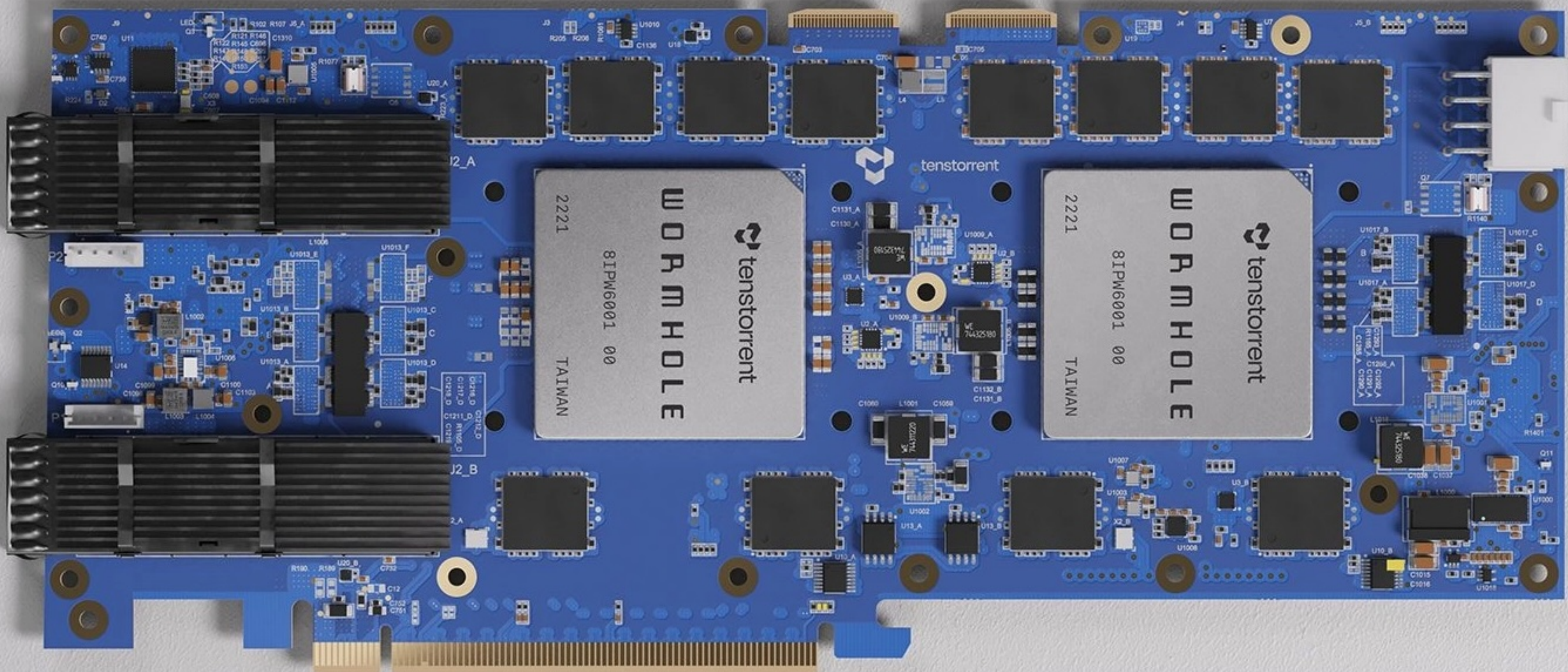


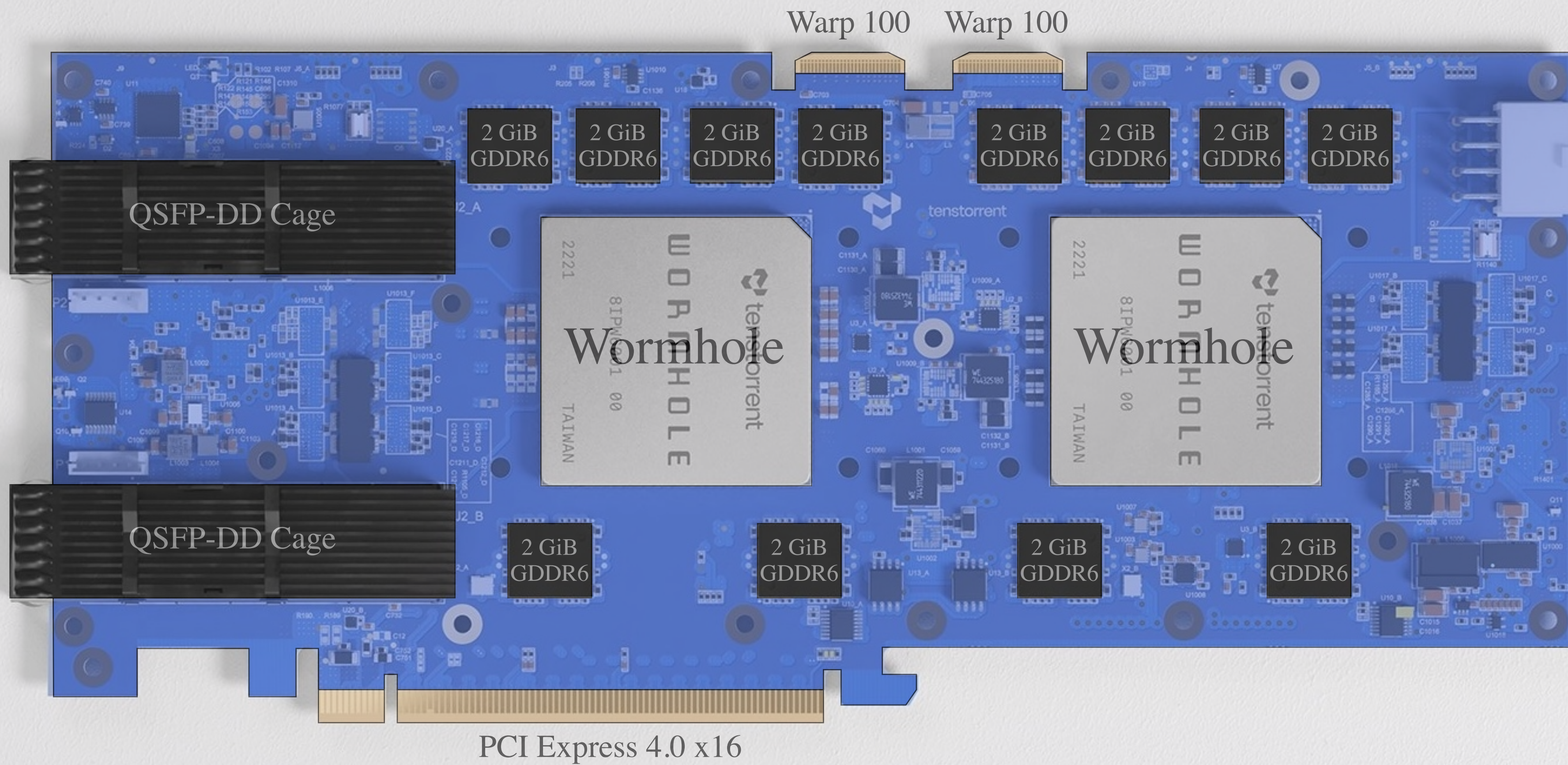
Scoping out the Tenstorrent Wormhole

