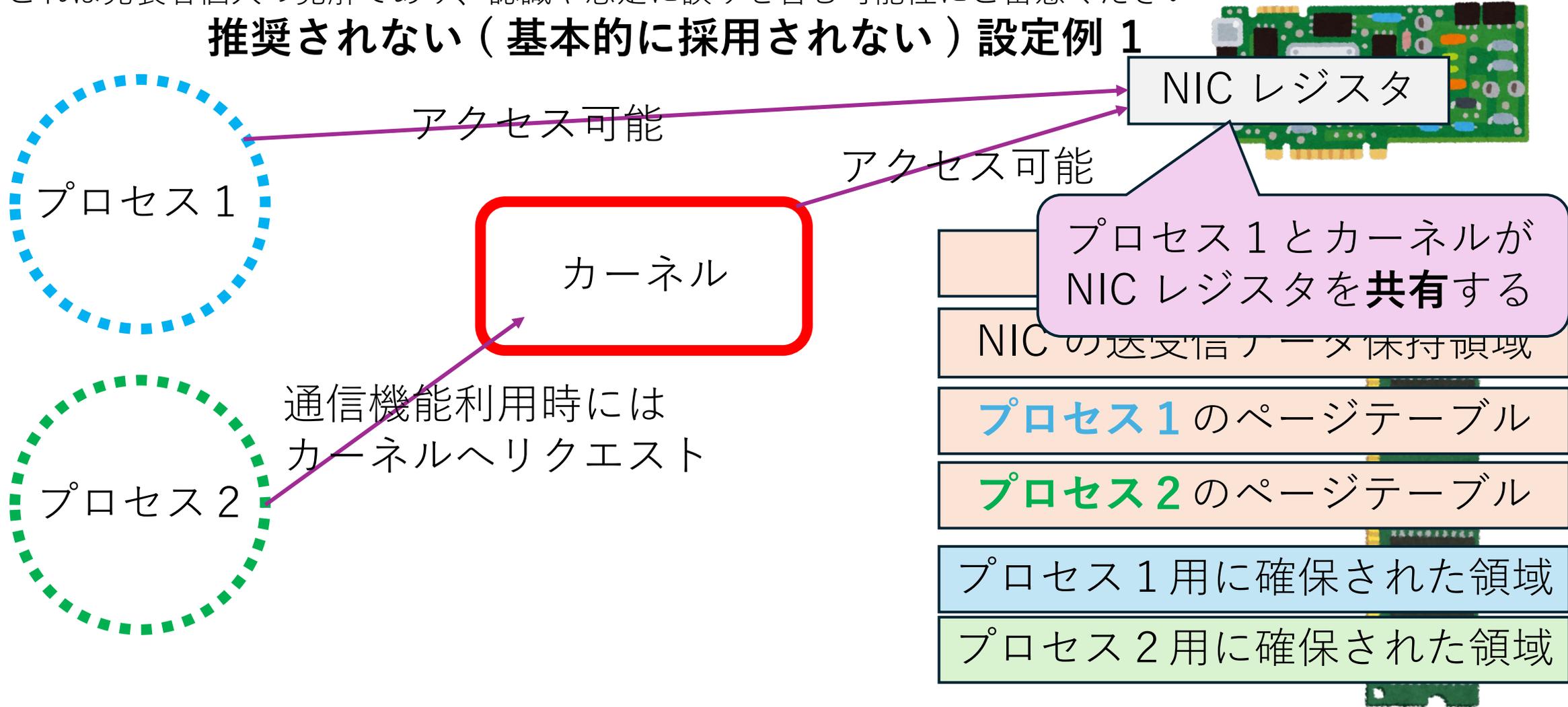
推奨されない（基本的に採用されない）設定例



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

推奨されない（基本的に採用されない）設定例 1



セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

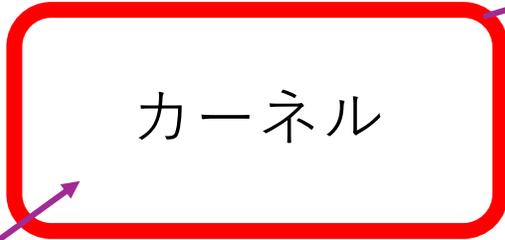
推奨されない（基本的に採用されない）設定例 1



NIC レジスタ

アクセス可能

アクセス可能



プロセス 1 とカーネルが
NIC レジスタを**共有**する

NIC の送受信データ保持領域

プロセス 1 のページテーブル

プロセス 2 のページテーブル

プロセス 1 用に確保された領域

プロセス 2 用に確保された領域

何故推奨されないか？

プロセス 1 上のアプリを
信頼しない想定だから

利用時には
ヘルプクエスト

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

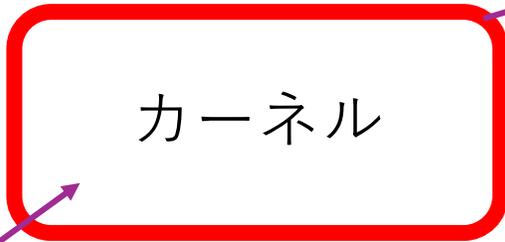
推奨されない（基本的に採用されない）設定例 1



アクセス可能

NIC レジスタ

アクセス可能

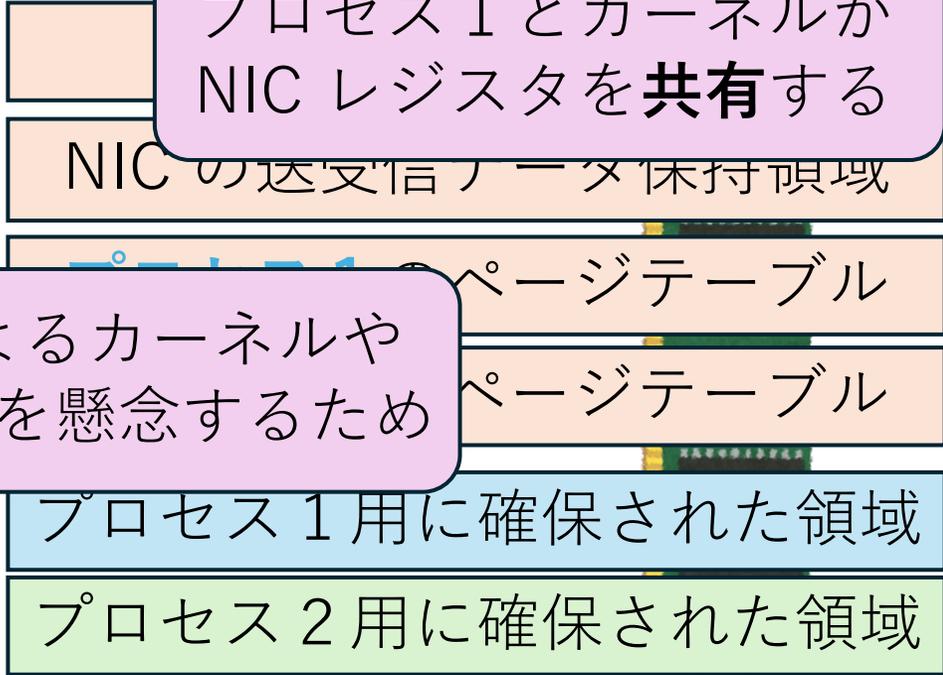


プロセス 1 とカーネルが NIC レジスタを **共有**する

何故推奨されないか？
プロセス 1 上のアプリを信頼しない想定だから

利用時には

バグや脆弱性によるカーネルや他の要素への波及を懸念するため



セキュリティについての考え方

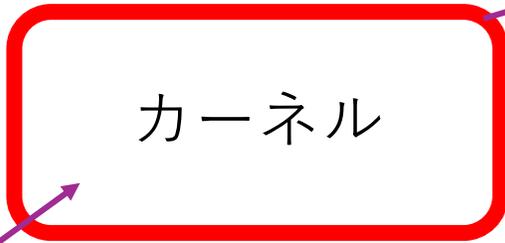
これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

推奨されない（基本的に採用されない）設定例 1



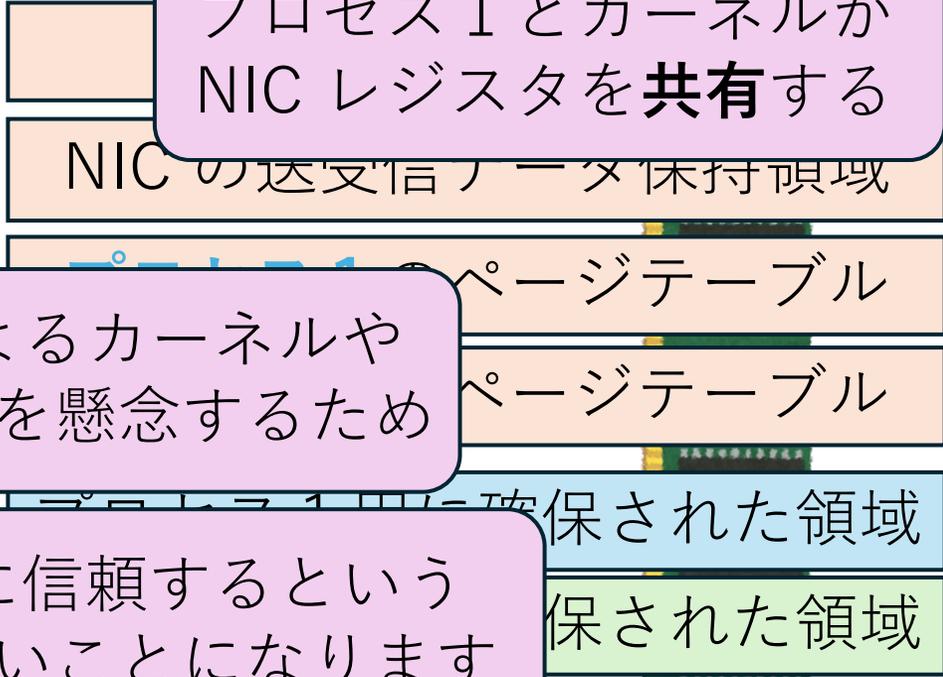
アクセス可能

アクセス可能



NIC レジスタ

プロセス 1 とカーネルが NIC レジスタを共有する



何故推奨されないか？
プロセス 1 上のアプリを信頼しない想定だから

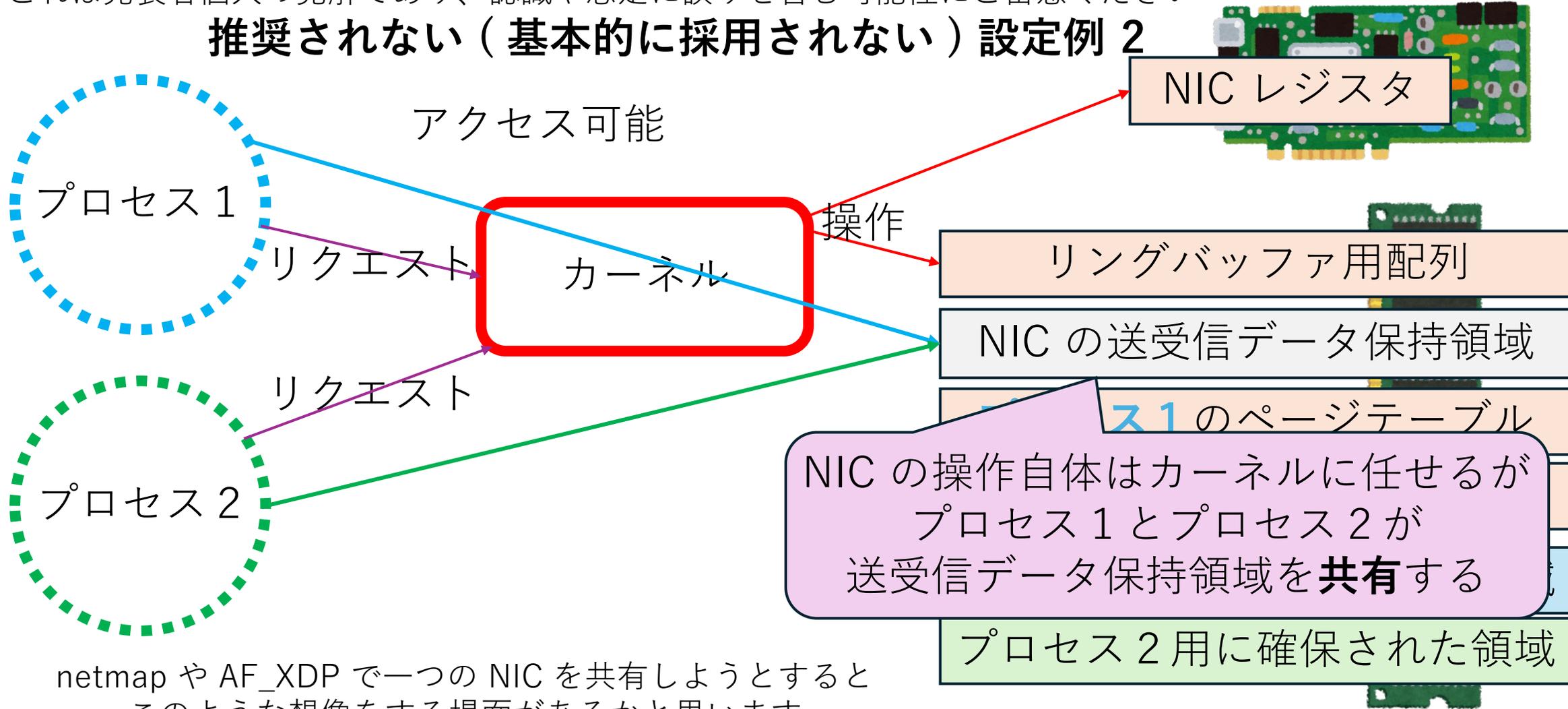
利用時には
バグや脆弱性によるカーネルや他の要素への波及を懸念するため

プロセス 1 上のアプリをカーネルと同様に信頼するという想定を適用するならこの設定には問題がないこととなります

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

推奨されない（基本的に採用されない）設定例 2

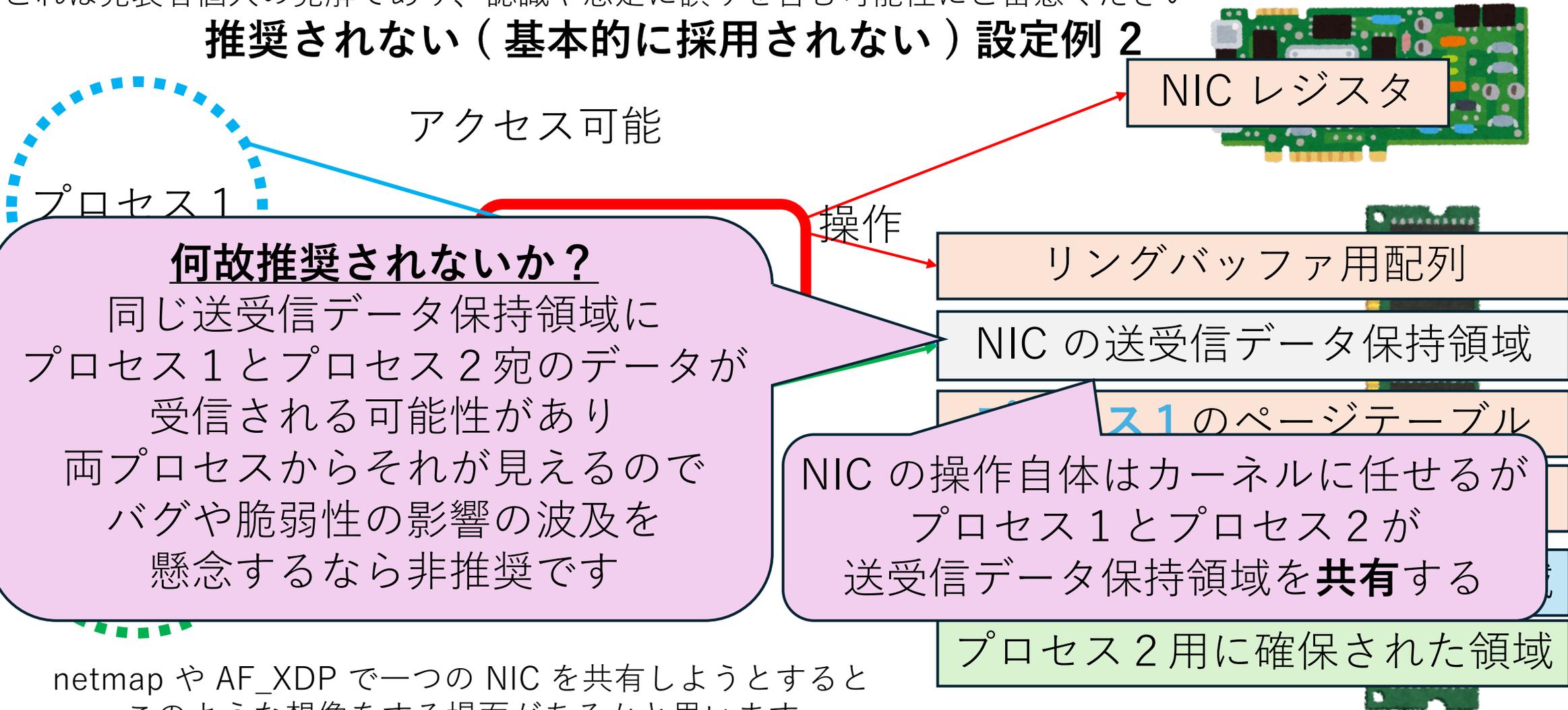


netmap や AF_XDP で一つの NIC を共有しようとするとき
このような想像をする場面があるかと思います

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

推奨されない（基本的に採用されない）設定例 2



プロセス 1

アクセス可能

操作

NIC レジスタ

リングバッファ用配列

NIC の送受信データ保持領域

プロセス 1 のページテーブル

NIC の操作自体はカーネルに任せるが
プロセス 1 とプロセス 2 が
送受信データ保持領域を**共有**する

プロセス 2 用に確保された領域

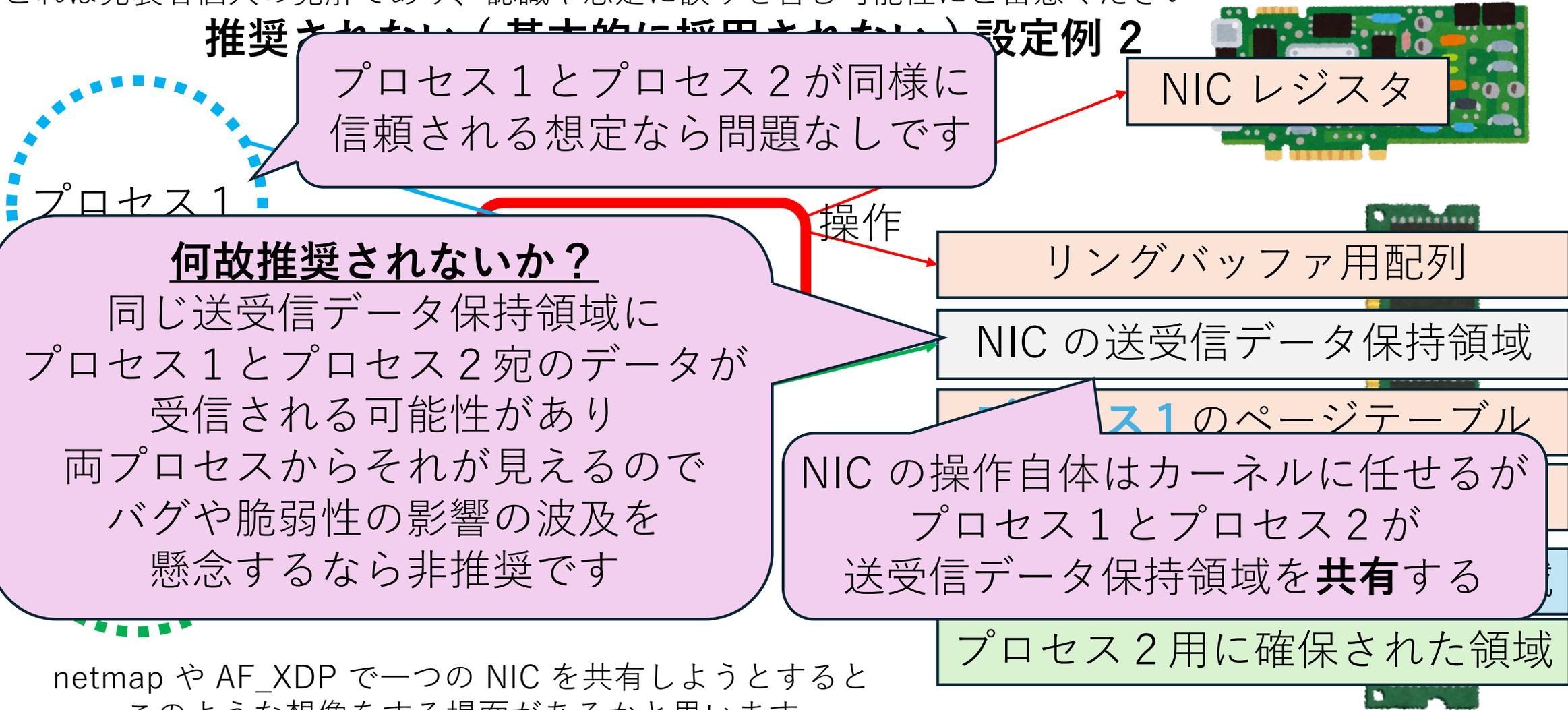
何故推奨されないか？
同じ送受信データ保持領域に
プロセス 1 とプロセス 2 宛のデータが
受信される可能性があり
両プロセスからそれが見えるので
バグや脆弱性の影響の波及を
懸念するなら非推奨です

netmap や AF_XDP で一つの NIC を共有しようとするとき
このような想像をする場面があるかと思います

セキュリティについての考え方

これは発表者個人の見解であり、認識や想定に誤りを含む可能性にご留意ください

推奨されない（基本的に採用されない）設定例 2



プロセス 1 とプロセス 2 が同様に信頼される想定なら問題なしです

NIC レジスタ

プロセス 1

何故推奨されないか？

同じ送受信データ保持領域にプロセス 1 とプロセス 2 宛のデータが受信される可能性があり
両プロセスからそれが見えるので
バグや脆弱性の影響の波及を懸念するなら非推奨です