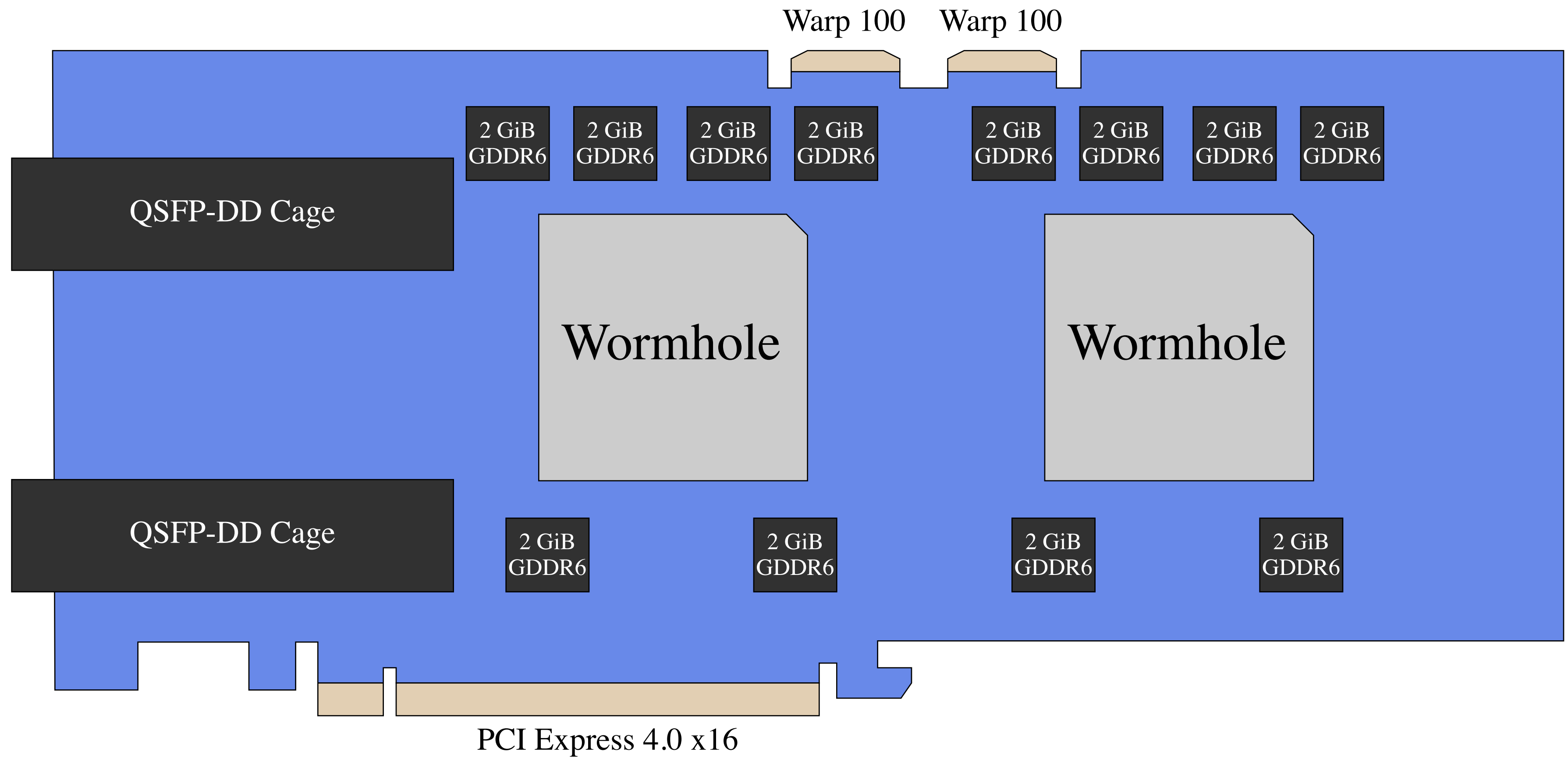
Peter Cawley (not affiliated with Tenstorrent ... yet)

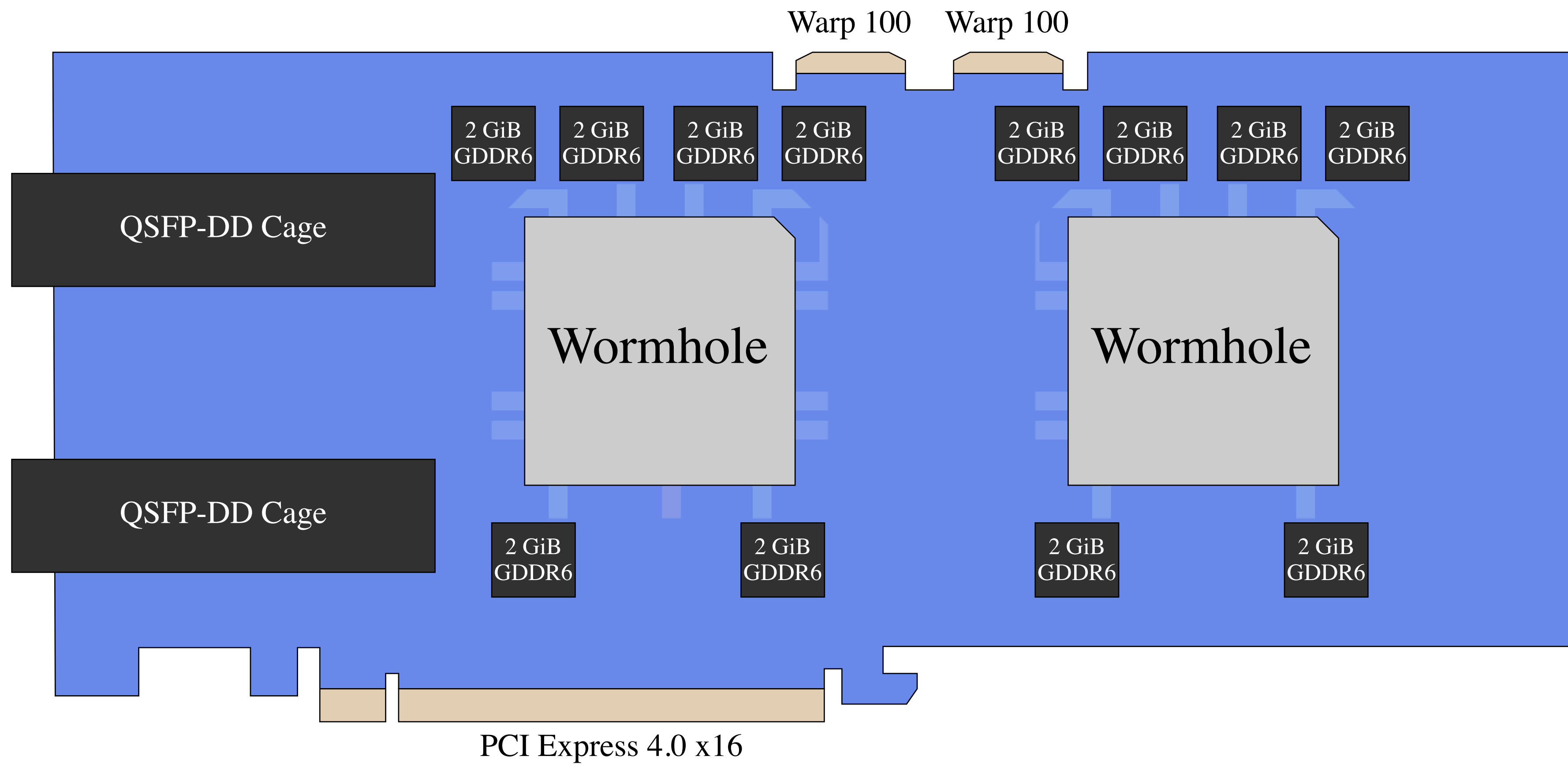
February 2nd, 2025 @ FOSDEM

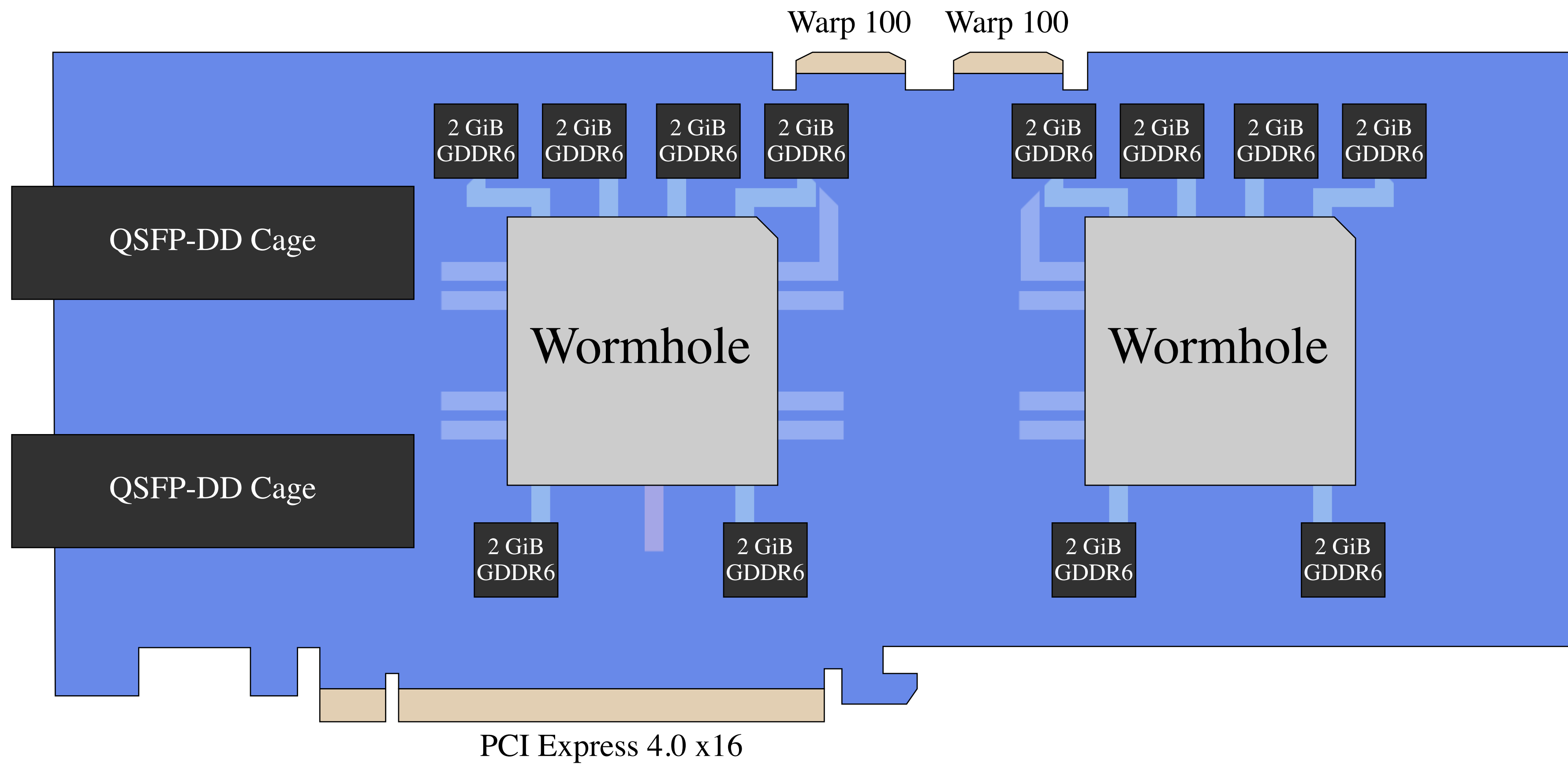
X @corsix