操作

リングバッファ用配列

NIC の送受信データ保持領域

プロセス 1 のページテーブル

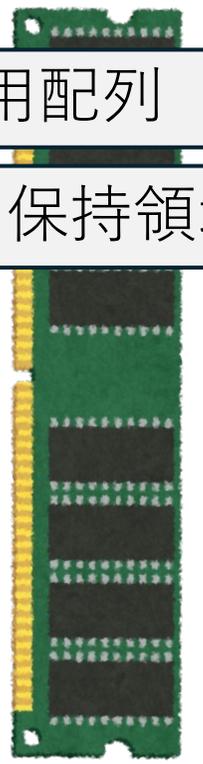
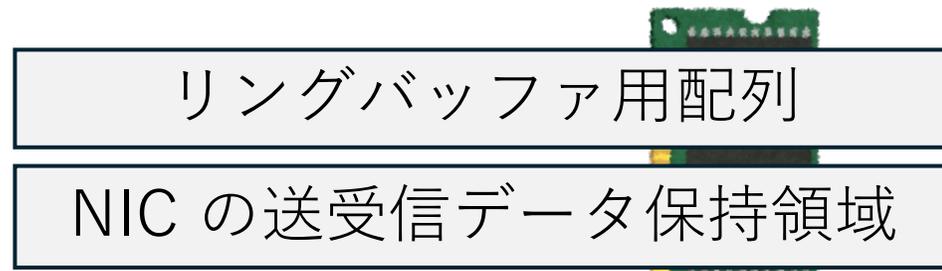
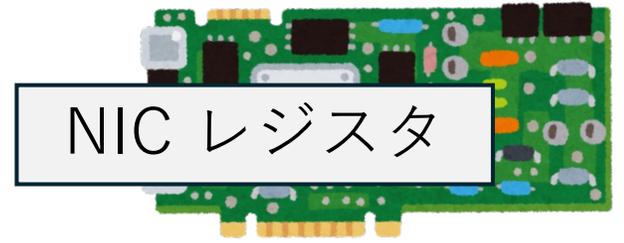
NIC の操作自体はカーネルに任せるが
プロセス 1 とプロセス 2 が送受信データ保持領域を共有する

プロセス 2 用に確保された領域

netmap や AF_XDP で一つの NIC を共有しようとするとき
このような想像をする場面があるかと思います

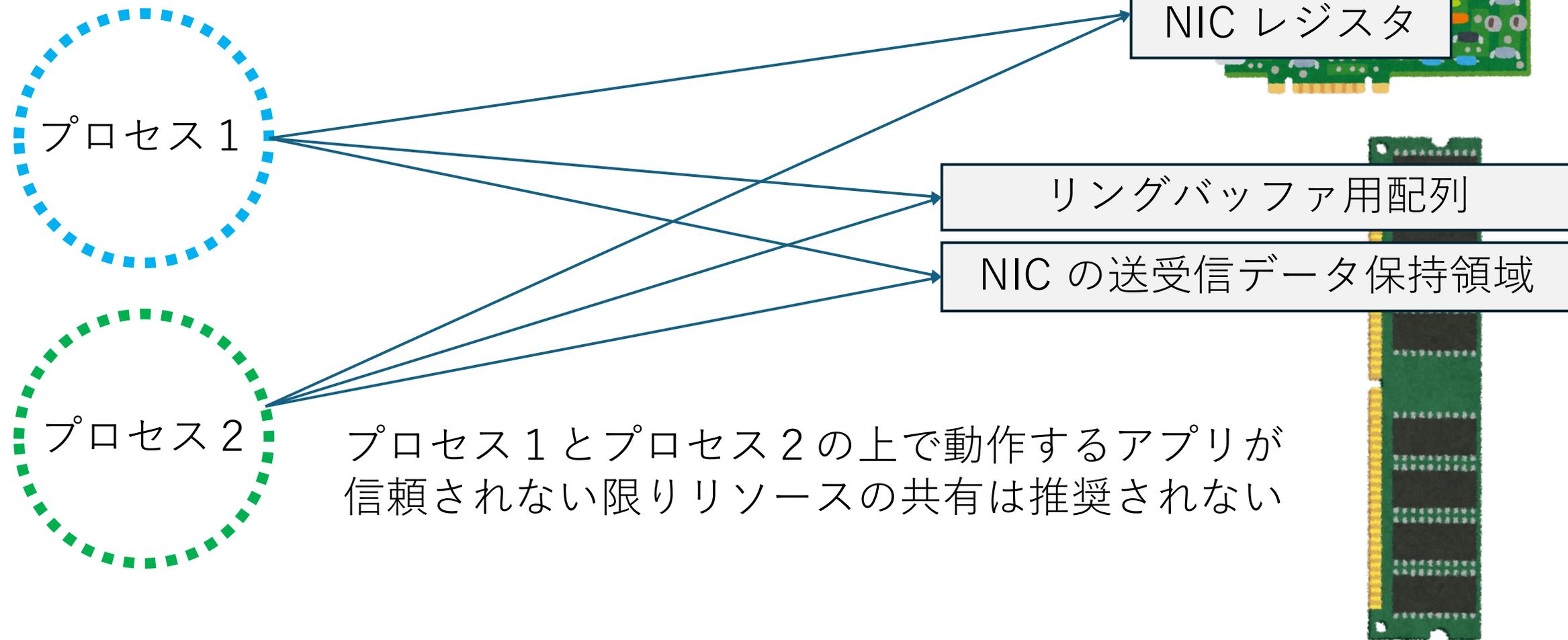
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



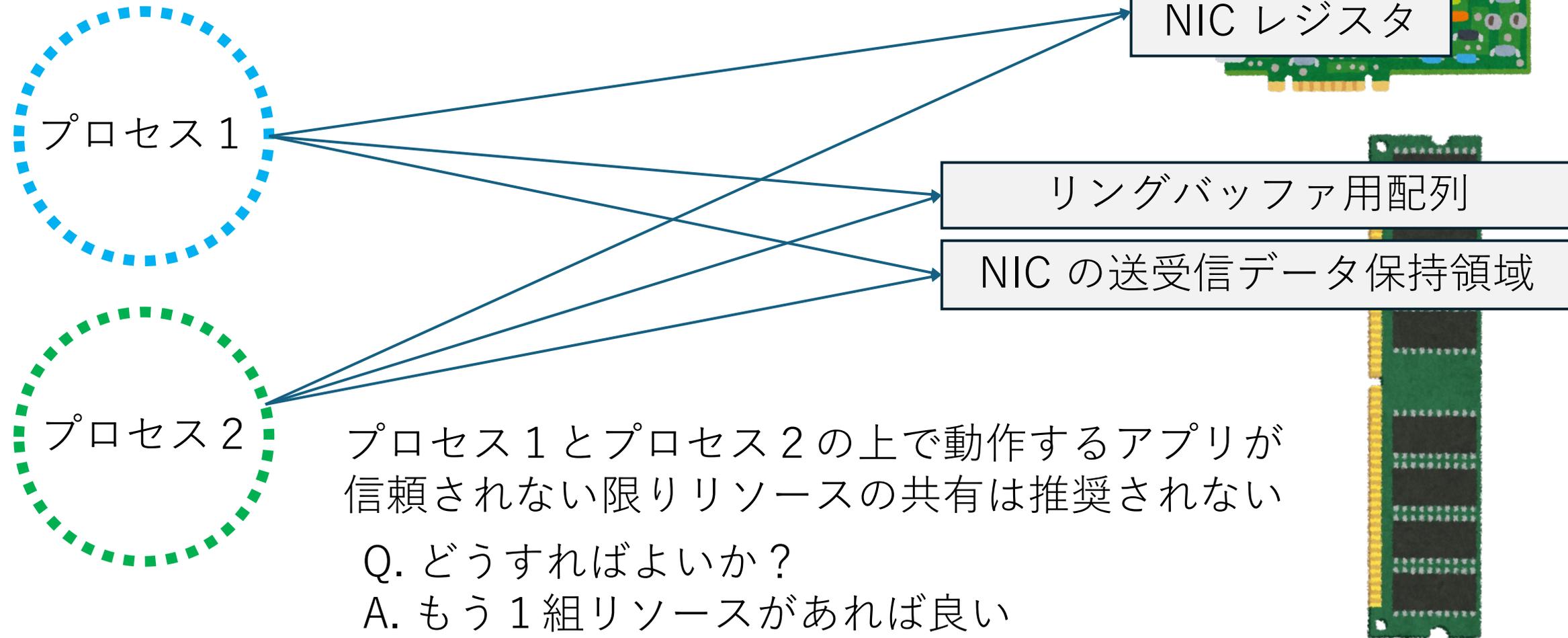
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



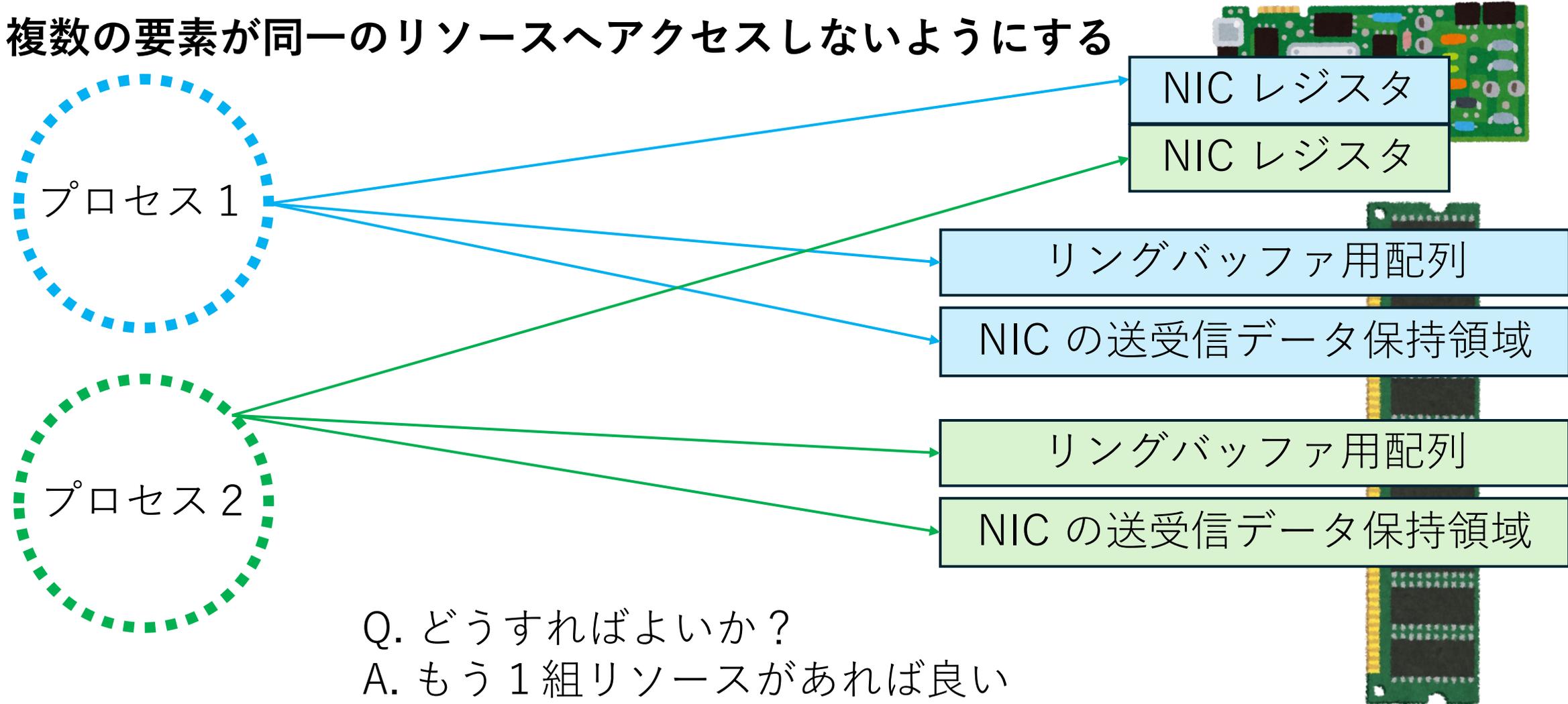
プロセス 1 とプロセス 2 の上で動作するアプリが信頼されない限りリソースの共有は推奨されない

Q. どうすればよいか？

A. もう 1 組リソースがあれば良い

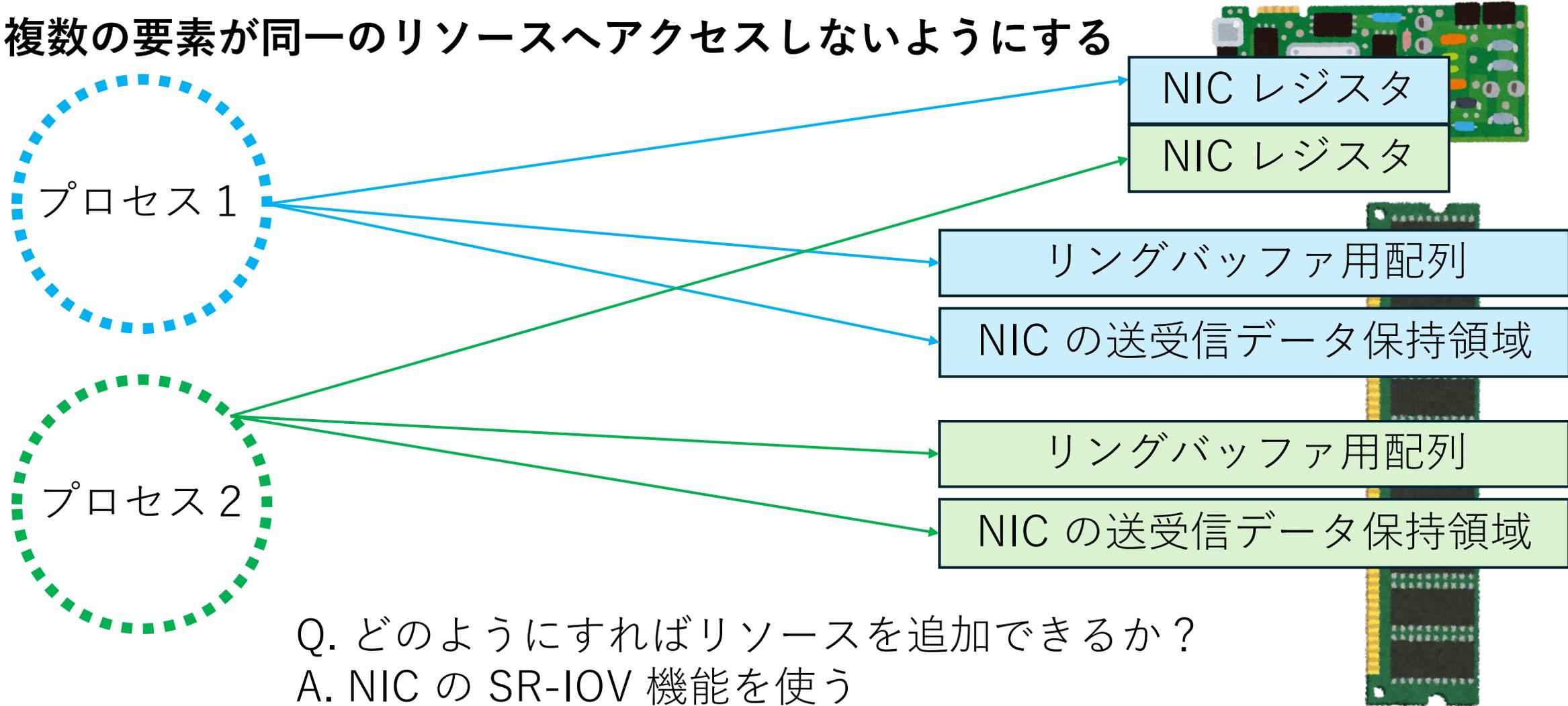
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



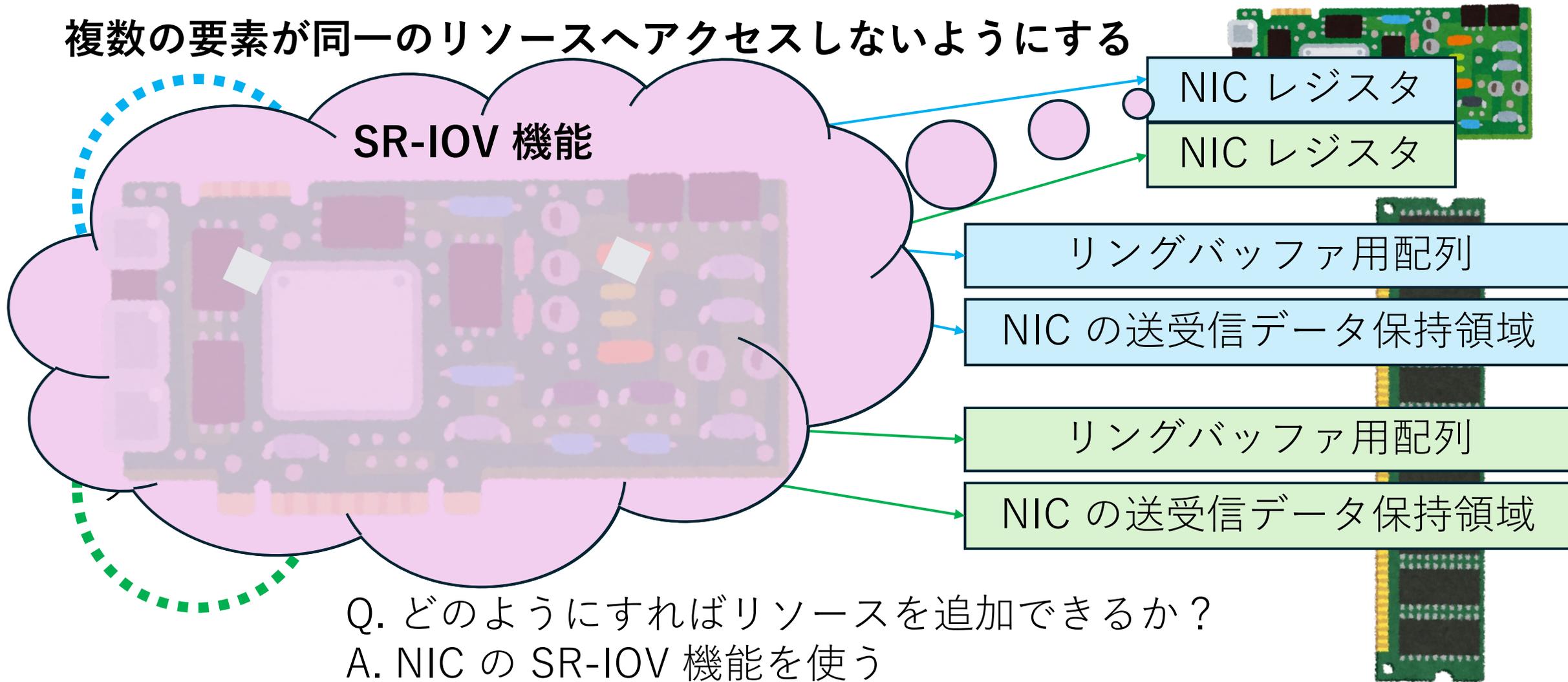
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



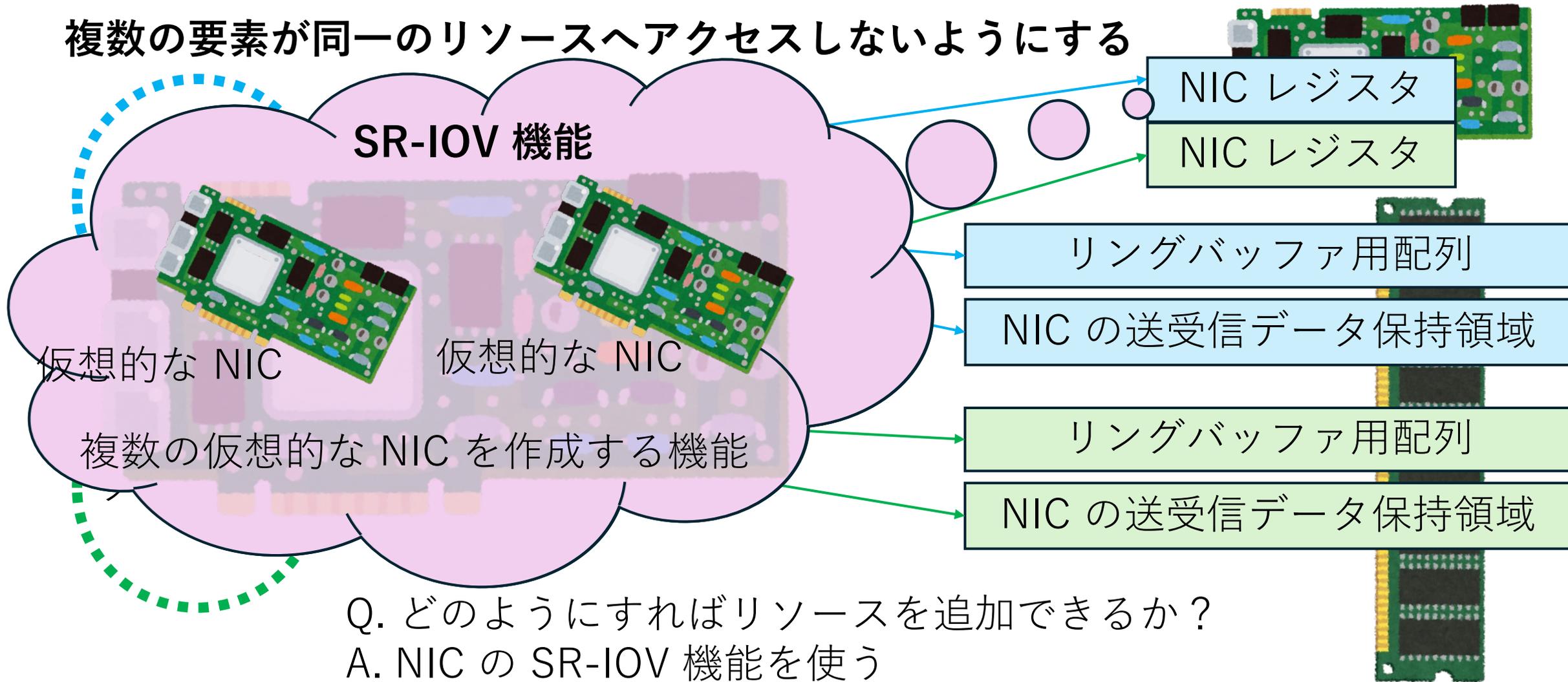
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



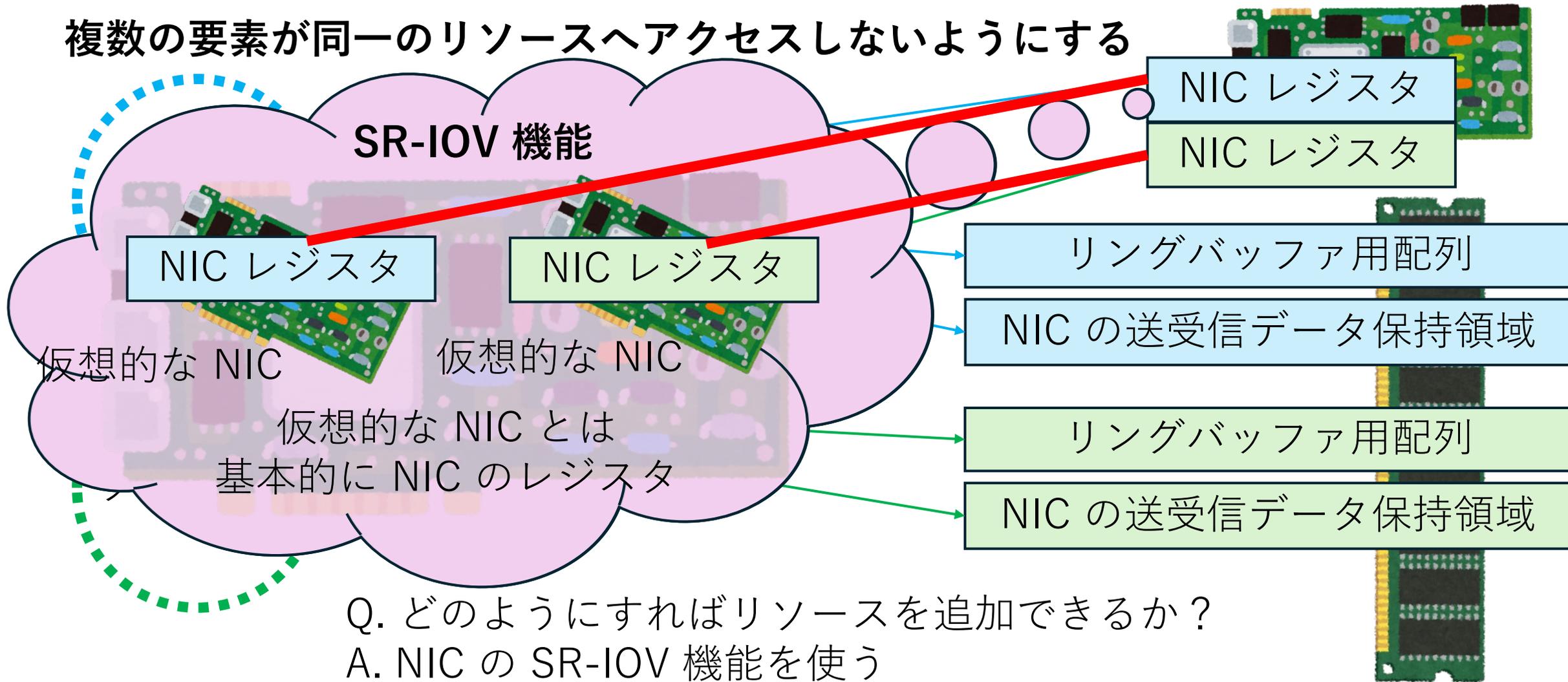
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



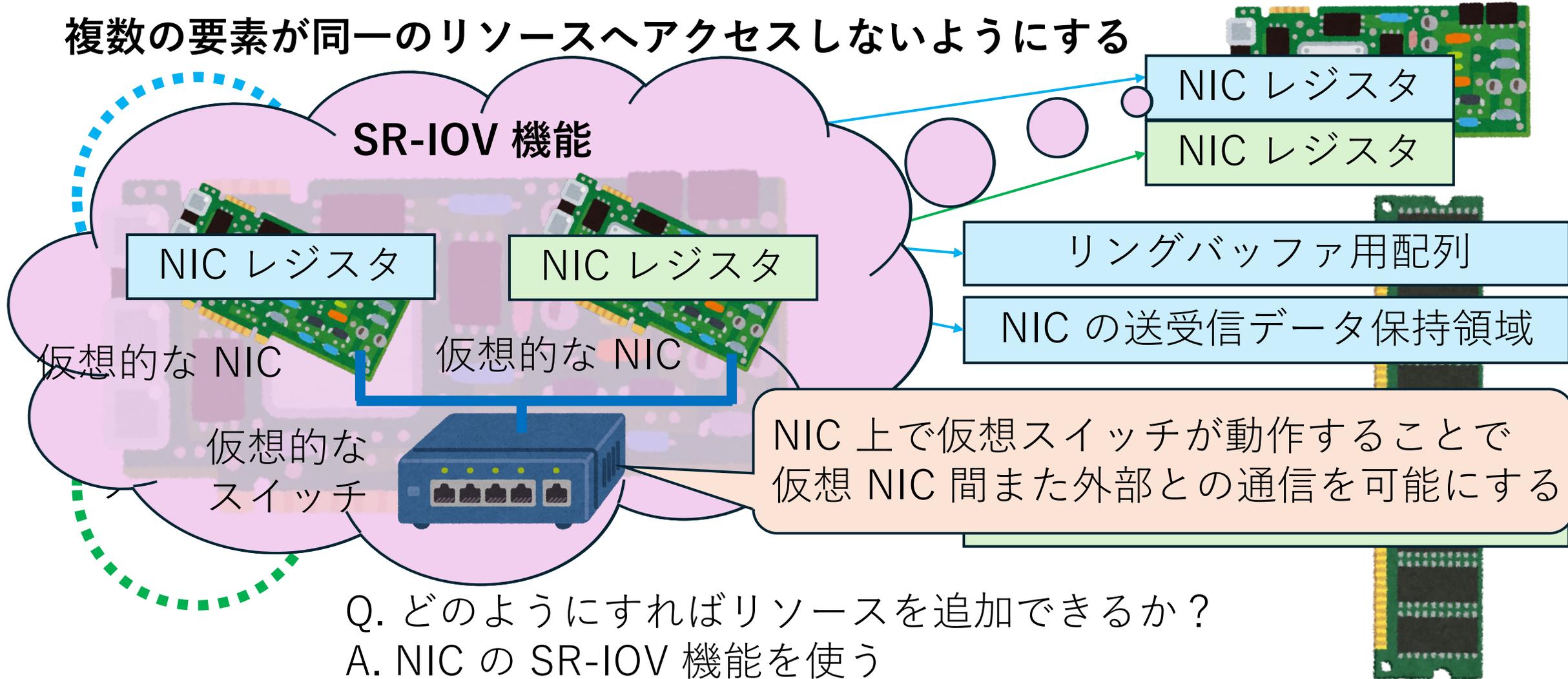
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



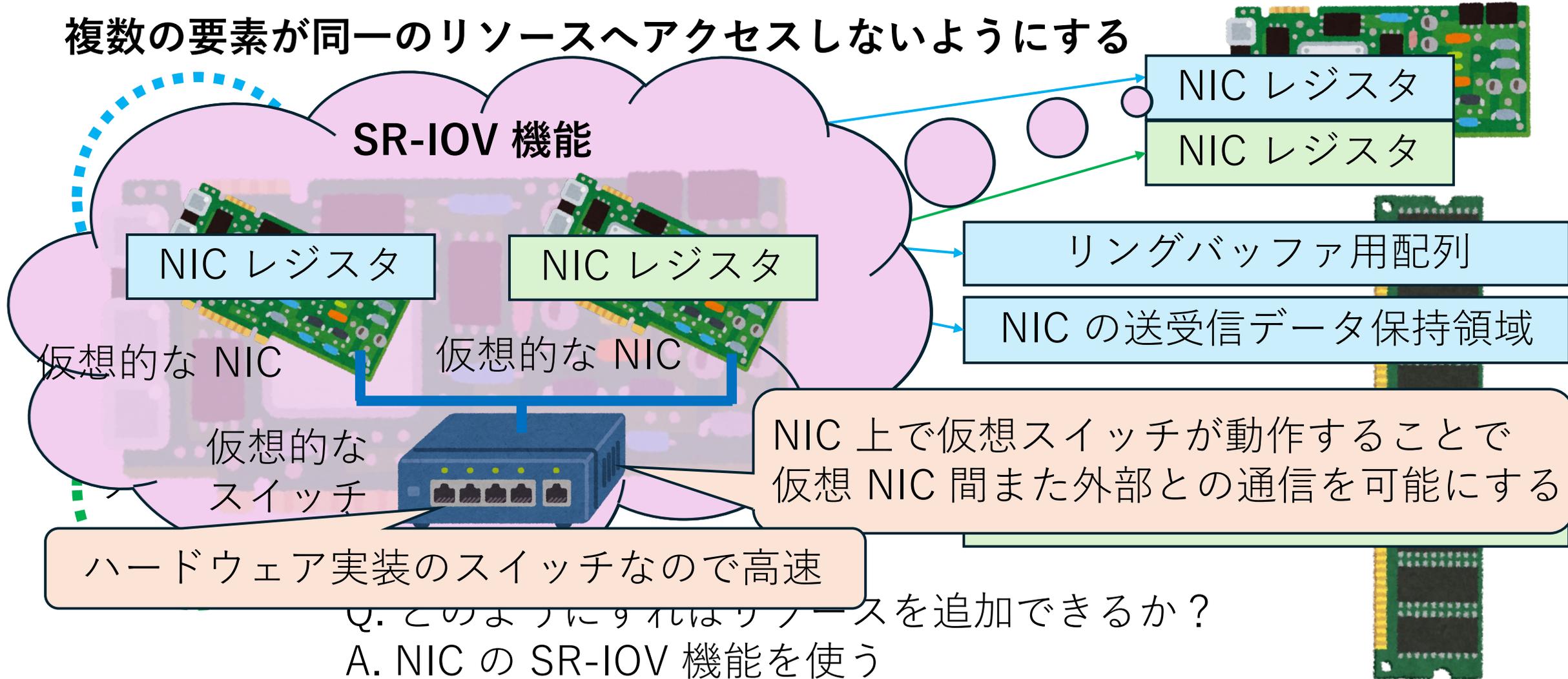
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする

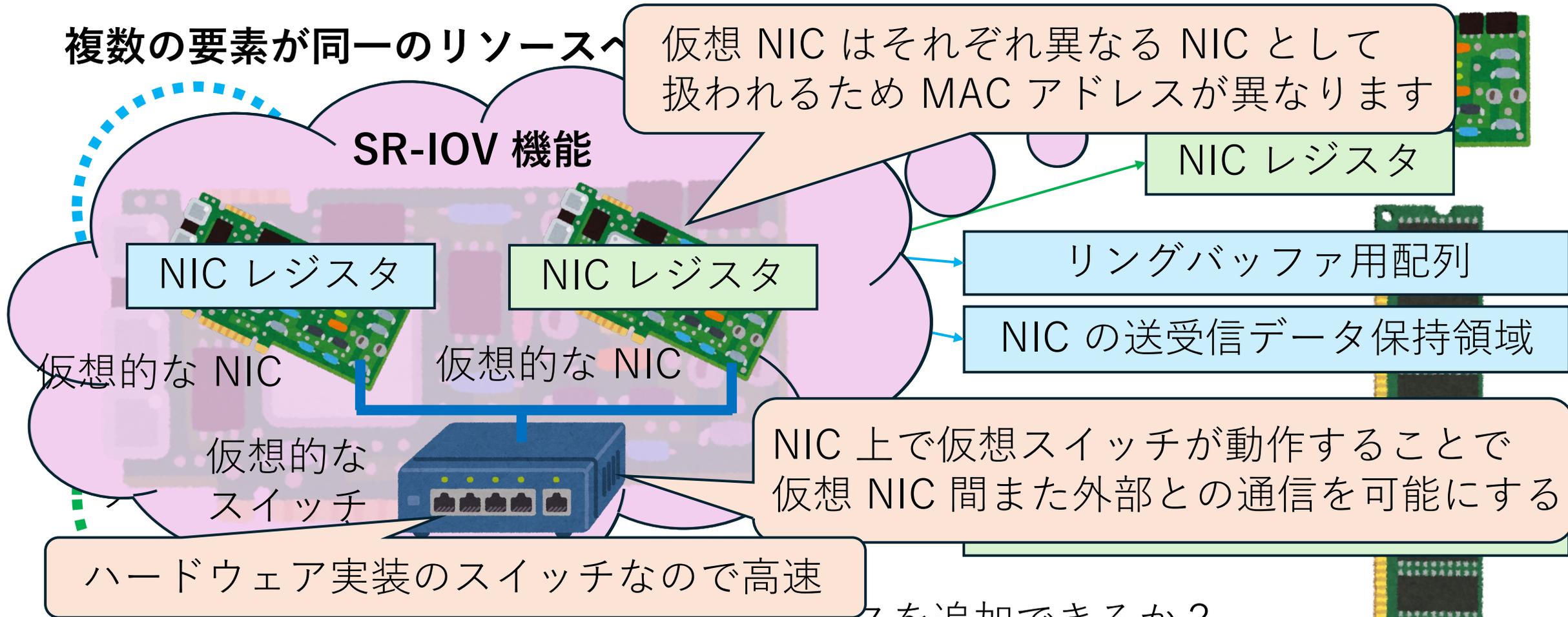


一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



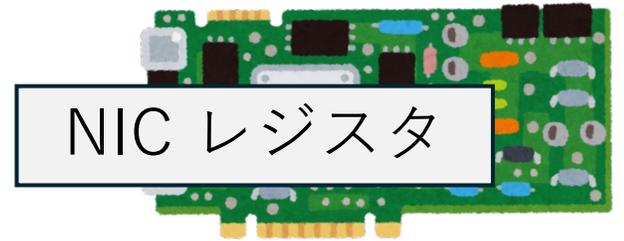
一つの NIC を共有する方法



Q. このようにすれば仮想 NIC を追加できるか？
A. NIC の SR-IOV 機能を使う

一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

一つの NIC を共有する方法

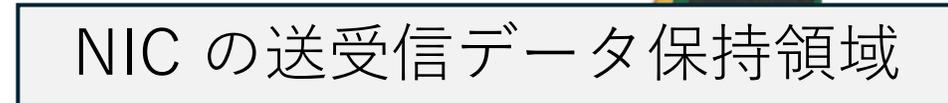
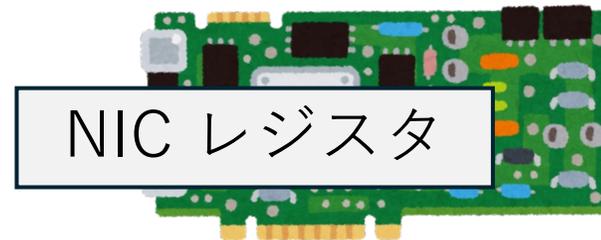
複数の要素が同一のリソースへアクセスしないようにする



仮想 NIC



仮想 NIC



Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

一つの NIC を共有する方法

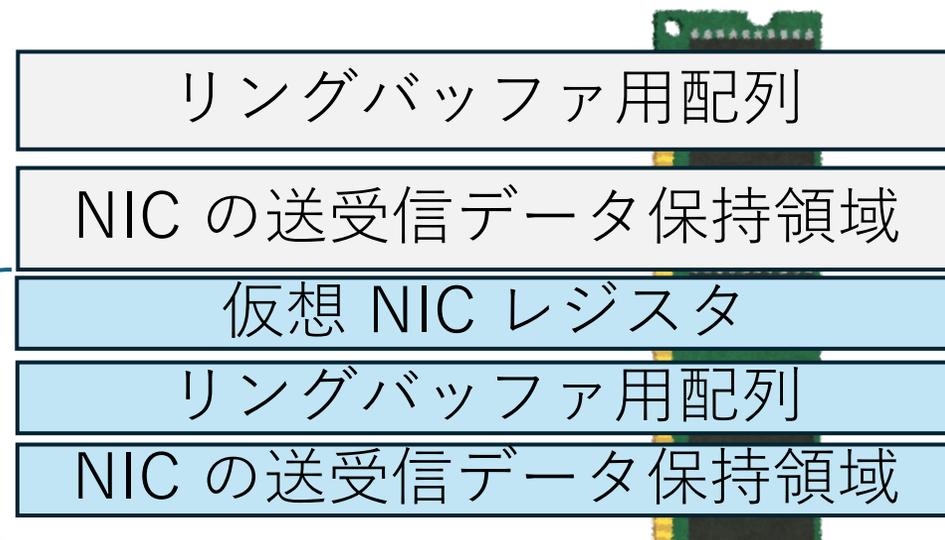
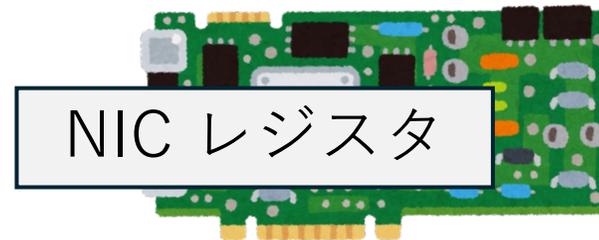
複数の要素が同一のリソースへアクセスしないようにする



仮想 NIC



仮想 NIC



Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

一つの NIC を共有する方法

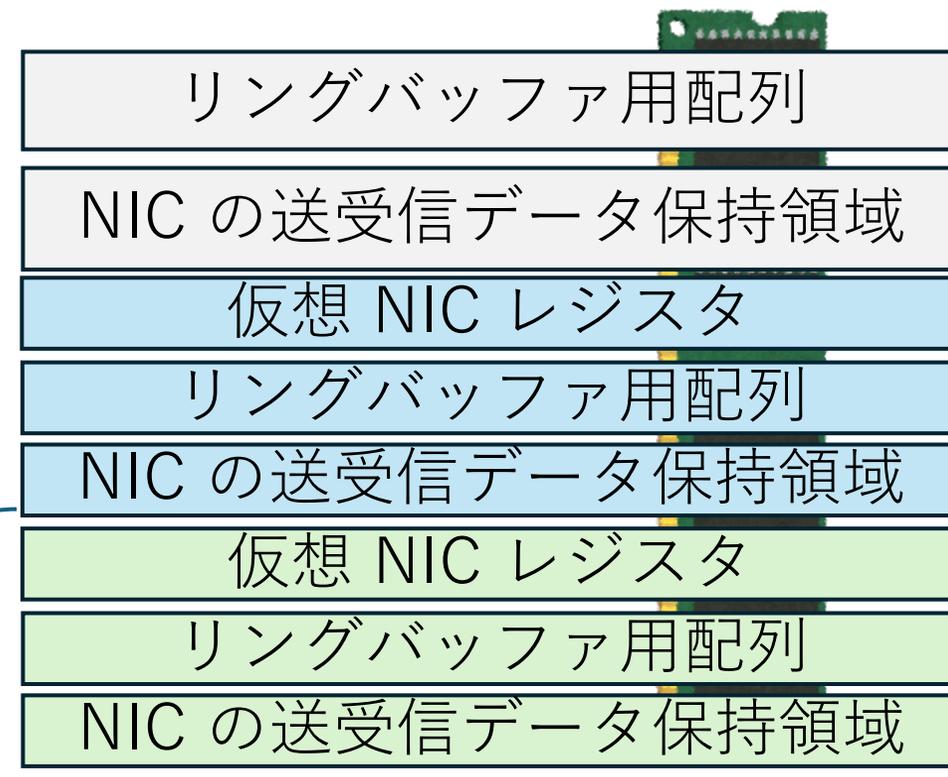
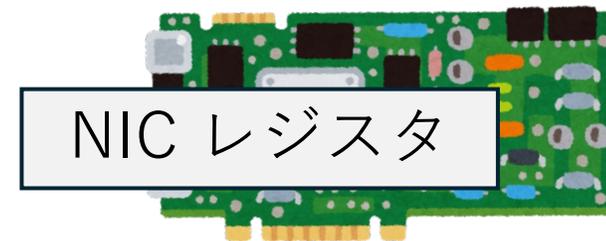
複数の要素が同一のリソースへアクセスしないようにする