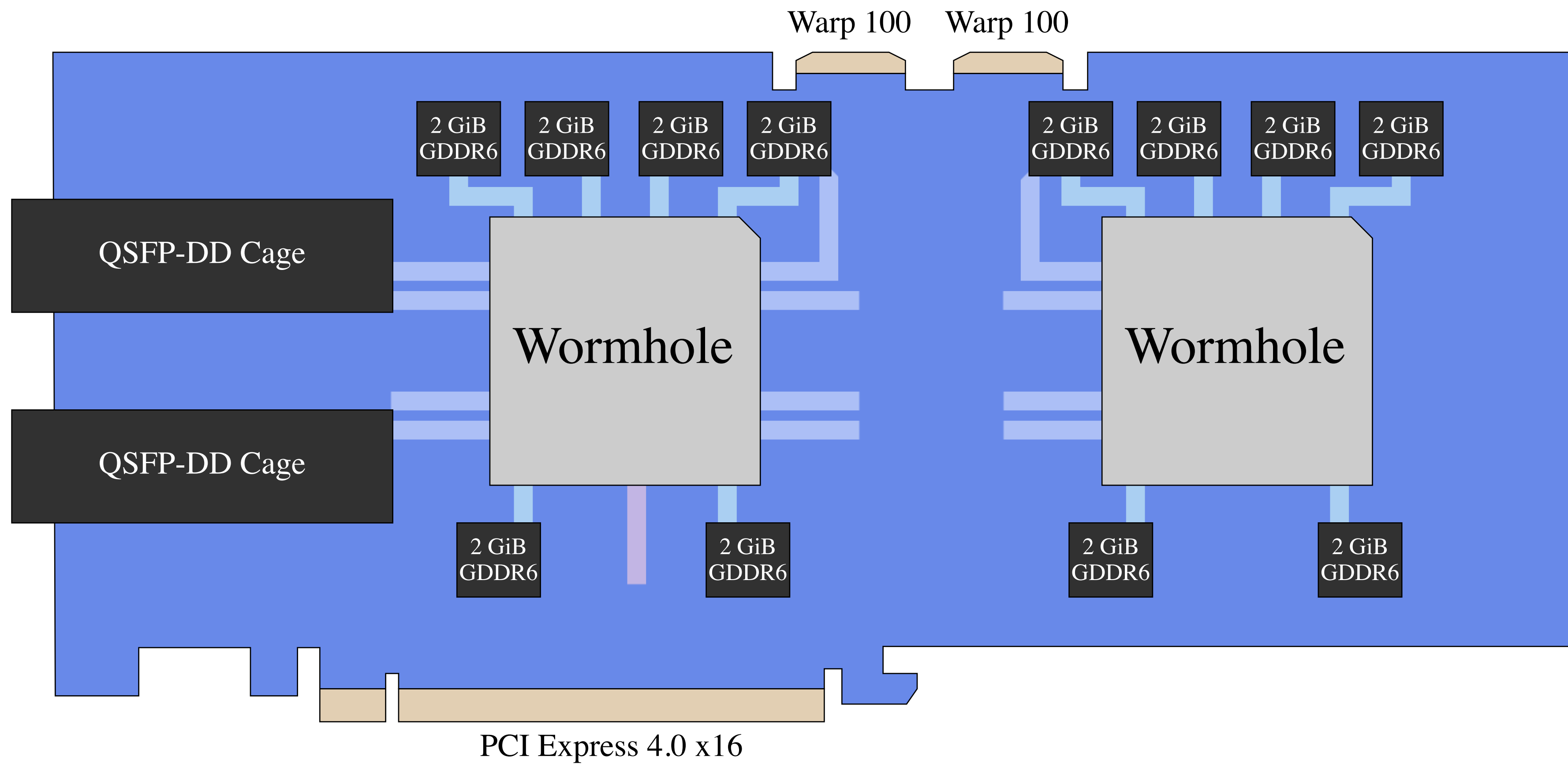


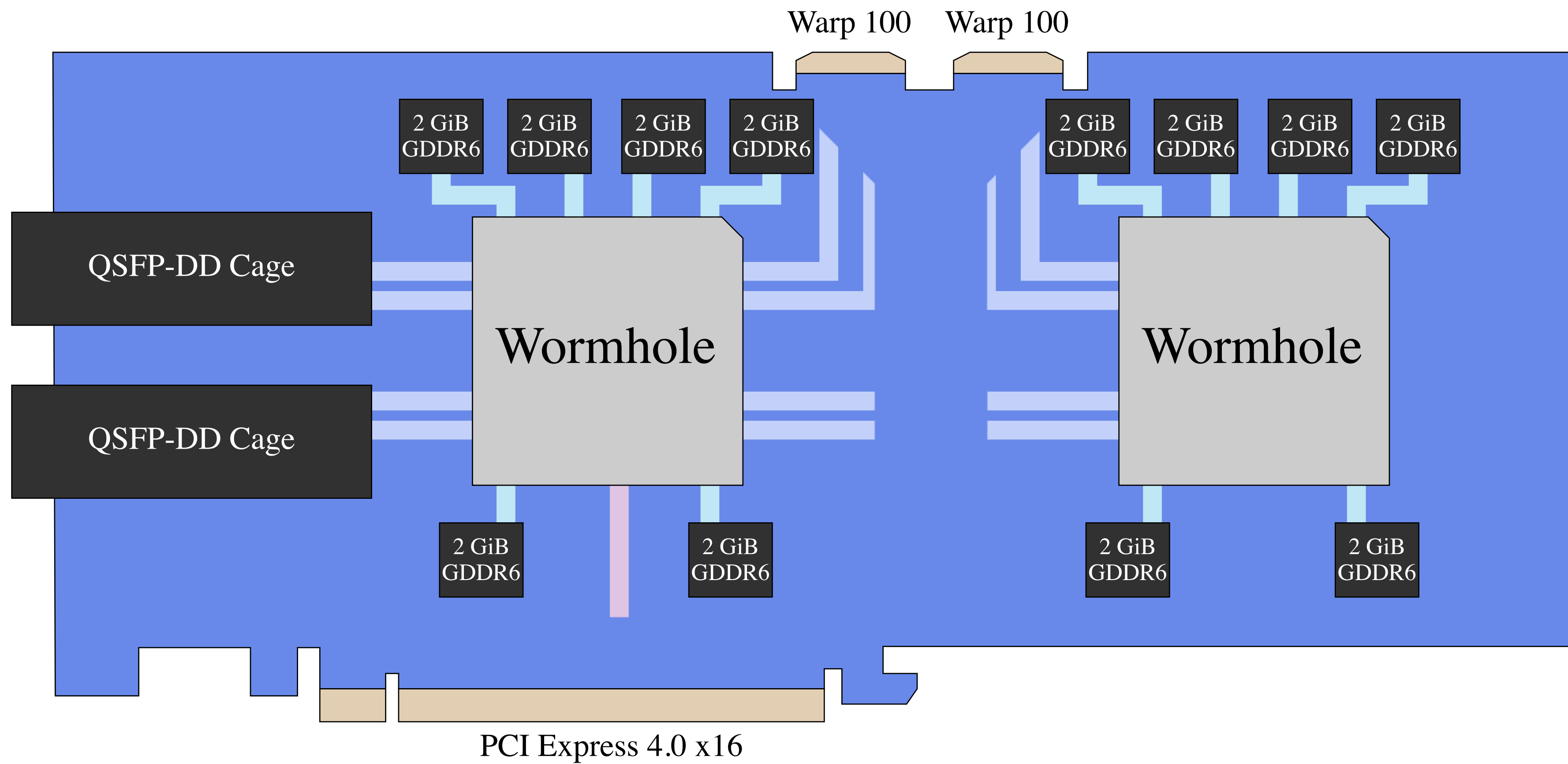


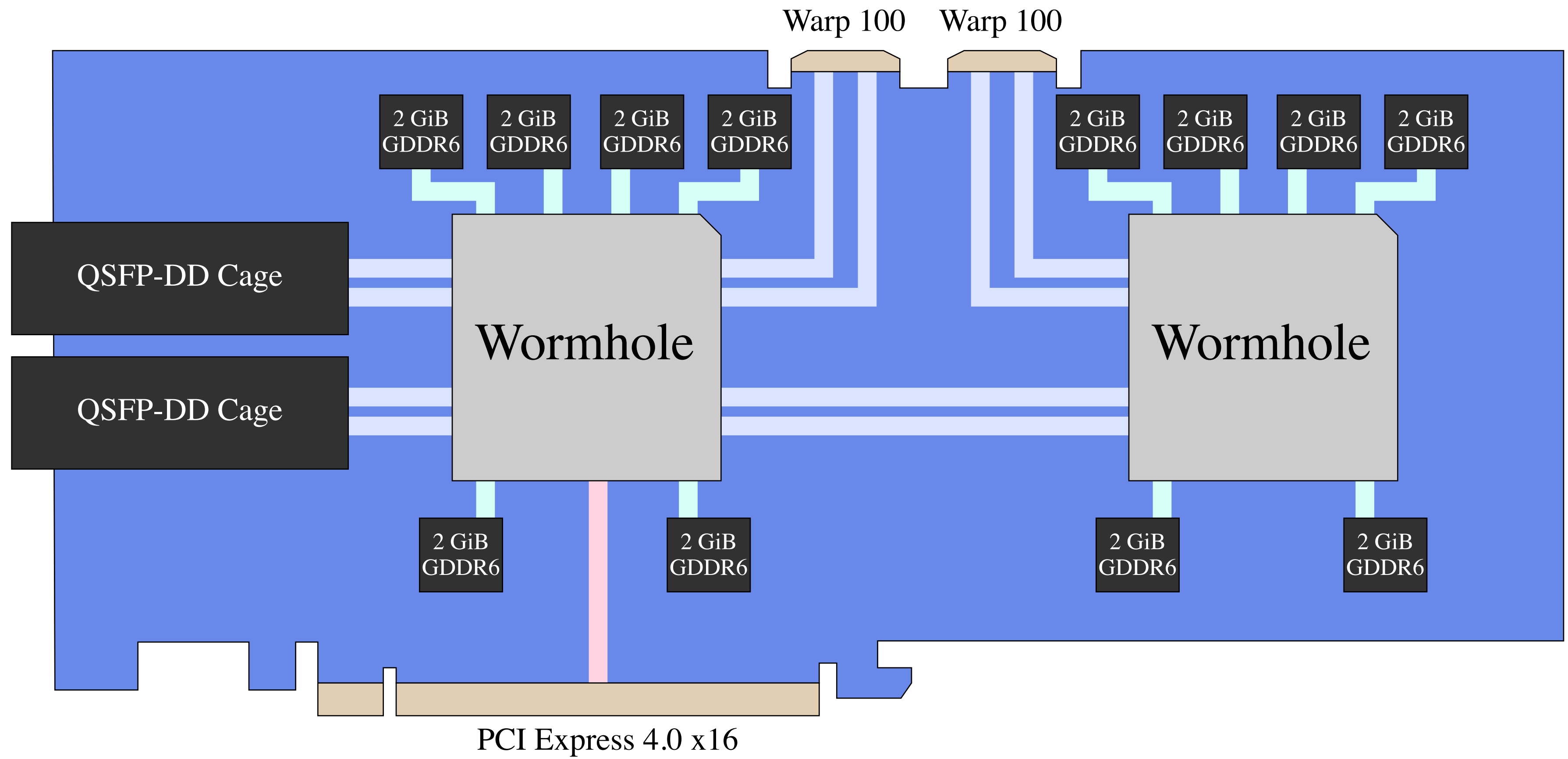