仮想 NIC



仮想 NIC

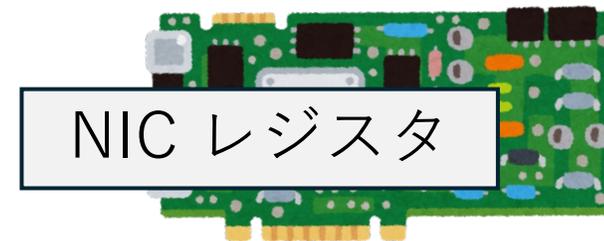


Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする

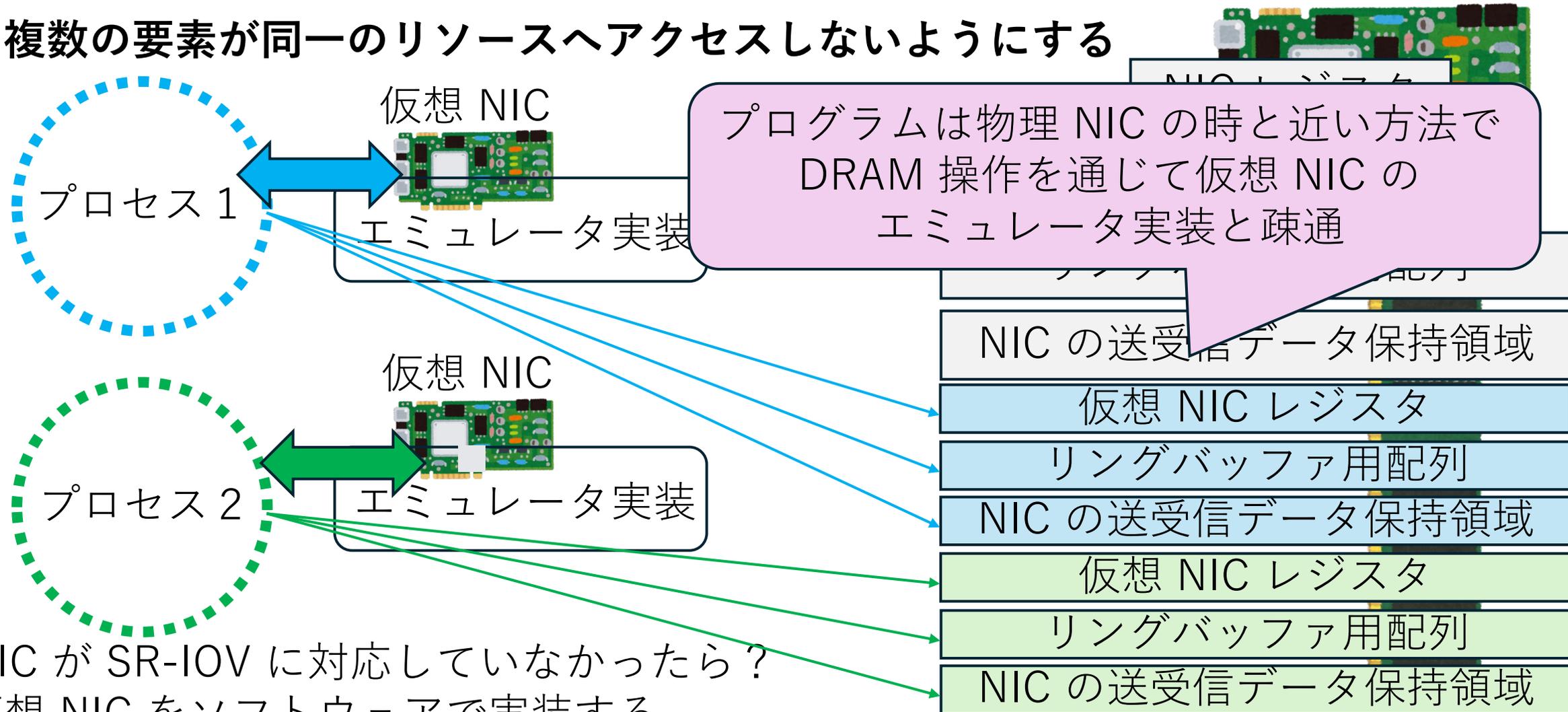


Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする

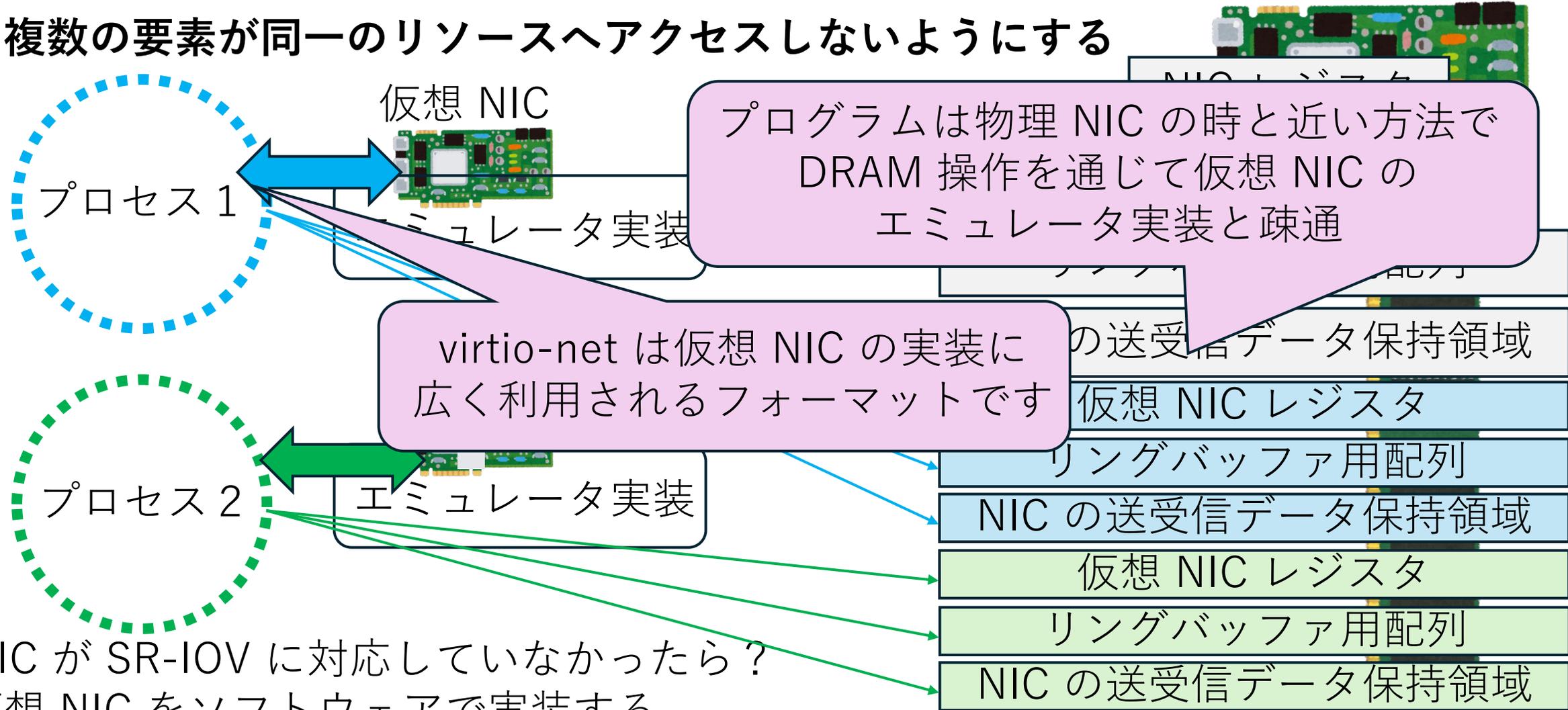


Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする

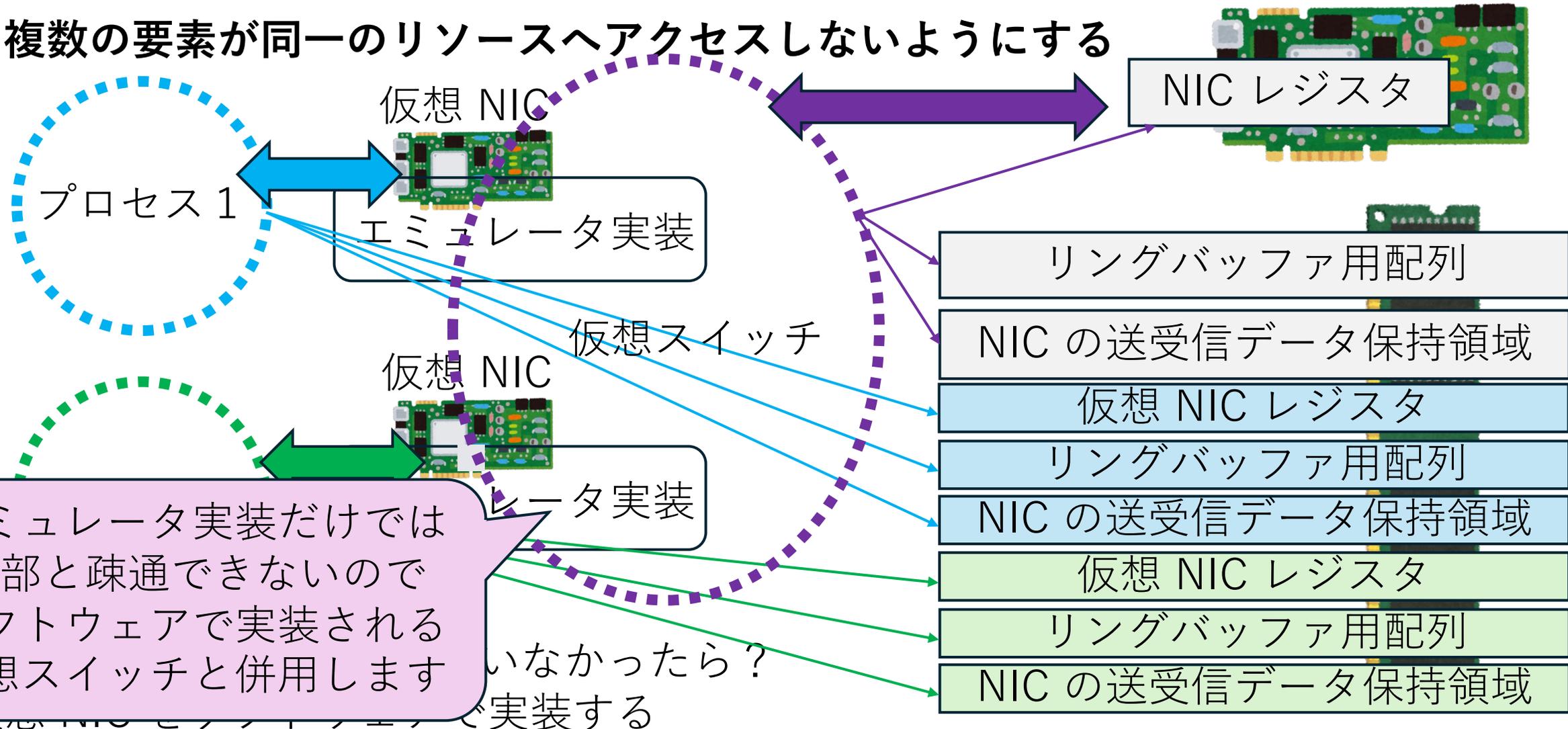


Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

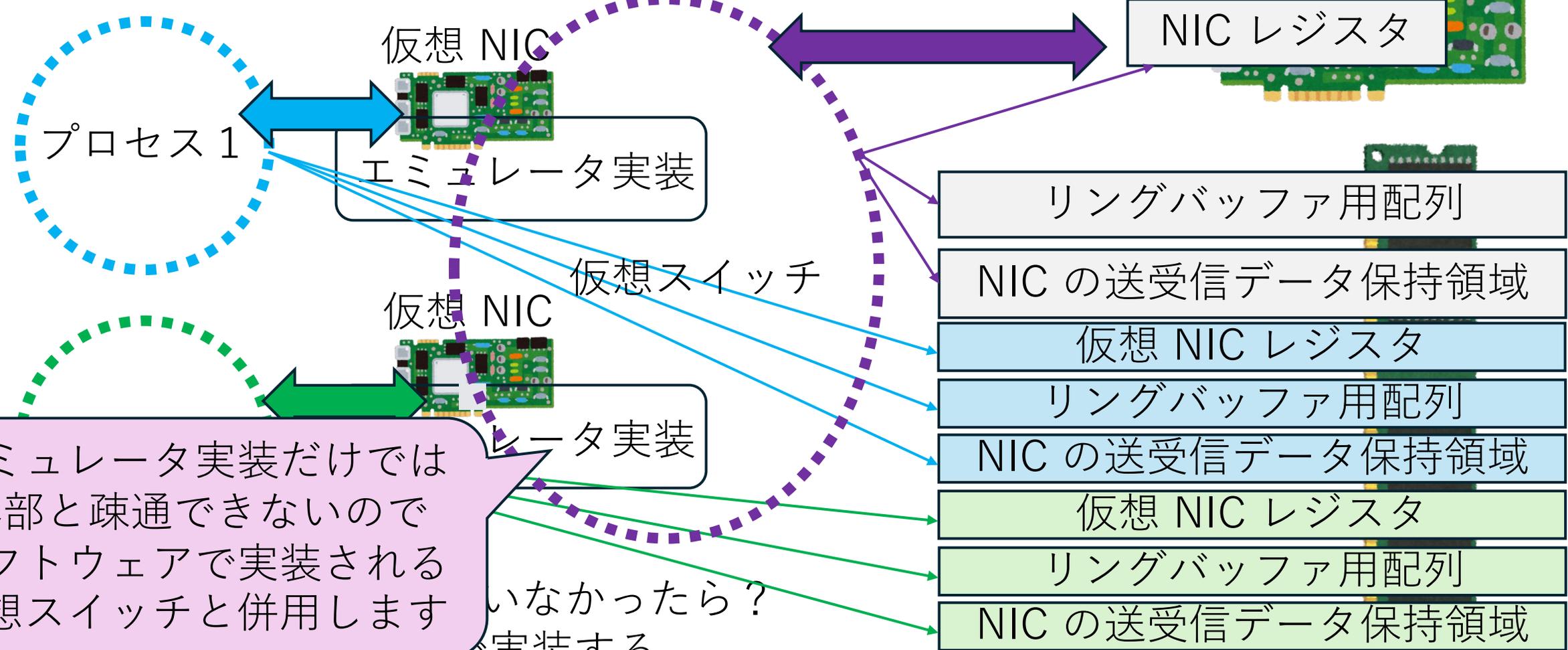
一つの NIC を共有する方法

複数の要素が同一のリソースへアクセスしないようにする



DPDK を使って仮想スイッチを実装する場合の例

複数の要素が同一のリソースへアクセスしないようにする

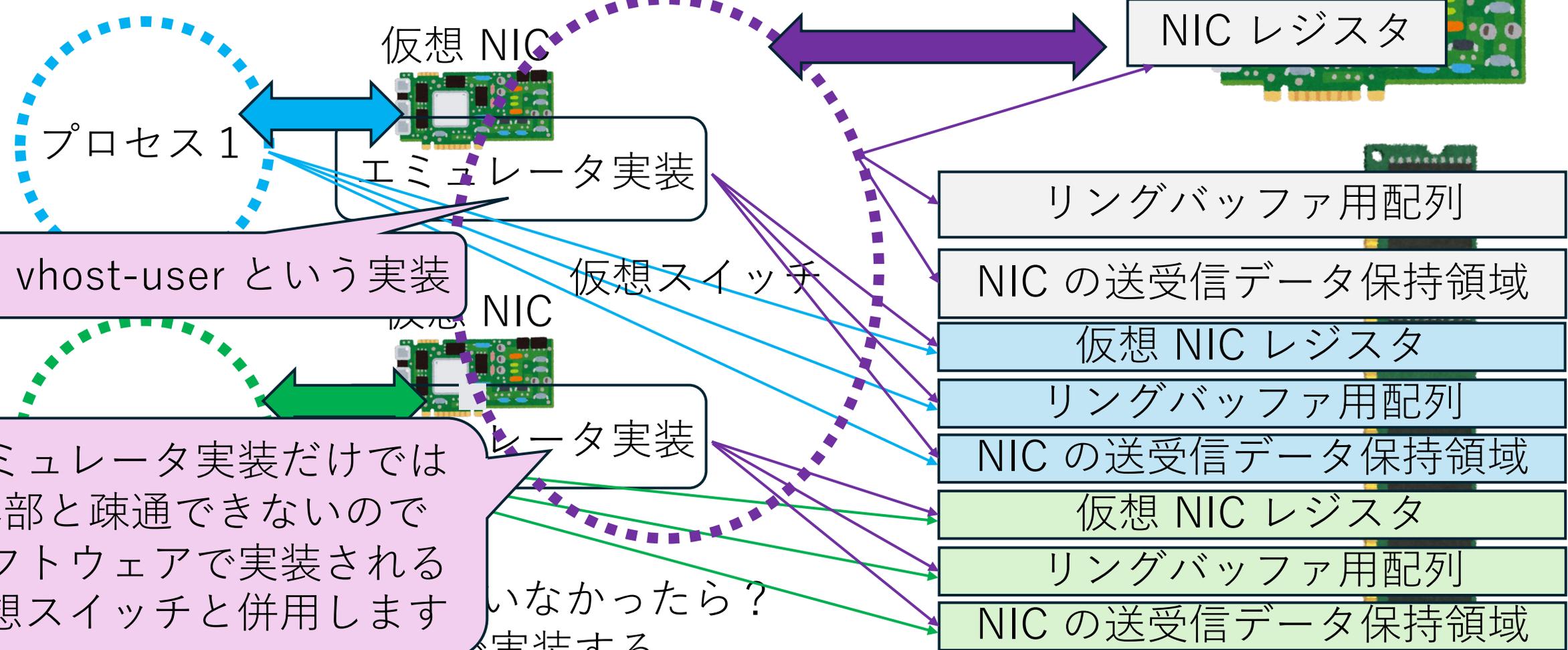


エミュレータ実装だけでは外部と疎通できないのでソフトウェアで実装される仮想スイッチと併用します

いなかったら？
で実装する

DPDK を使って仮想スイッチを実装する場合の例

複数の要素が同一のリソースへアクセスしないようにする



vhost-user という実装

エミュレータ実装だけでは外部と疎通できないのでソフトウェアで実装される仮想スイッチと併用します

いなかったら？

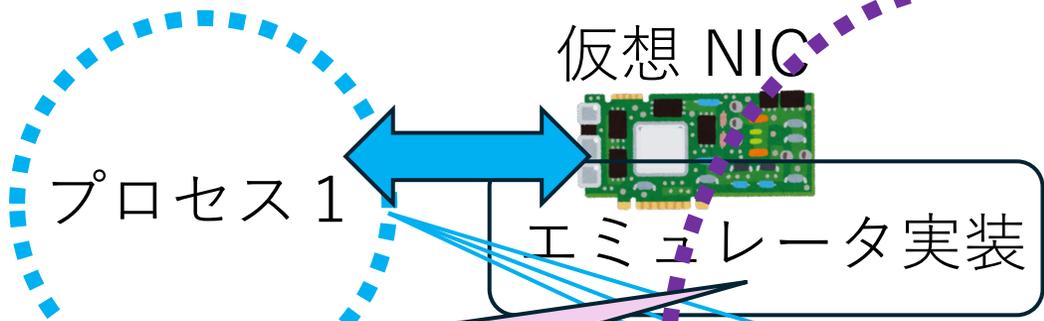
で実装する

DPDK を使って仮想スイッチを実装する場合の例

複数の要素が同一のリソースへアクセスしないようにする



NIC レジスタ



エミュレータ実装を動かす
仮想スイッチプロセスは
仮想NIC関連のDRAM上領域に
アクセスできる必要があります

- 仮想NICレジスタ
- リングバッファ用配列
- NICの送受信データ保持領域
- 仮想NICレジスタ
- リングバッファ用配列
- NICの送受信データ保持領域

vhost-user という実装

仮想スイッチ

エミュレータ実装だけでは
外部と疎通できないので
ソフトウェアで実装される
仮想スイッチと併用します



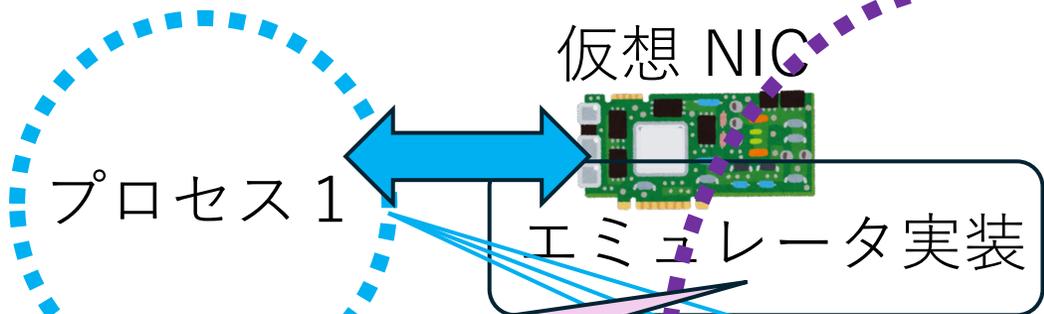
いなかったら？

で実装する

DPDK を使って仮想スイッチを実装する場合の例

複数の要素が同一のリソースへアクセスしないようにする

NIC レジスタ



エミュレータ実装を動かす
仮想スイッチプロセスは
仮想 NIC 関連の DRAM 上領域に
アクセスできる必要があります

vhost-user という実装

仮想スイッチ



エミュレータ実装だけでは
外部と疎通できないので
ソフトウェアで実装される
仮想スイッチと併用します

プロセス間の共有メモリで実装可能です

仮想 NIC レジスタ
リングバッファ用配列
NIC の送受信データ保持領域

いなかったら？

で実装する

一般的な OS での運用

ちなみに、敢えて異なるプロセスが
同じ物理アドレスを参照できるように
設定するのがプロセス間の**共有メモリ**です

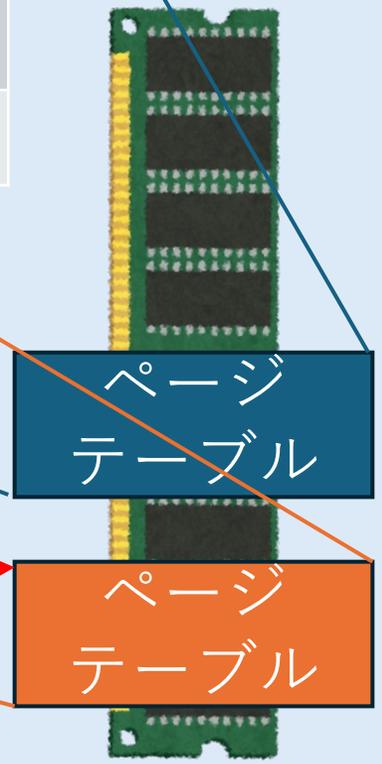
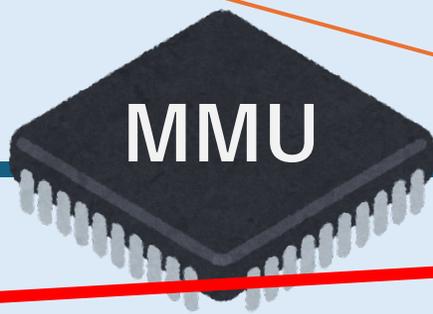
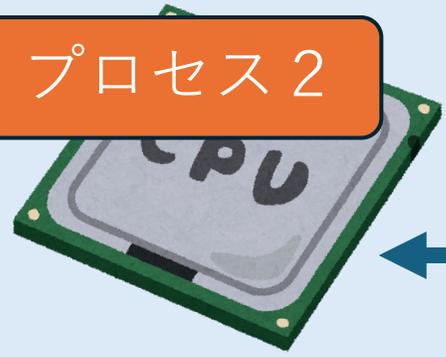
仮想	物理
0x0000	0x4000
0x1000	0x3000
...	仮想アドレスはプロセス1とプロセス2で 一致する必要はありません
...	

プロセス 2 用ページテーブル

仮想	物理
0x0000	
0x1000	0x2000
0x2000	0x4000
...	

プロセス 1 用ページテーブル

プロセス 2



この場合
物理メモリアドレス 0x4000 ~ 0x4fff が
プロセス 1 とプロセス 2 で共有されます

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

cr3: **仮想** アドレス 0x1000 へアクセス

一般的な OS での運用

ちなみに、敢えて異なるプロセスが同じ物理アドレスを参照できるように設定するのがプロセス間の**共有メモリ**です

仮想	物理
0x0000	0x4000
0x1000	0x3000
...	仮想アドレスはプロセス1とプロセス2で一致する必要はありません
...	

プロセス 2 用ページテーブル

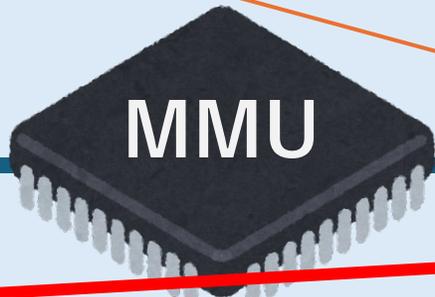
仮想	物理
0x0000	
0x1000	0x2000
0x2000	0x4000
...	

プロセス 1 用ページテーブル

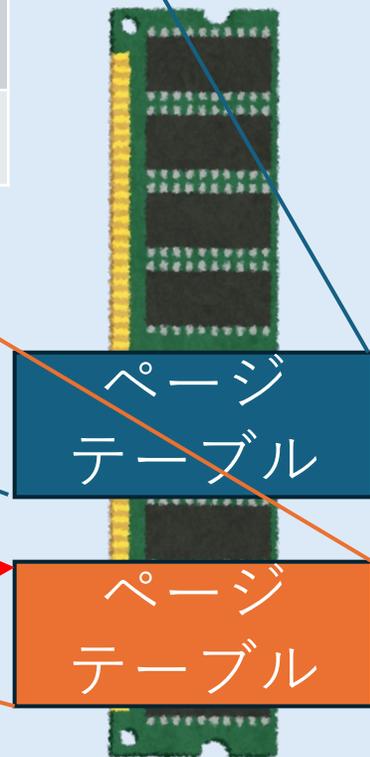
この場合
物理メモリアドレス 0x4000 ~ 0x4fff が
プロセス 1 とプロセス 2 で共有されます

カーネルは
各プログラム（プロセス）ごとに
ページテーブルを用意する

プロセス 2

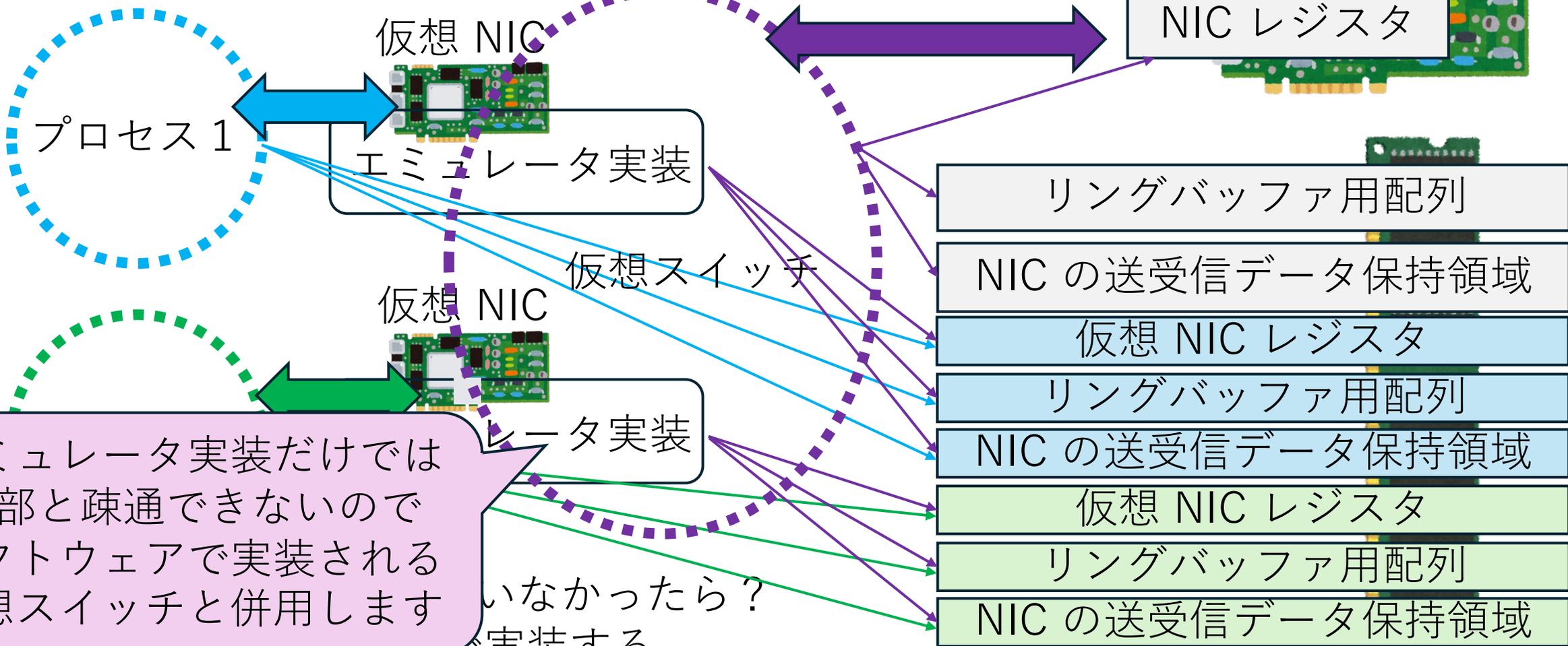


cr3: **仮想** アドレス 0x1000 へアクセス



DPDK を使って仮想スイッチを実装する場合の例

複数の要素が同一のリソースへアクセスしないようにする



エミュレータ実装だけでは外部と疎通できないのでソフトウェアで実装される仮想スイッチと併用します

仮想スイッチを実装する

DPDK を使って仮想スイッチを実装する場合の例

複数の要素が同一のリソースへアクセスしないようにする

ポイント

物理 NIC へアクセスできるのは
仮想スイッチを実行するプロセスのみ

NIC レジスタ



リングバッファ用配列

NIC の送受信データ保持領域

仮想 NIC レジスタ

リングバッファ用配列

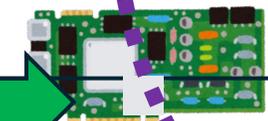
NIC の送受信データ保持領域

仮想 NIC レジスタ

リングバッファ用配列

NIC の送受信データ保持領域

仮想 NIC



レジスタ実装

いなかったら？

で実装する

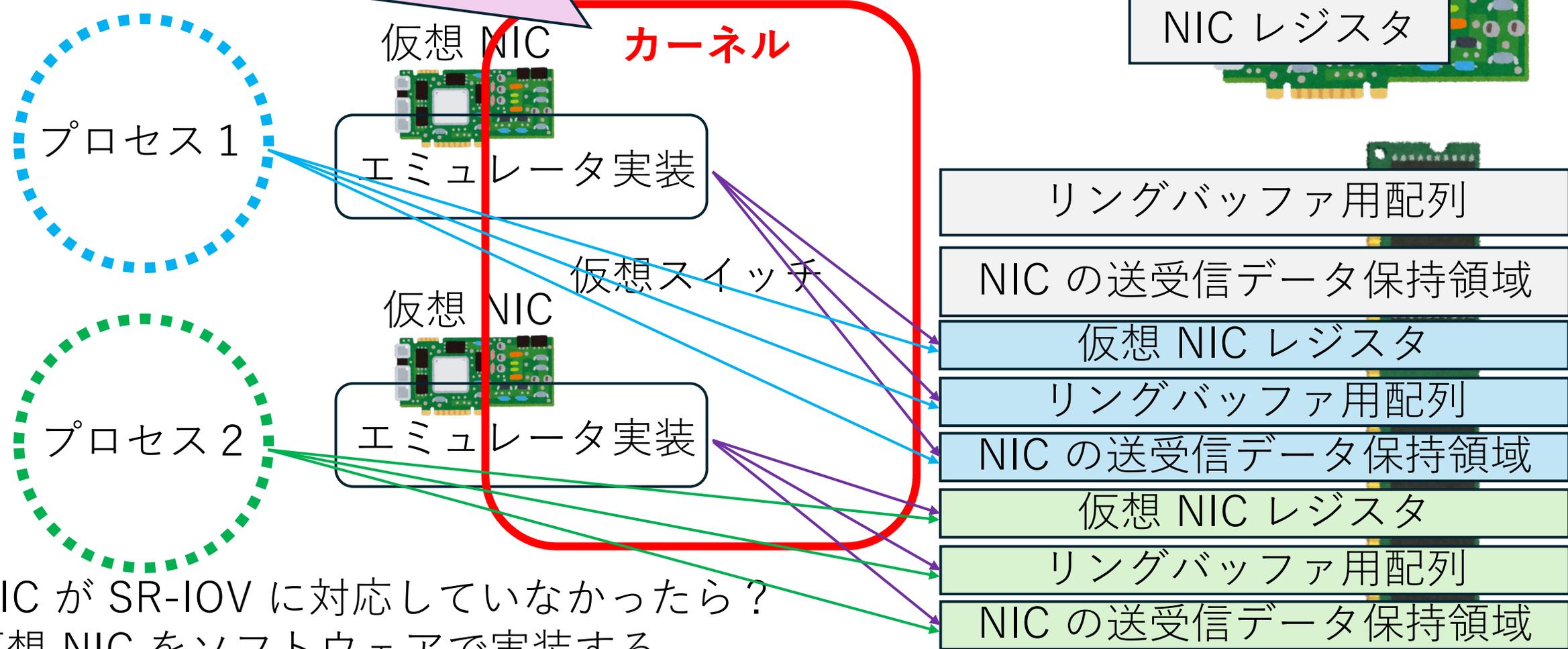
エミュレータ実装だけでは
外部と疎通できないので
ソフトウェアで実装される
仮想スイッチと併用します

仮想 NIC をソフトウェアで実装する

一つのNICを共有する仕組み

このような構成になっていれば仮想スイッチはカーネルに実装されていても大丈夫です (例: netmap / VALE)

複数の仮想NIC



Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

DPDK のインストール方法

- コマンド：

<https://github.com/yasukata/jumpstart-on-docker/tree/master#dpdk-installation>

- おすすめしたいポイント

- 実行するマシンごとにソースコードからコンパイルすること

- 何故？：コンパイル時に、利用可能な CPU 命令などに合わせた最適化が適用されるため

- 別のマシンでコンパイルしたライブラリファイルを利用する場合に対応していない CPU 命令が含まれているとプログラムが停止したりします

- インストール先のディレクトリを指定してコンパイルすること

- 何故？：複数の異なる DPDK ライブラリを保持できるようにするため

- インストール先が一般ユーザーのディレクトリであればインストール自体に root 権限が不要なのも良いです

DPDK の使い方の基本ポイント

- ポイント 1 : hugepages を設定する
 - DPDK 付属のツールで設定できます
 - 例 : `./dpdk-version/usertools/dpdk-hugepages.py -p 2M -r 2G`
 - `-p`: huge page size
 - `-r`: huge pages として扱う合計のメモリサイズ
- ポイント 2 : 利用するデバイスは PCI バス・デバイス・ファンクション番号で指定する
 - `lspci` コマンドで確認できます
 - 基本的にはカーネルのドライバと紐付けられているので、DPDK 付属のツールで DPDK から扱えるようにします
 - 例 : `./dpdk-version/usertools/dpdk-devbind.py -b vfio-pci XXXX:XX:XX.X`
 - `-b`: 紐付けるデバイスドライバです
 - `vfio-pci` は DPDK が扱えるようにするもので、Linux のカーネルに返す場合は、対応する NIC のドライバの名前を指定します
 - 設定可能なドライバの名前は `-s` オプションで確認できます

DPDK を使ってアプリを作る場合の おすすめの想定

- DPDK のライブラリは実行環境を操作する権限のある人からの提供を想定すべき
 - DPDK を利用するアプリケーションを配布する人が DPDK も一緒に配布しないようにした方が良い
- 異なるバージョンや利用者固有の改変が含まれた DPDK ライブラリと併用されることを想定しておく方が良い

DPDK のインストール方法

- コマンド：

<https://github.com/yasukata/jumpstart-on-docker/tree/master#dpdk-installation>

- おすすめしたいポイント

- 実行するマシンごとにソースコードからコンパイルすること

- 何故？：コンパイル時に、利用可能な CPU 命令などに合わせた最適化が適用されるため

- 別のマシンでコンパイルしたライブラリファイルを利用する場合に対応していない CPU 命令が含まれているとプログラムが停止したりします

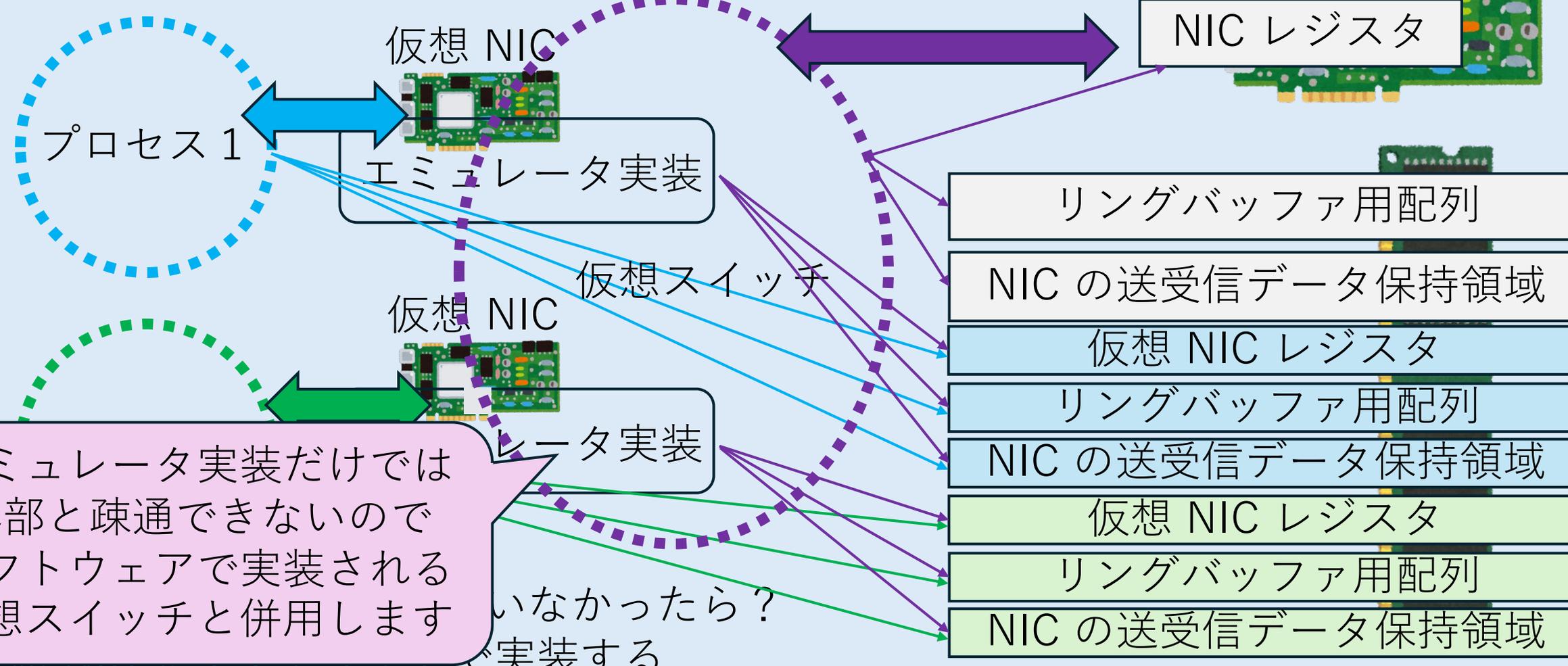
- インストール先のディレクトリを指定してコンパイルすること

- 何故？：複数の異なる DPDK ライブラリを保持できるようにするため

- インストール先が一般ユーザーのディレクトリであればインストール自体に root 権限が不要なのも良いです

DPDK を使って仮想スイッチを実装する場合の例

複数の要素が同一のリソースへアクセスしないようにする

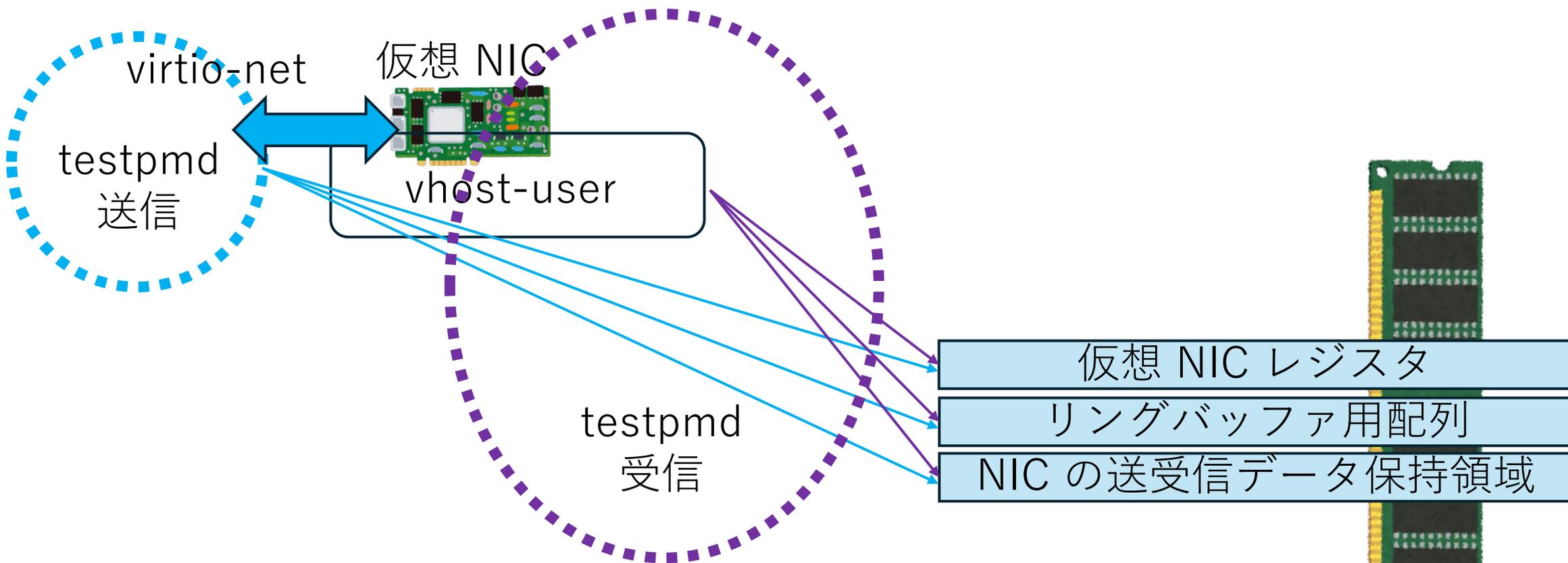


エミュレータ実装だけでは外部と疎通できないのでソフトウェアで実装される仮想スイッチと併用します

いなかったら？
ソフトウェアで実装する

DPDK のインストール方法

<https://github.com/yasukata/jumpstart-on-docker/tree/master#dpdk-installation>



testpmd プログラムを実行する二つのプロセスが片方からもう一方へ virtio-net/vhost-user インターフェースを通じてパケットを送信します

DPDK の使い方のポイント

- DPDK の API である `rte_eal_init()` の引数の渡し方だけ覚える
- よく使う引数
 - `-l` : 実行に利用する CPU コアのリスト
 - `--proc-type` : だいたい `primary` を指定しておけば大丈夫
 - `--file-prefix` : 起動するプロセスごとに変えることで、複数の DPDK を扱うアプリが実行できる
 - PCI デバイス関連
 - `--allow` : `lspci` コマンドで確認できる PCI アドレスを指定 (この指定がないと全てのアクセス可能な PCI デバイスを探索する)
 - `--no-pci` : PCI デバイスを探索しない
 - 仮想デバイス設定
 - `--vdev=net_` デバイスの種類, デバイスドライバ依存情報, ...
 - tap デバイスの例 : `--vdev=net_tap,iface=tap001`