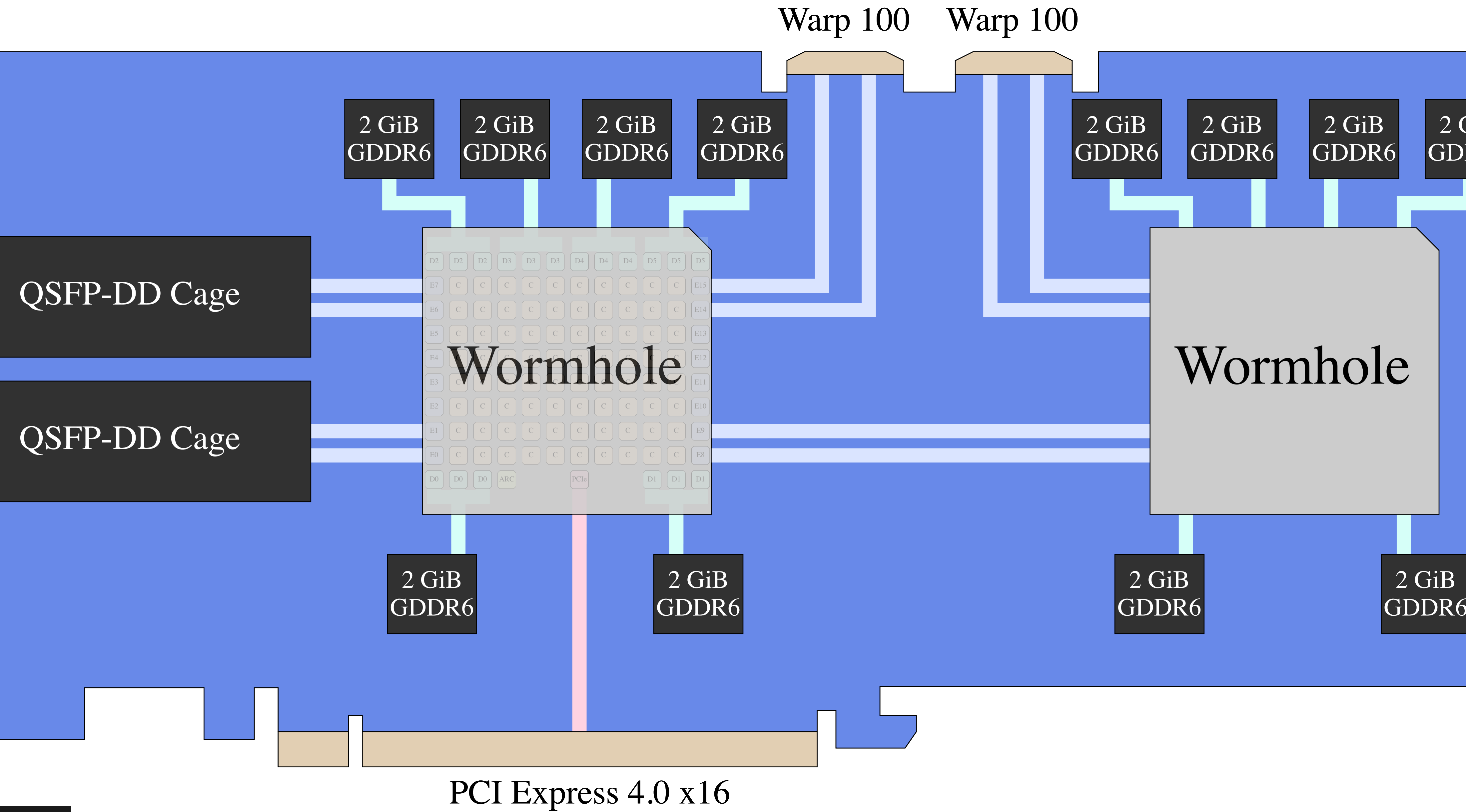


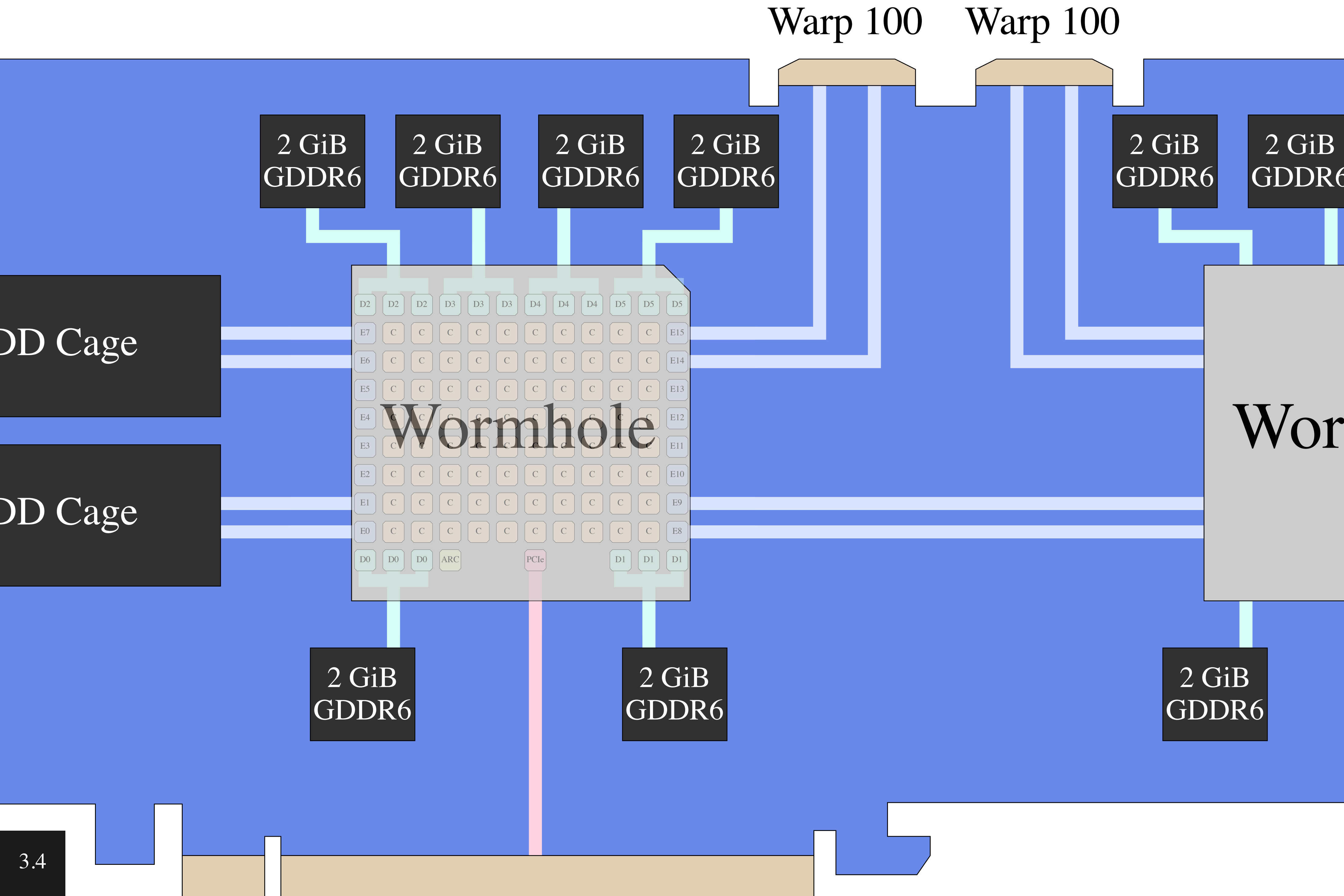












Warp 100

Warp 100

2