DPDK の使い方のポイント

- DPDK の API である `rte_eal_init()` の引数の渡し方だけ覚える
- `testpmd` の例

```
./dpdk-testpmd ¥
```

```
-l 0,1 ¥
```

```
--proc-type=primary ¥
```

```
--file-prefix=pmd1 ¥
```

```
--vdev=net_vhost0,iface=/var/run/dpdk-app/vhost0,client=1 ¥
```

```
--no-pci ¥
```

```
--single-file-segments ¥
```

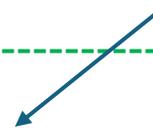
```
--
```

```
--nb-cores=1 ¥
```

```
--forward-mode=rxonly ¥
```

```
--stats-period=1"
```

`rte_eal_init()` に渡される



testpmd 固有の引数

簡単な仮想スイッチ l2fwd を使う

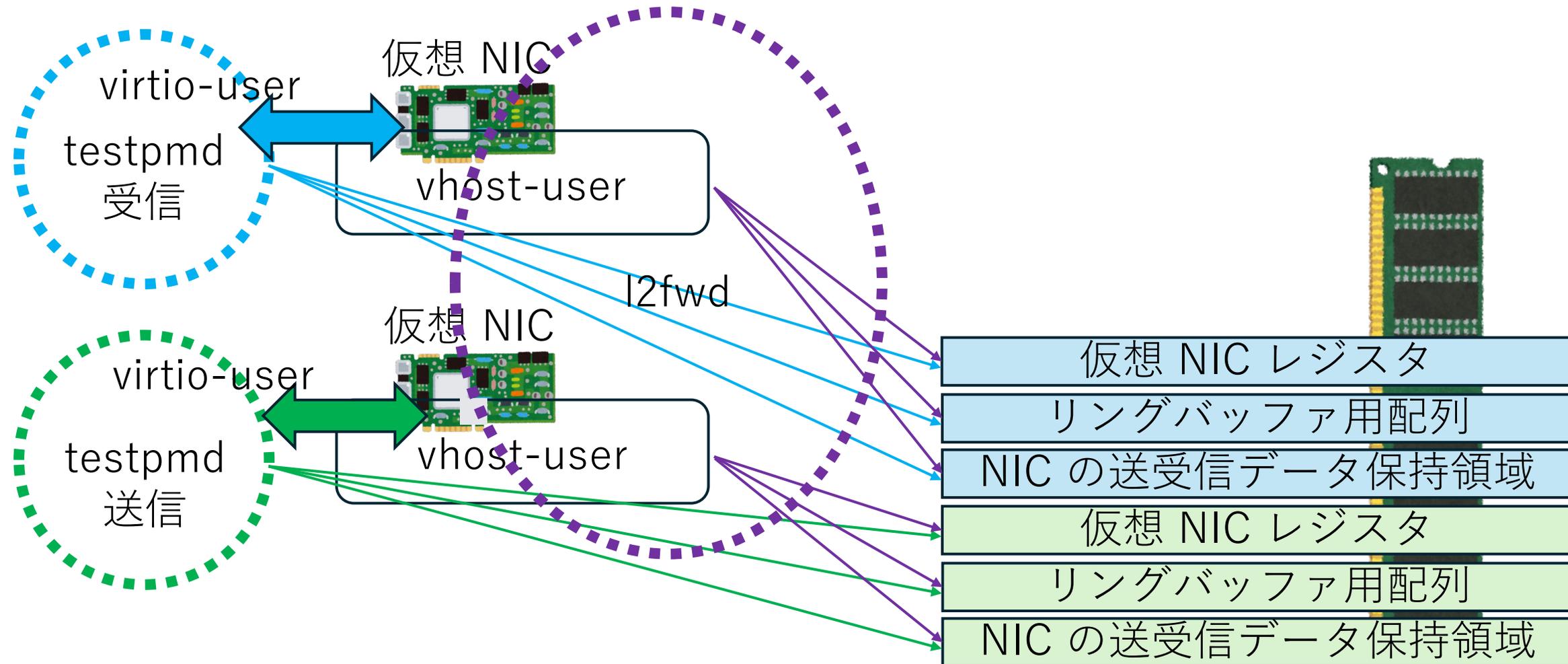
- コマンド：

<https://github.com/yasukata/jumpstart-on-docker#dpdk-l2fwd-bridging-dpdk-testpmd-containers>

- DPDK に付属している l2fwd という簡易的な仮想スイッチアプリで testpmd 間を接続することができます

簡単な仮想スイッチ l2fwd を使う

<https://github.com/yasukata/jumpstart-on-docker#dpdk-l2fwd-bridging-dpdk-testpmd-containers>



OVS-DPDK のインストール方法

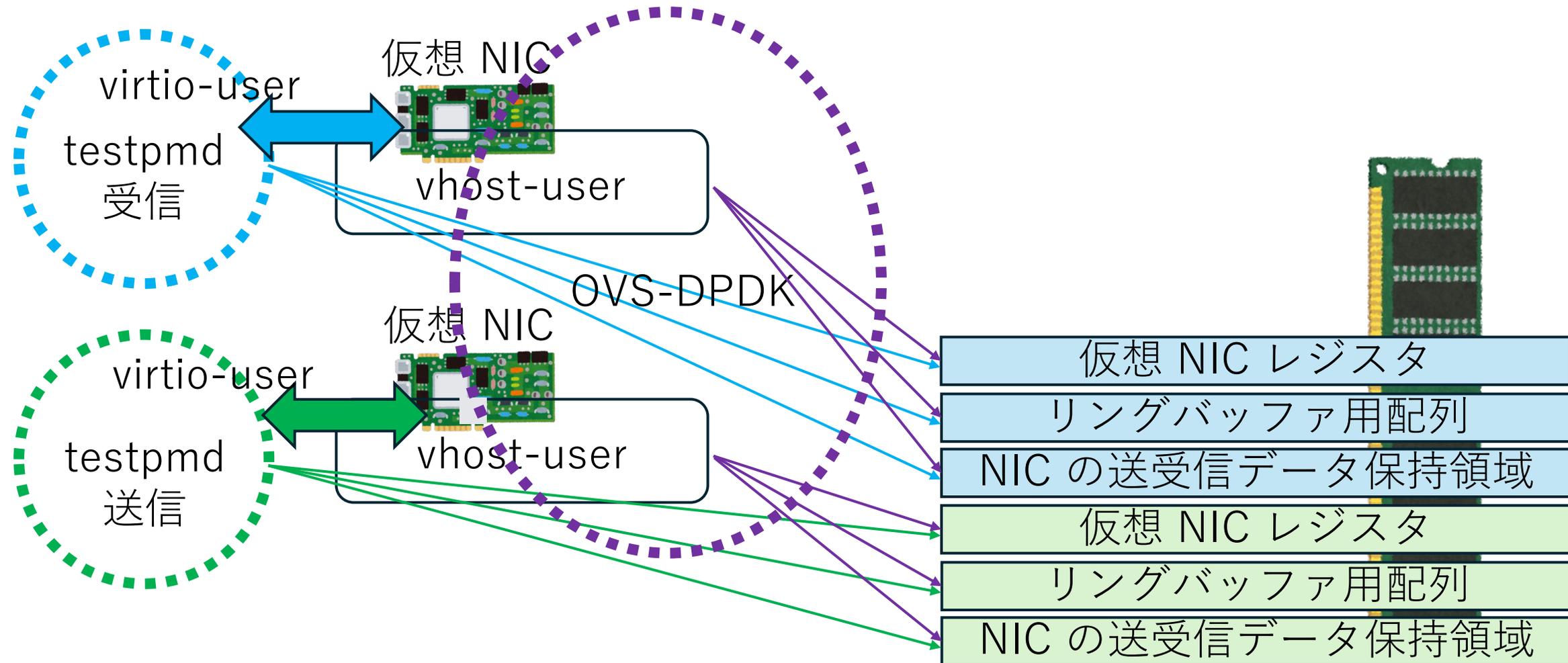
- コマンド :

<https://github.com/yasukata/jumpstart-on-docker#ovs-dpdk-installation>

- Open vSwitch で DPDK アプリ間を接続することができます

OVS-DPDK のインストール方法

<https://github.com/yasukata/jumpstart-on-docker#ovs-dpdk-installation>



VPP のインストール方法

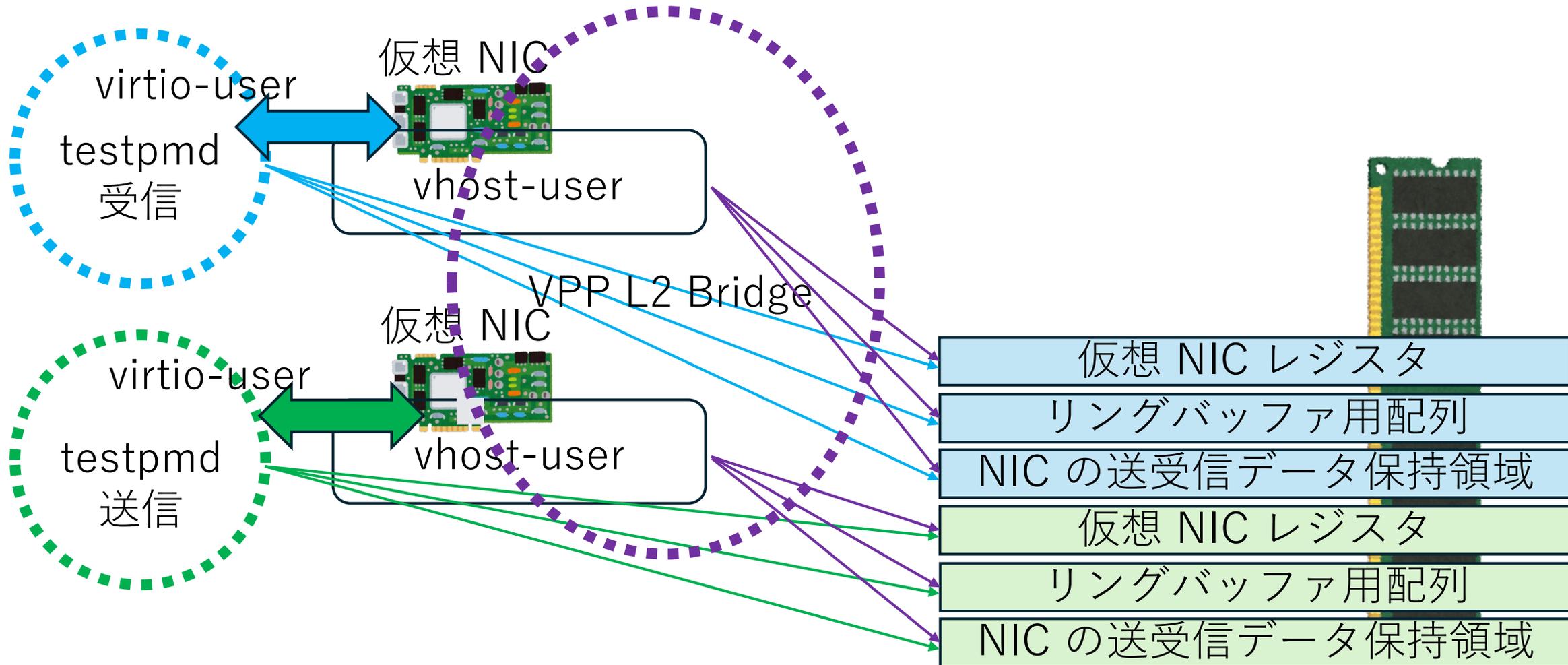
- コマンド：

<https://github.com/yasukata/jumpstart-on-docker#vpp-installation>

- VPP の L2 Bridge 機能 で DPDK アプリ間を接続します

VPP のインストール方法

<https://github.com/yasukata/jumpstart-on-docker#ovs-dpdk-installation>



TCP/IP 通信を行う方法



TCP/IP 通信を行う方法

- 通常 TCP/IP スタックはカーネルに実装されているので、カーネルをバイパスすると TCP/IP 通信ができない



TCP/IP 通信を行う方法

- 通常 TCP/IP スタックはカーネルに実装されているので、カーネルをバイパスすると TCP/IP 通信ができない

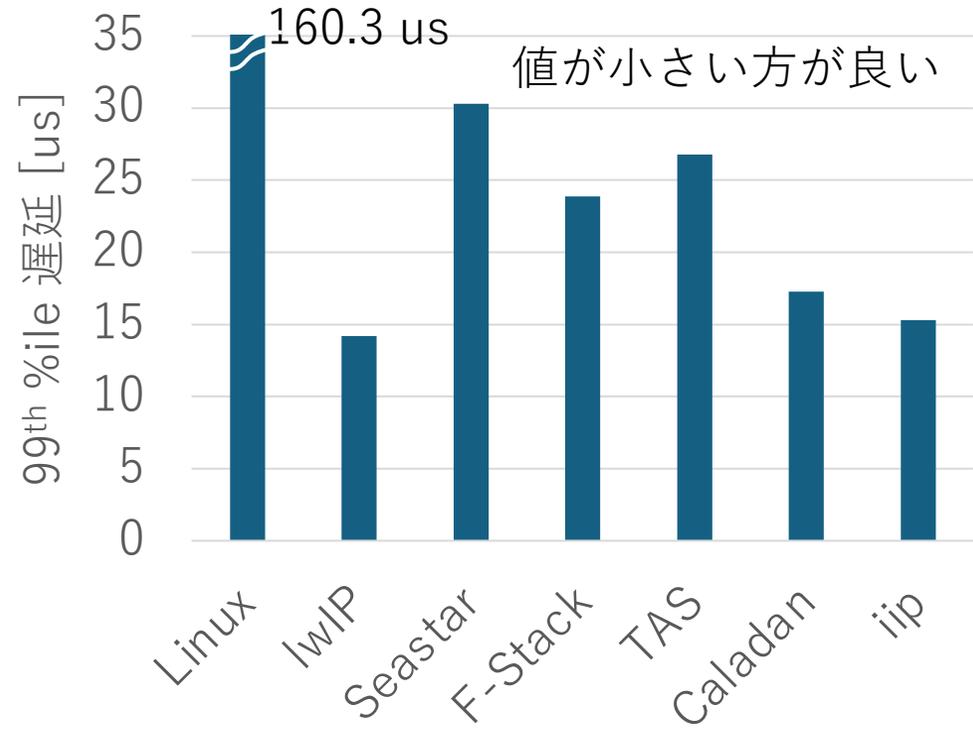
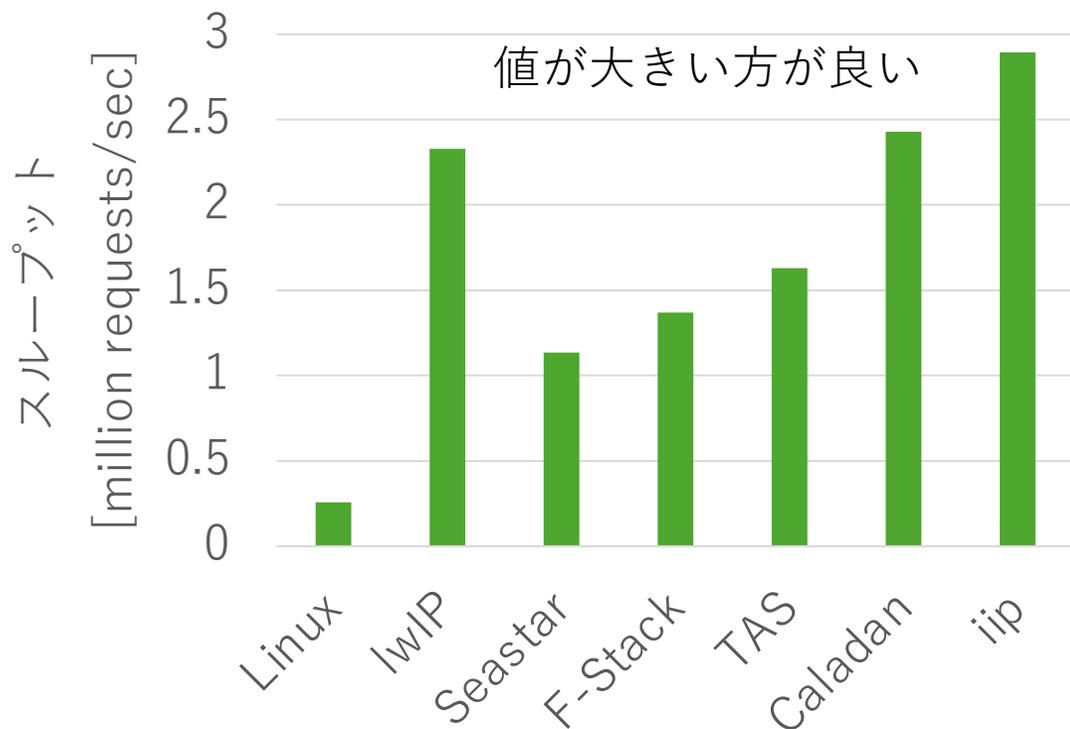


- アプリの一部として TCP/IP スタックを動かさせばカーネルをバイパスしながら TCP/IP 通信ができます

TCP/IP スタック実装の選択肢

アプリケーションが1 CPU コアを利用して 32 並列接続を通して1バイトのメッセージを往復させる場合の性能

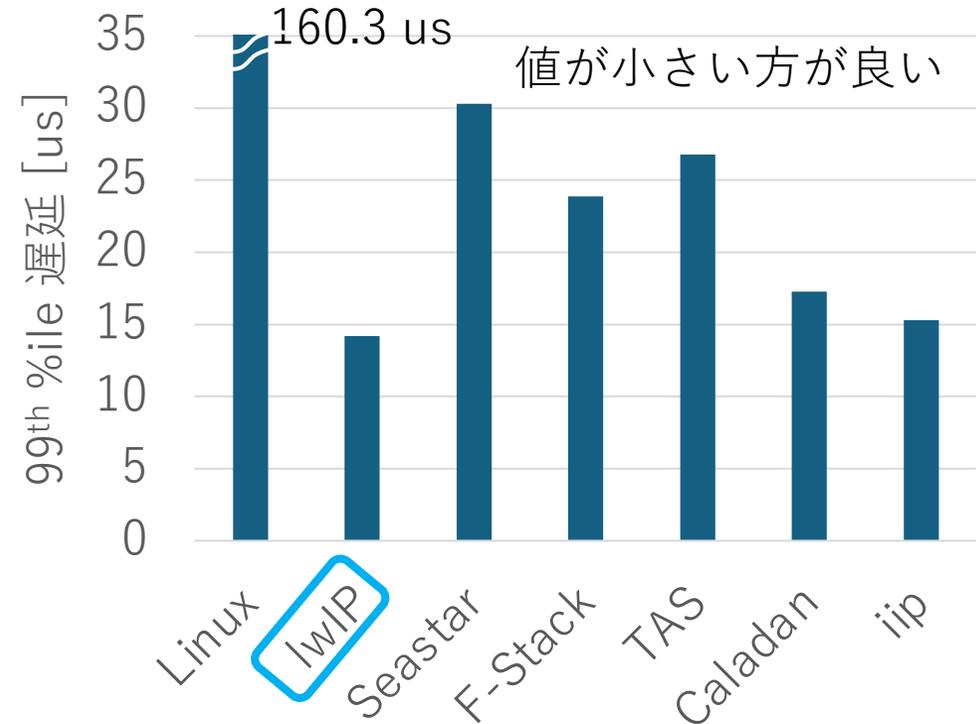
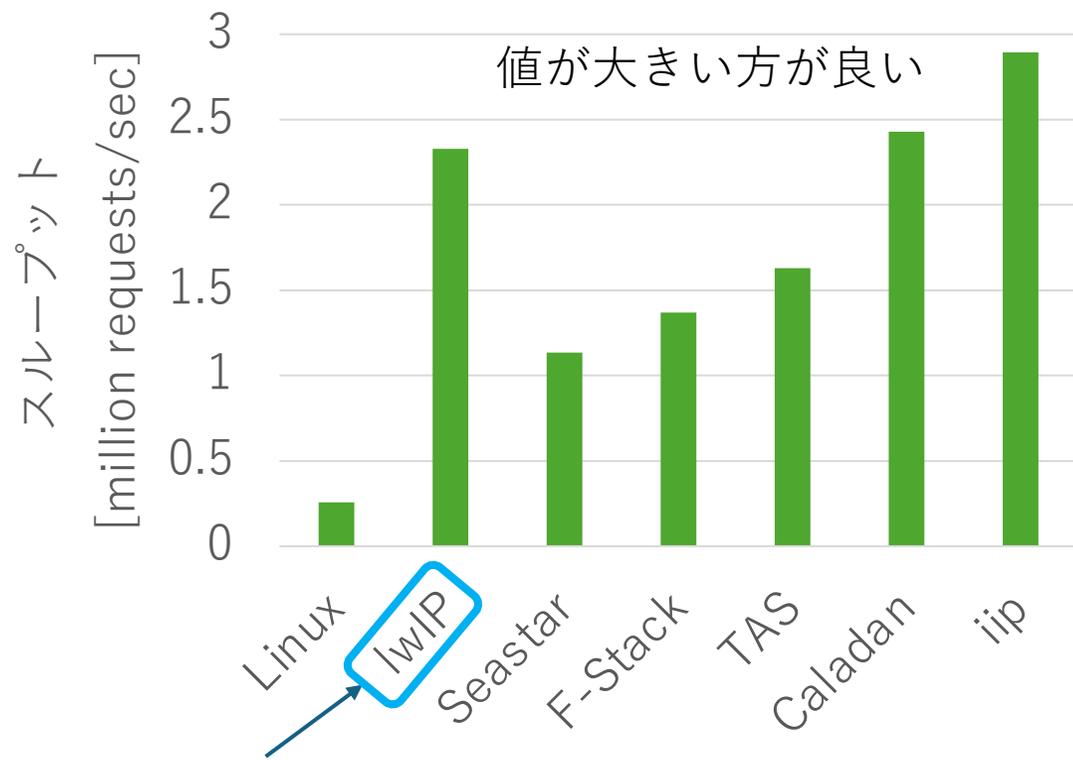
注意：それぞれの実装の機能が違うため公平な比較ではありません。目安としてご覧ください



TCP/IP スタック実装の選択肢

アプリケーションが1 CPU コアを利用して 32 並列接続を通して1バイトのメッセージを往復させる場合の性能

注意：それぞれの実装の機能が違うため公平な比較ではありません。目安としてご覧ください



個人的なおすすめ：lwIP は組み込み用途で長らく広く利用され、最高速を目指す研究でも採用されています

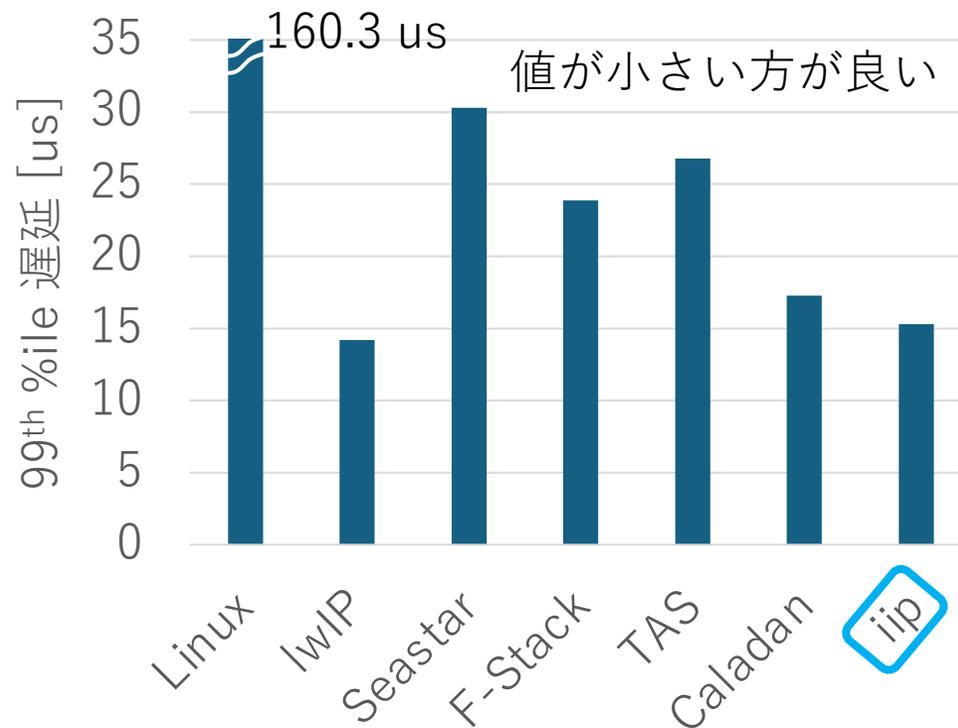
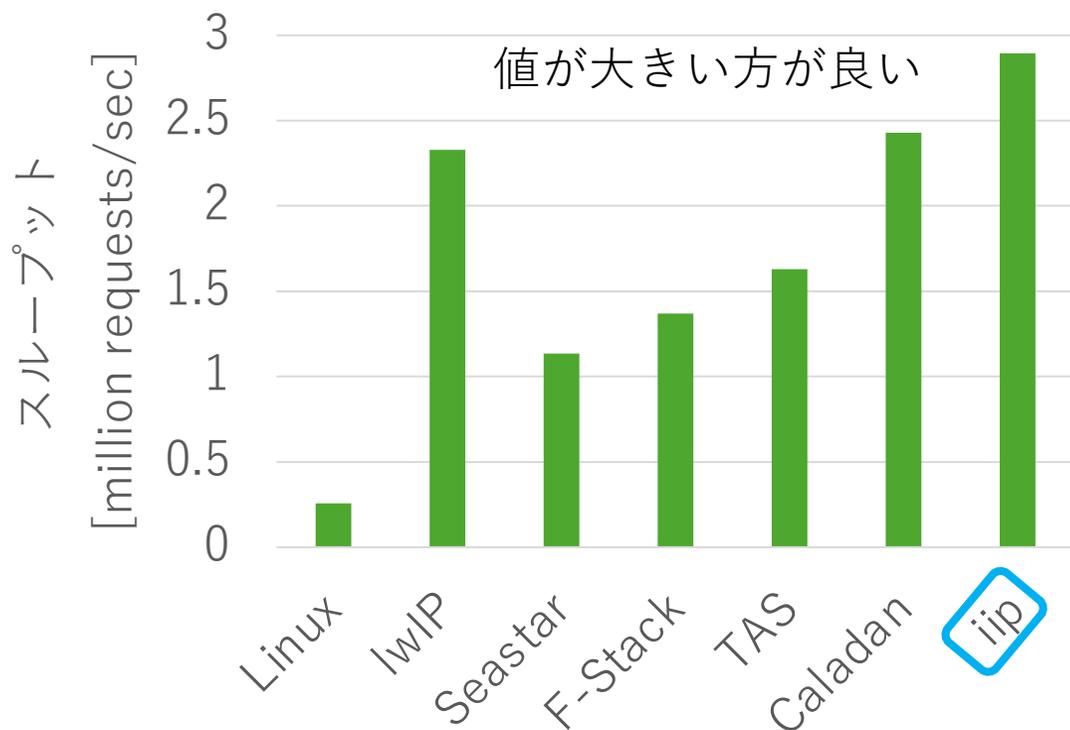
(ですが、lwIP は1 CPU コアでしか利用できず NIC のオフロード機能にも対応していないという制約があります)

実験の設定とコマンド：<https://github.com/yasukata/bench-iip#performance-numbers-of-other-tcpip-stacks>

TCP/IP スタック実装の選択肢

アプリケーションが1 CPU コアを利用して 32 並列接続を通して1バイトのメッセージを往復させる場合の性能

注意：それぞれの実装の機能が違うため公平な比較ではありません。目安としてご覧ください



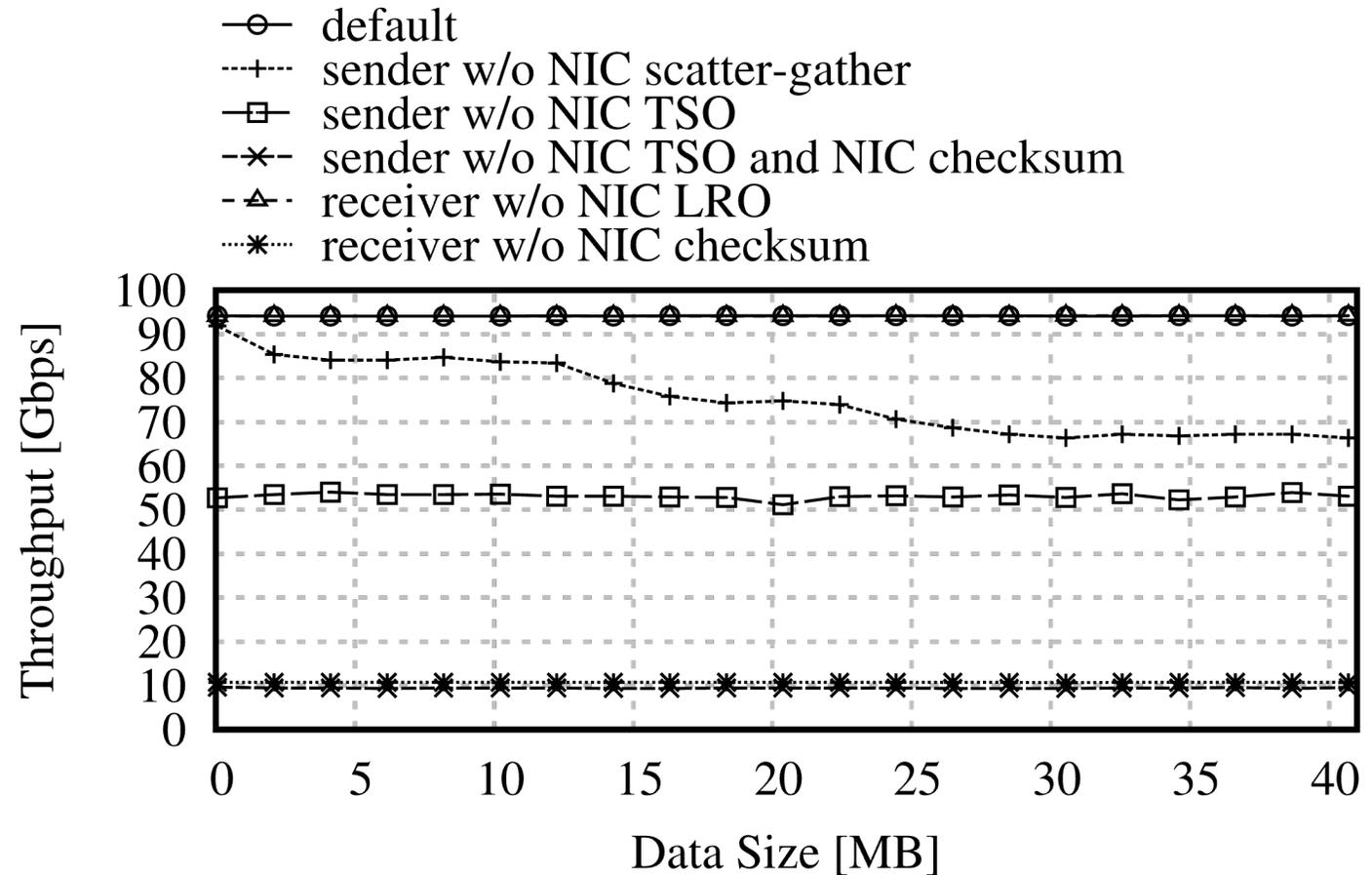
iip はこれら制約がないように作られているので、lwIP で物足りなくなったら試してみてください

(ですが、lwIP は1 CPU コアでしか利用できず NIC のオフロード機能にも対応していないという制約があります)

実験の設定とコマンド：<https://github.com/yasukata/bench-iip#performance-numbers-of-other-tcpip-stacks>

NIC のオフロード機能の効果

- 100 Gbps NIC を通じて大きなデータを片方のマシンからもう一方のマシンへ送る場合

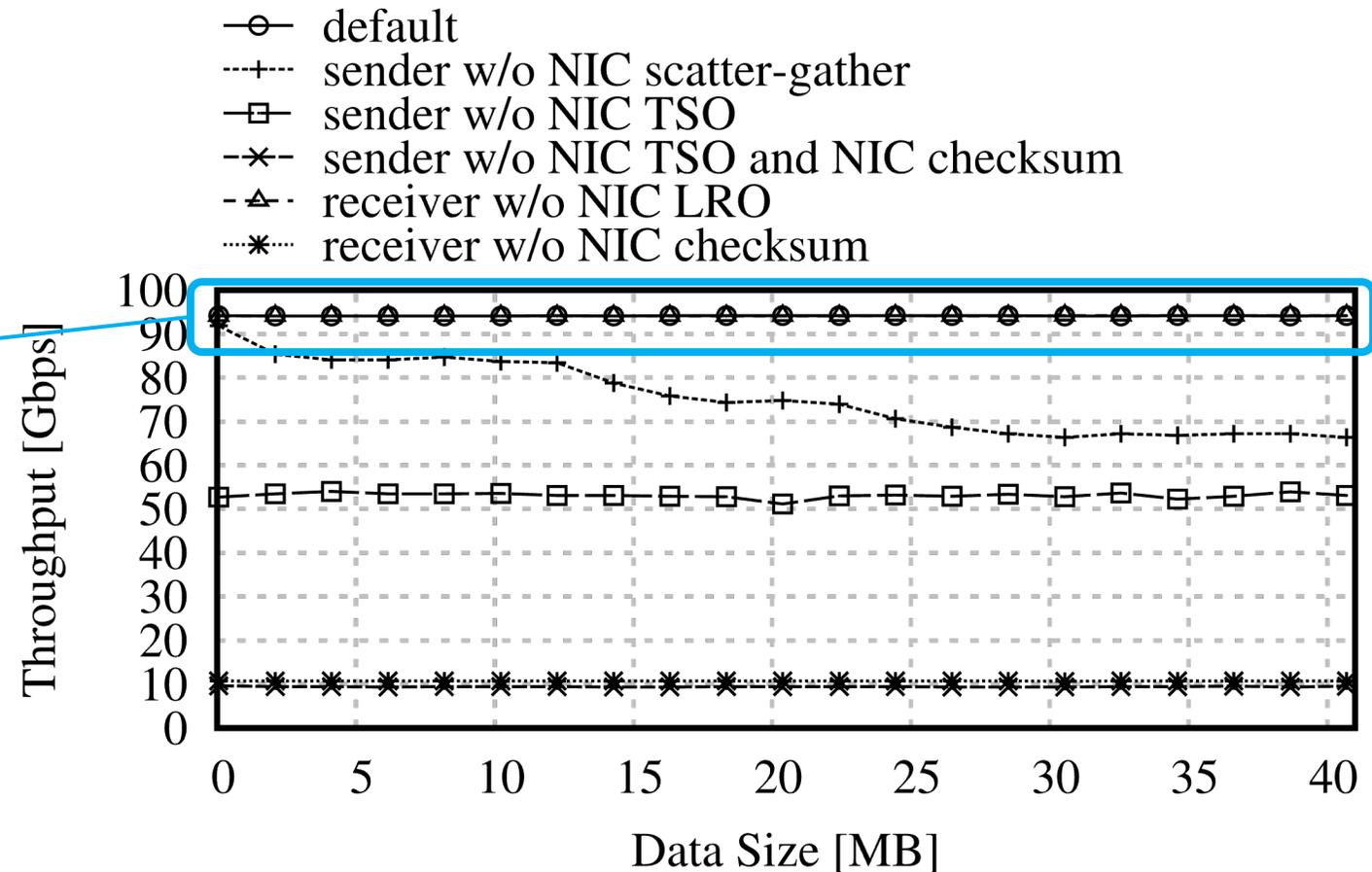


NIC のオフロード機能の効果

- 100 Gbps NIC を通じて大きなデータを片方のマシンからもう一方のマシンへ送る場合

全部のオフロードが有効だと約 100 Gbps

NIC の速度が 10 Gbps くらいまでであればオフロード機能はなくても大丈夫かもしれません



NIC のオフロード機能の効果

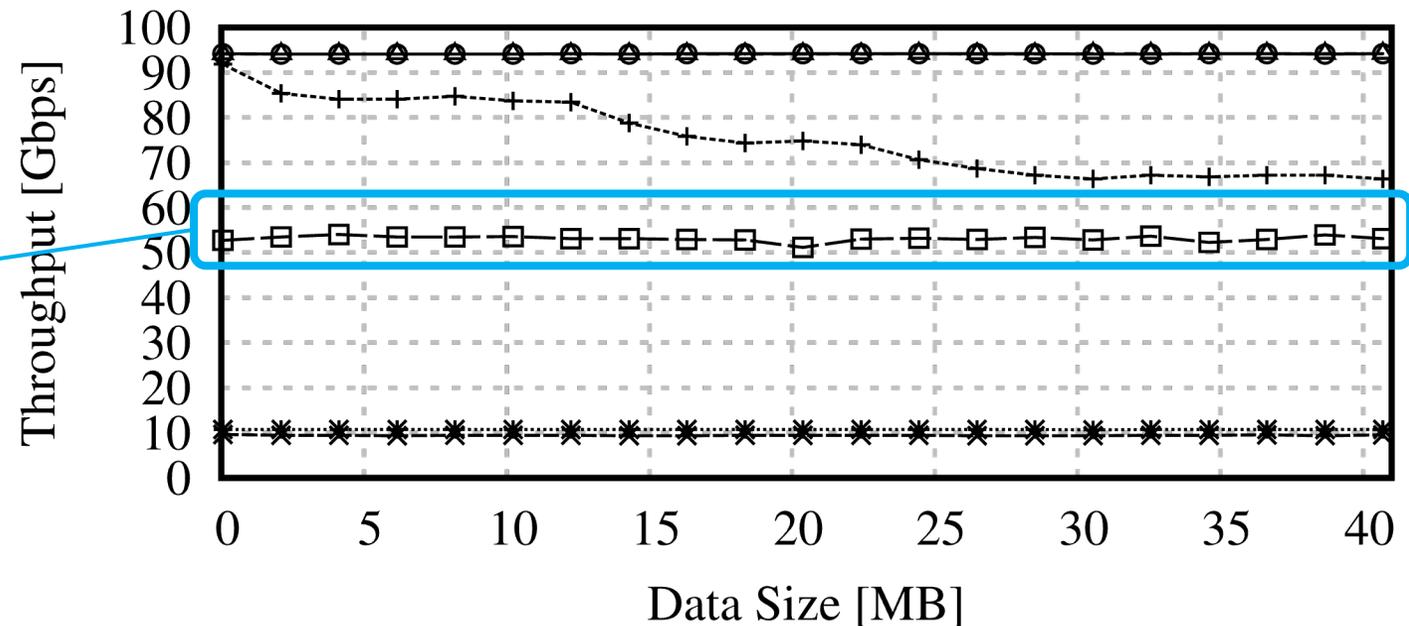
- 100 Gbps NIC を通じて大きなデータを片方のマシンからもう一方のマシンへ送る場合

全部のオフロードが有効だと約 100 Gbps

TSO が無効だと約 50 Gbps

NIC の速度が 10 Gbps くらいまでであればオフロード機能はなくても大丈夫かもしれません

- default
- + sender w/o NIC scatter-gather
- sender w/o NIC TSO
- * sender w/o NIC TSO and NIC checksum
- △ receiver w/o NIC LRO
- * receiver w/o NIC checksum



NIC のオフロード機能の効果

- 100 Gbps NIC を通じて大きなデータを片方のマシンからもう一方のマシンへ送る場合

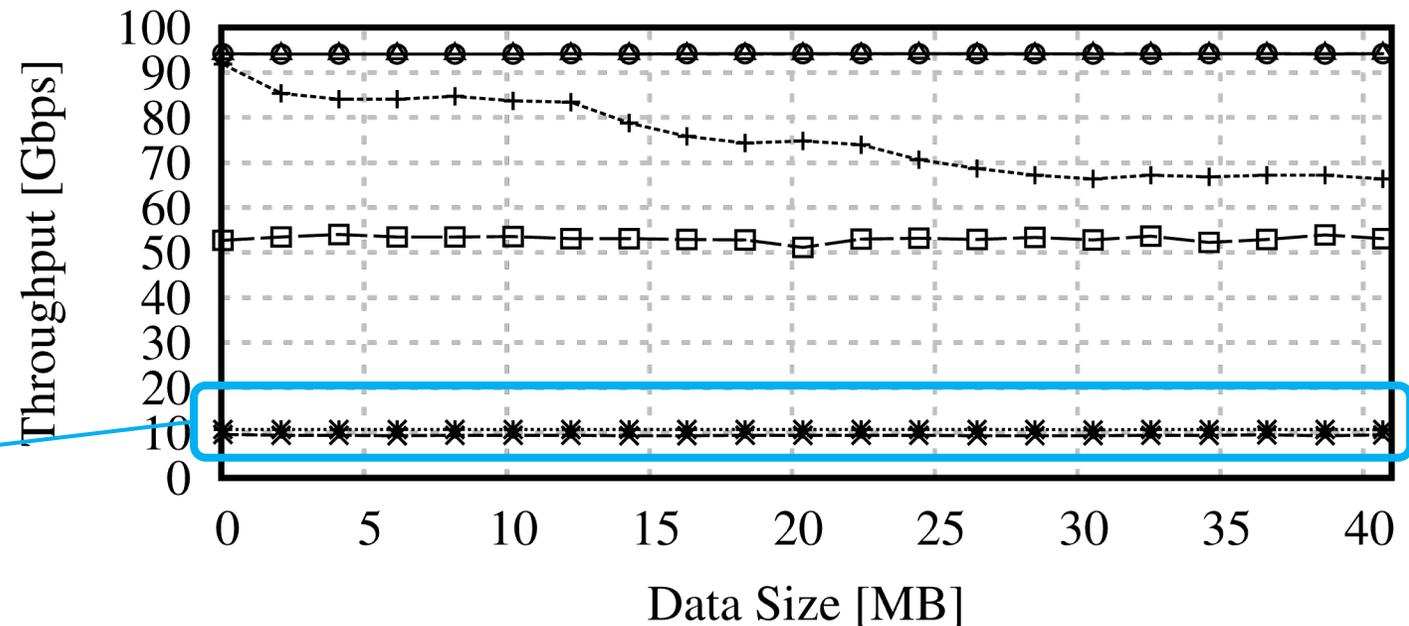
全部のオフロードが有効だと約 100 Gbps

TSO が無効だと約 50 Gbps

チェックサムオフロードが無効だと約 10 Gbps

NIC の速度が 10 Gbps くらいまでであればオフロード機能はなくても大丈夫かもしれません

- default
- + sender w/o NIC scatter-gather
- sender w/o NIC TSO
- * sender w/o NIC TSO and NIC checksum
- △ receiver w/o NIC LRO
- * receiver w/o NIC checksum