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6

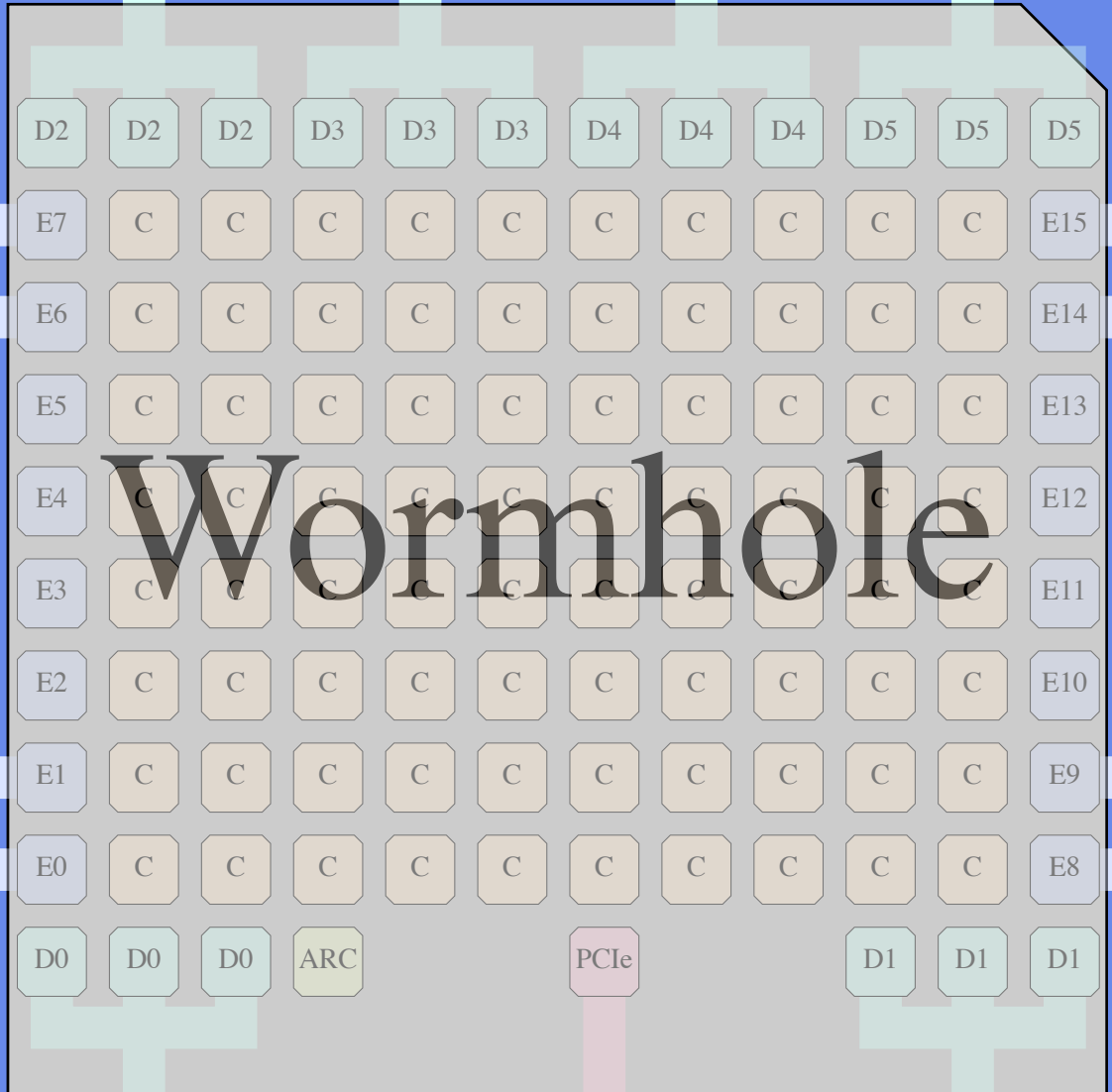
2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6

GDD Cage

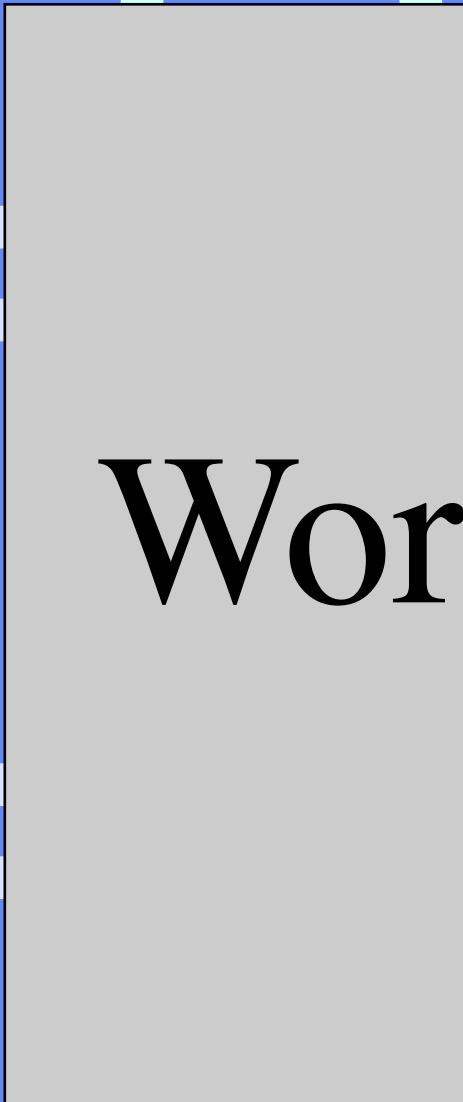
GDD Cage

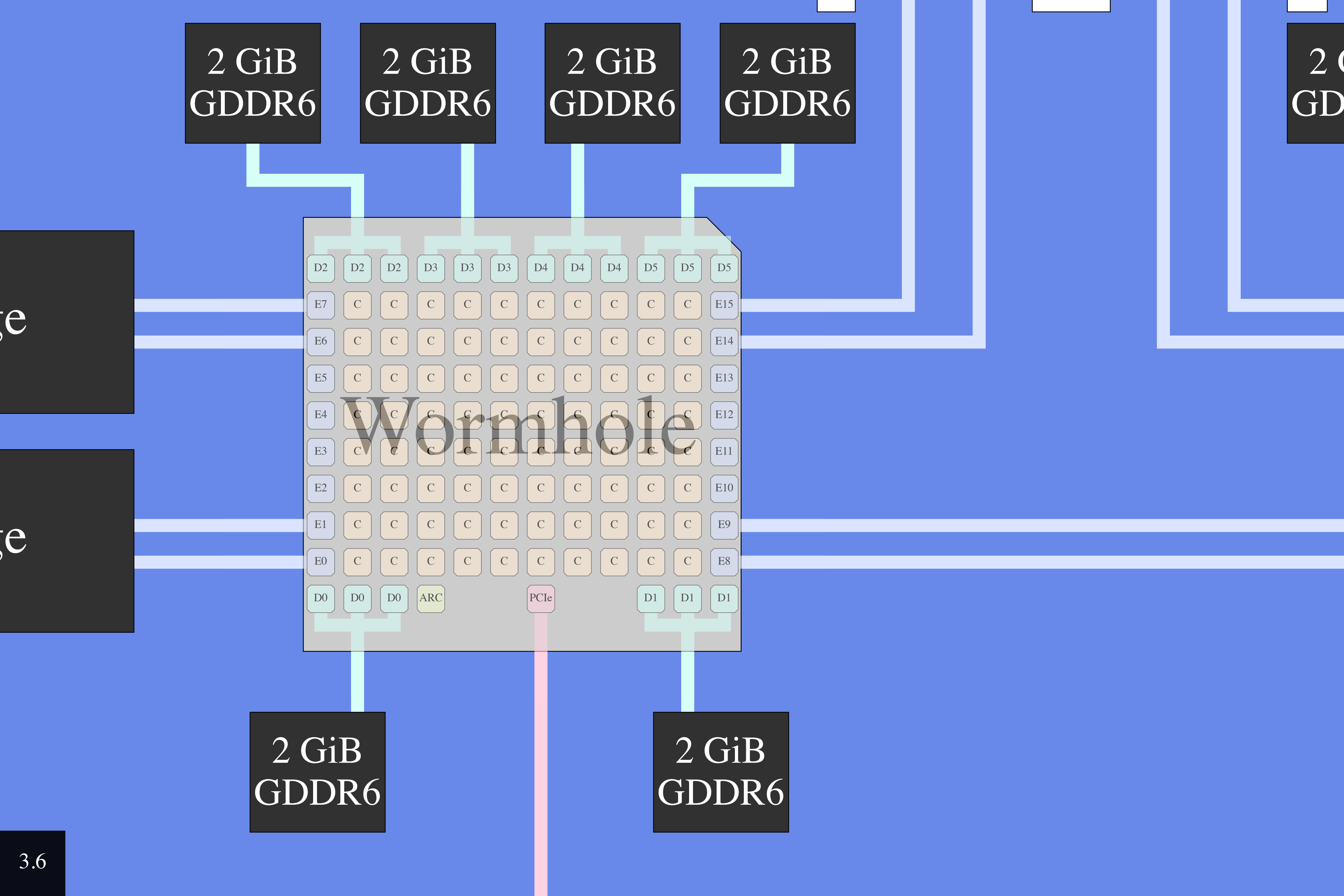


2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6