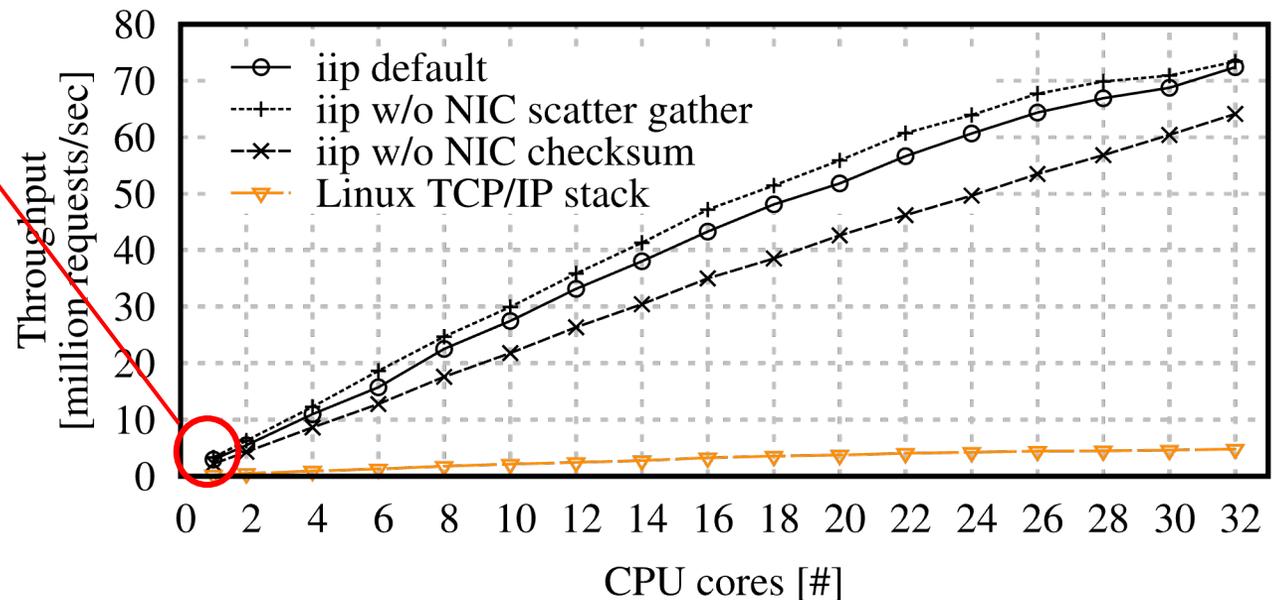


マルチコア環境での性能

- 100 Gbps NIC を経由して 2 つのマシン間で TCP 接続を通じて 1 バイトのメッセージを往復させる場合のスループット
 - 1 CPU コアが 32 並列接続 TCP を処理するように並列数を調整

1 CPU コアしか利用できない場合は
2~3 million requests / sec くらい

複数コアをうまく使えらると
もっとスループットを伸ばせます

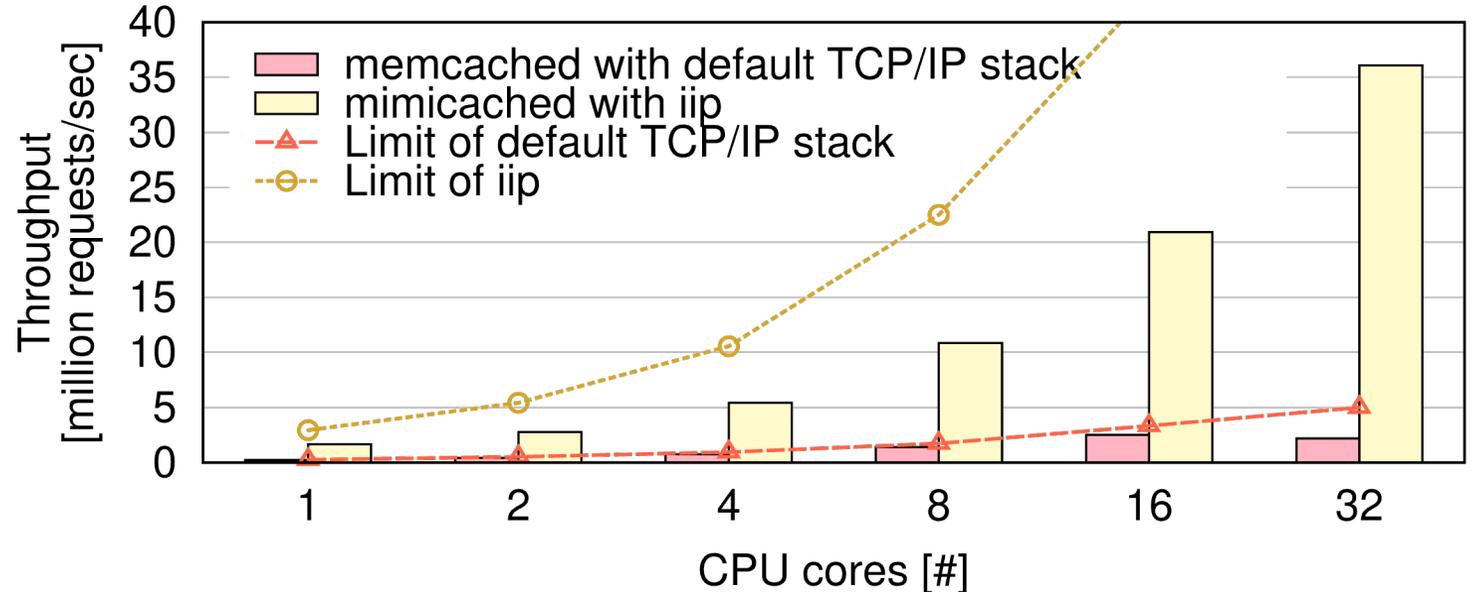


マルチコア環境での性能

- DPDK と iip の上で memcached 互換サーバーを動かすと
公式の memcached 実装よりも大幅に高い性能を発揮できます
 - <https://github.com/yasukata/mimicached>

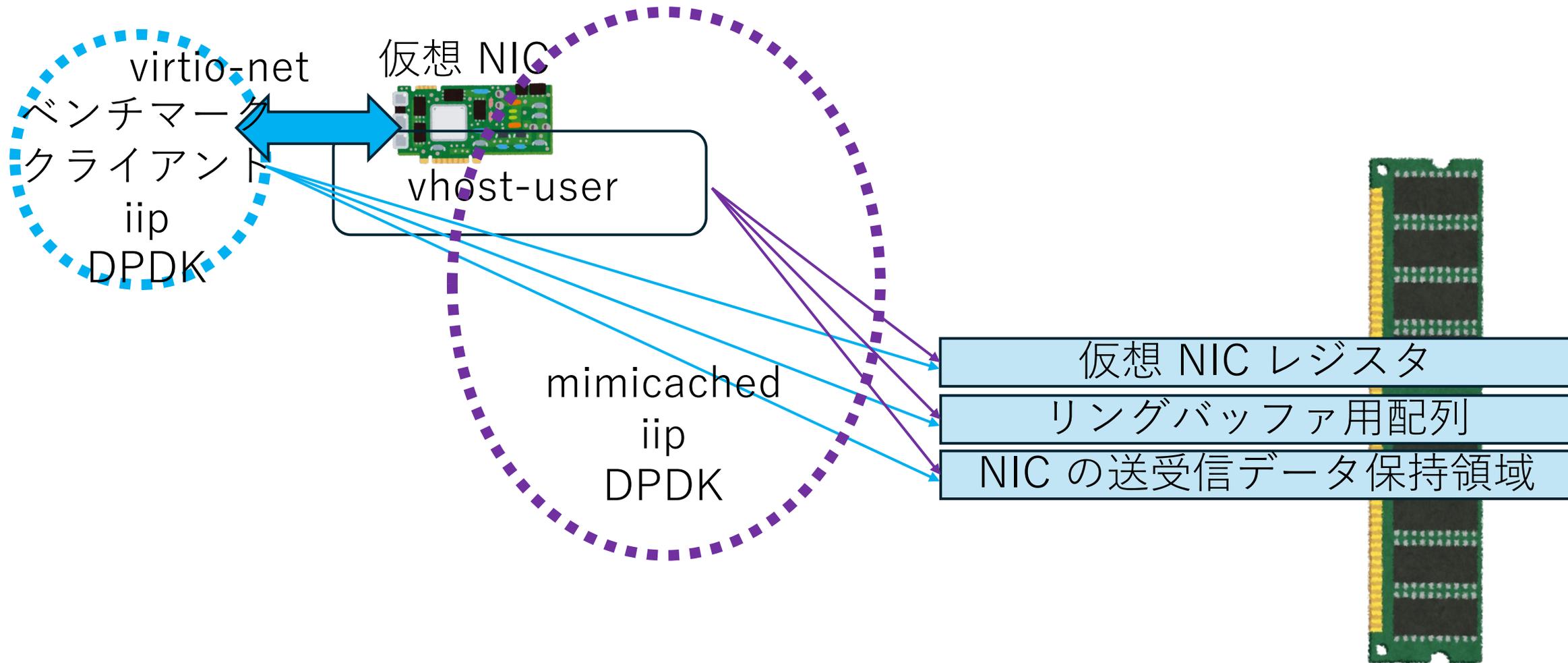
赤い線が Linux の TCP/IP スタックの
最高性能を表しており
インメモリストレージ実装が
速くてもこれ以上高速化できない
要因になっています

DPDK と iip を使うと
この要因を大幅に緩和できます



TCP/IP アプリの使い方

コマンド : <https://github.com/yasukata/jumpstart-on-docker#mimicached-on-dpdk>



AF_XDP のインストール方法

- コマンド：`apt install libxdp-dev`
 - AF_XDP のインストールというよりは、Linux に実装されている AF_XDP 機能を利用するためのライブラリのインストール
- AF_XDP の基本的な機能は Linux カーネルに含まれているものであるのでコンパイルが必要な DPDK よりも用意自体は手軽

AF_XDP 初期化時の設定ポイント

- 送受信リングサイズ：十分に大きい値が設定されている方が良いかも？
- XDP_FLAGS_DRV_MODE：できたら有効の方が速い
- XDP_USE_NEED_WAKEUP
 - 有効にしないとカーネルスレッドが busy loop して 100% CPU を使い続けたりするので基本は有効で良さそう、だが busy loop のおかげで遅延が減る場合もある（次ページ rx-usecs との兼ね合いもありそう）
- XDP_ZEROCOPY：できたら有効の方が速い
- SO_PREFER_BUSY_POLL

NIC のドライバのサポートの不足により設定に失敗してエラーが返ってくる場合はオプションを無効にして再度設定を試みるように実装するのが良いようです

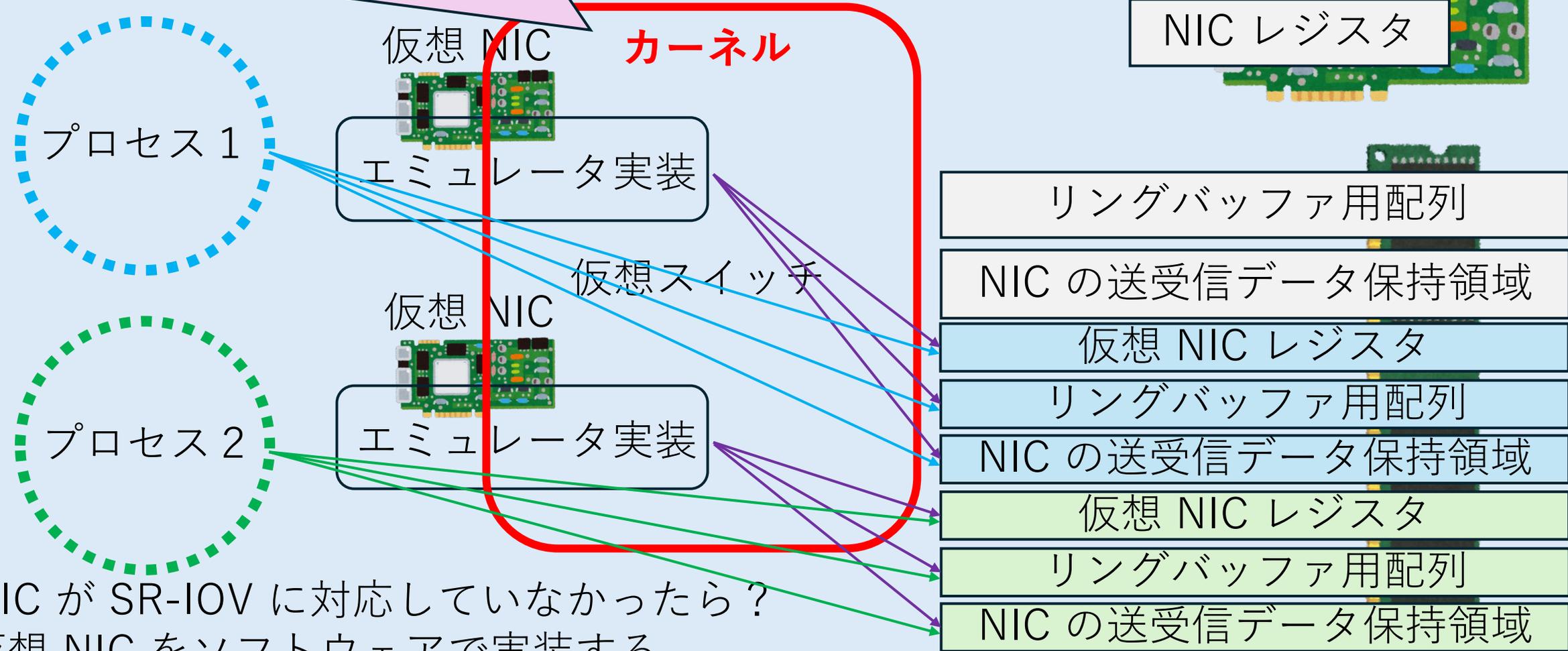
AF_XDP の設定のポイント

- 物理 NIC 利用時に ethtool で事前に設定すると良さそうな項目
 - NIC のキューの数をプログラムが使う数と揃える (-L)
 - 複数のキューが有効で NIC が受信したパケットを分散して振り分けている場合には、プログラムは全てのキューを確認しないと全てのパケットを受け取れない
 - 割り込み頻度 (-C)
 - rx-usecs 0 にしないと遅い場合がある？
 - カーネルの更新で改善される可能性もあると思います
 - 送受信リングサイズ (-G)
 - 十分に大きい値が設定されている方が良いかも？
- 仮想 NIC を利用したい場合はサポートの充実度の観点から veth を利用するのが良さそうです

一つのNICを共有する仕組み

このような構成になっていれば仮想スイッチはカーネルに実装されていても大丈夫です (例: netmap / VALE)

複数の仮想NIC

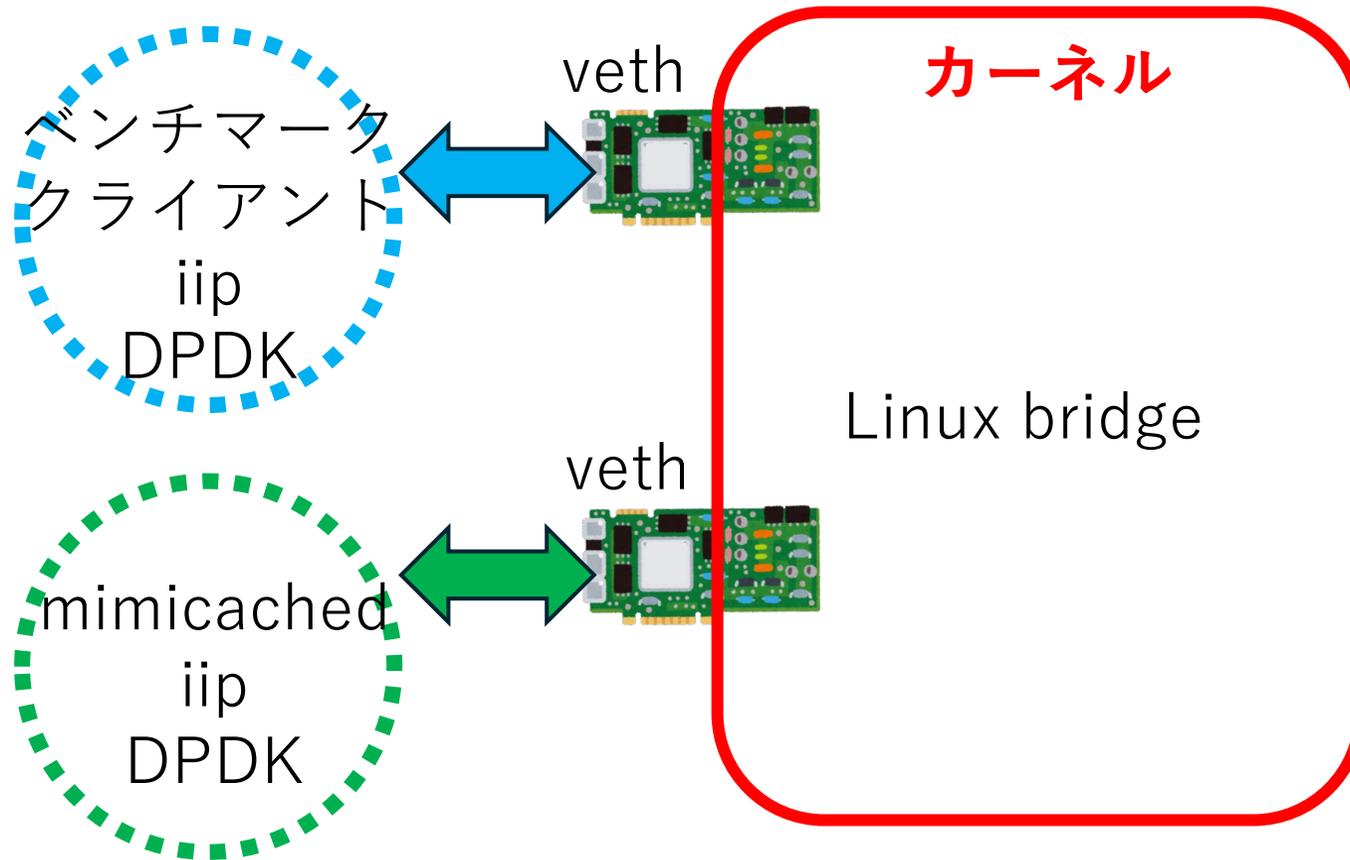


Q. NIC が SR-IOV に対応していなかったら？

A. 仮想 NIC をソフトウェアで実装する

TCP/IP アプリの使い方

コマンド : https://github.com/yasukata/jumpstart-on-docker?tab=readme-ov-file#mimicached-on-af_xdp



その他

- 高い性能を発揮するためには、頻繁に実行される箇所が適切に実装されている必要があります
- 特にデバイスドライバは重要で性能に大きく影響するとともに NIC ごとに DPDK 付属の実装と AF_XDP が扱う実装で精度にばらつきがあります
 - 必ずしも全てのデバイスドライバが最高の性能を発揮できるように実装されていなさそうです
- なので、単一のパケット I/O 機構に依存せず、複数のパケット I/O 機構の利用をサポートしておいて、環境に合わせて高い性能が発揮できる設定を選べるようにすることがおすすです

まとめ

- カーネルバイパス構成の理解に重要と思われるポイント
 1. プログラム（CPU 命令）によるメモリアクセス方法
 2. ページテーブルを通じたメモリアクセスの制限方法
 3. NIC のリングバッファの構成上記 3 点を基本とするとイメージがしやすくなると思います
- セキュリティのポイント
 - 複数の要素でリソースを共有しないこと
 - 信頼しない実装が NIC のレジスタへアクセスにする場合には IOMMU などを利用して NIC がアクセス可能な DRAM 領域を制限すること
- 一つの NIC を共有する場合のポイント
 - 仮想 NIC を使う（SR-IOV かソフトウェアでのエミュレーション）
 - 仮想スイッチで仮想 NIC に外部との疎通性を提供する