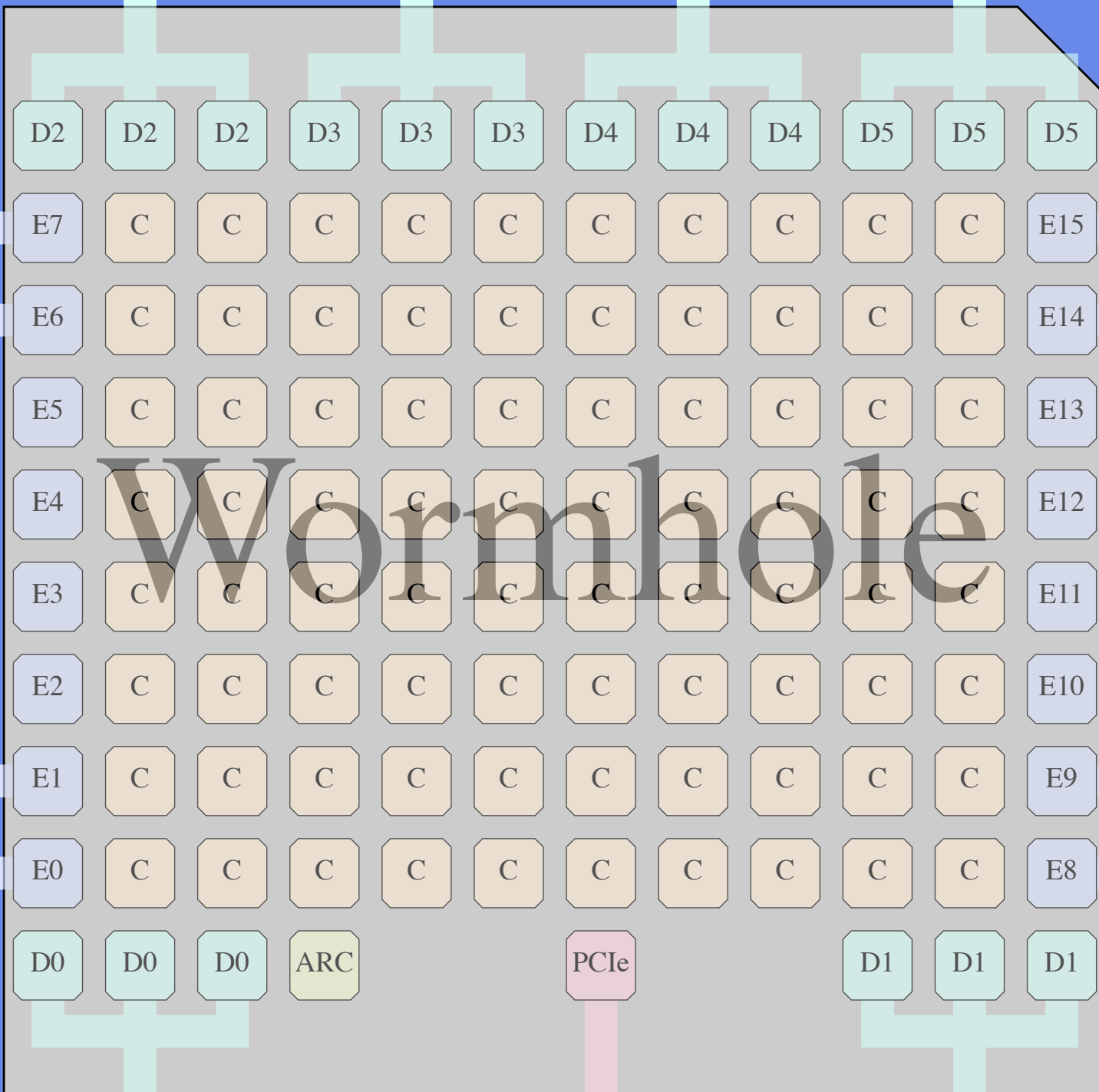
2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6

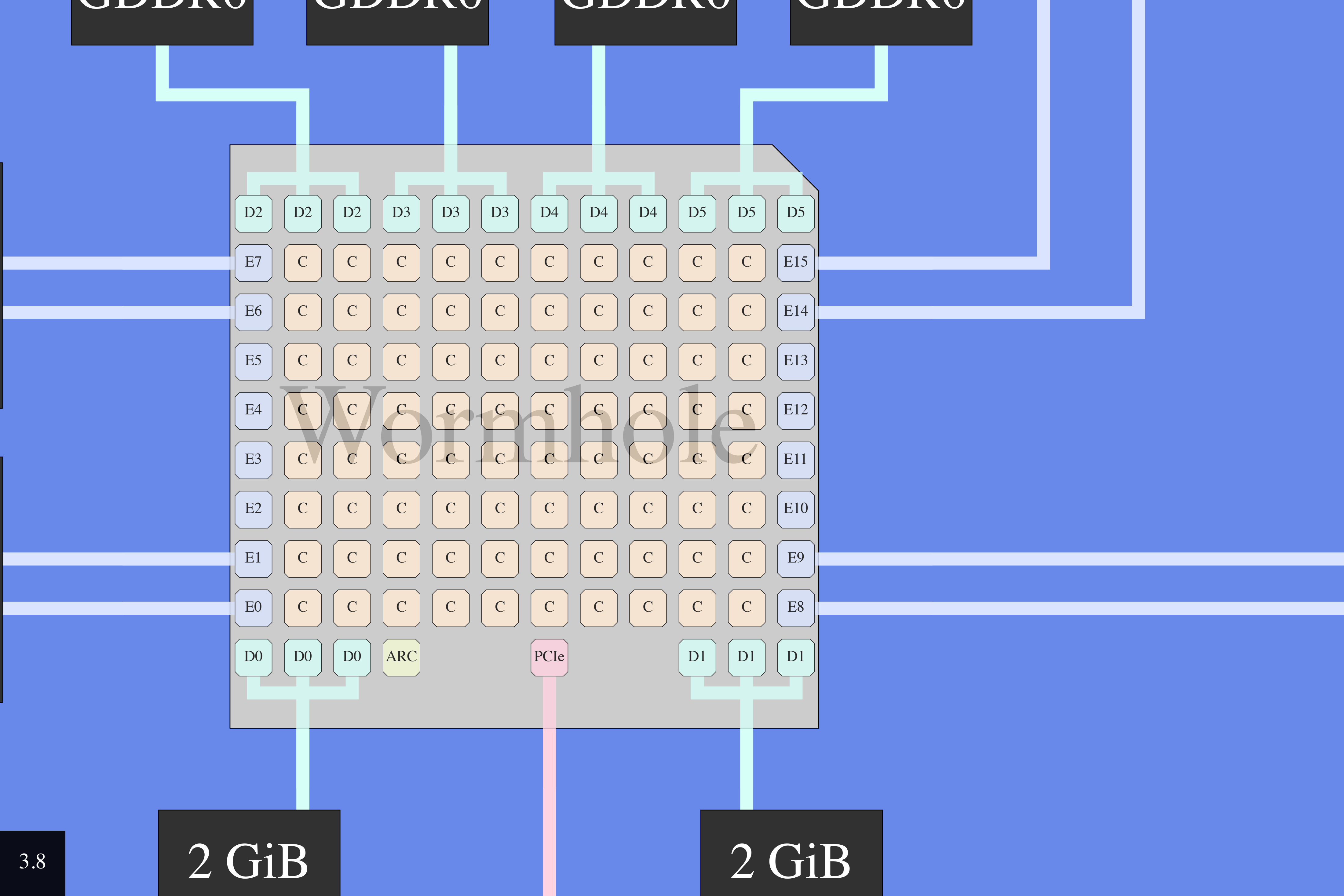
2 GiB
GDDR6

2 GiB
GDDR6



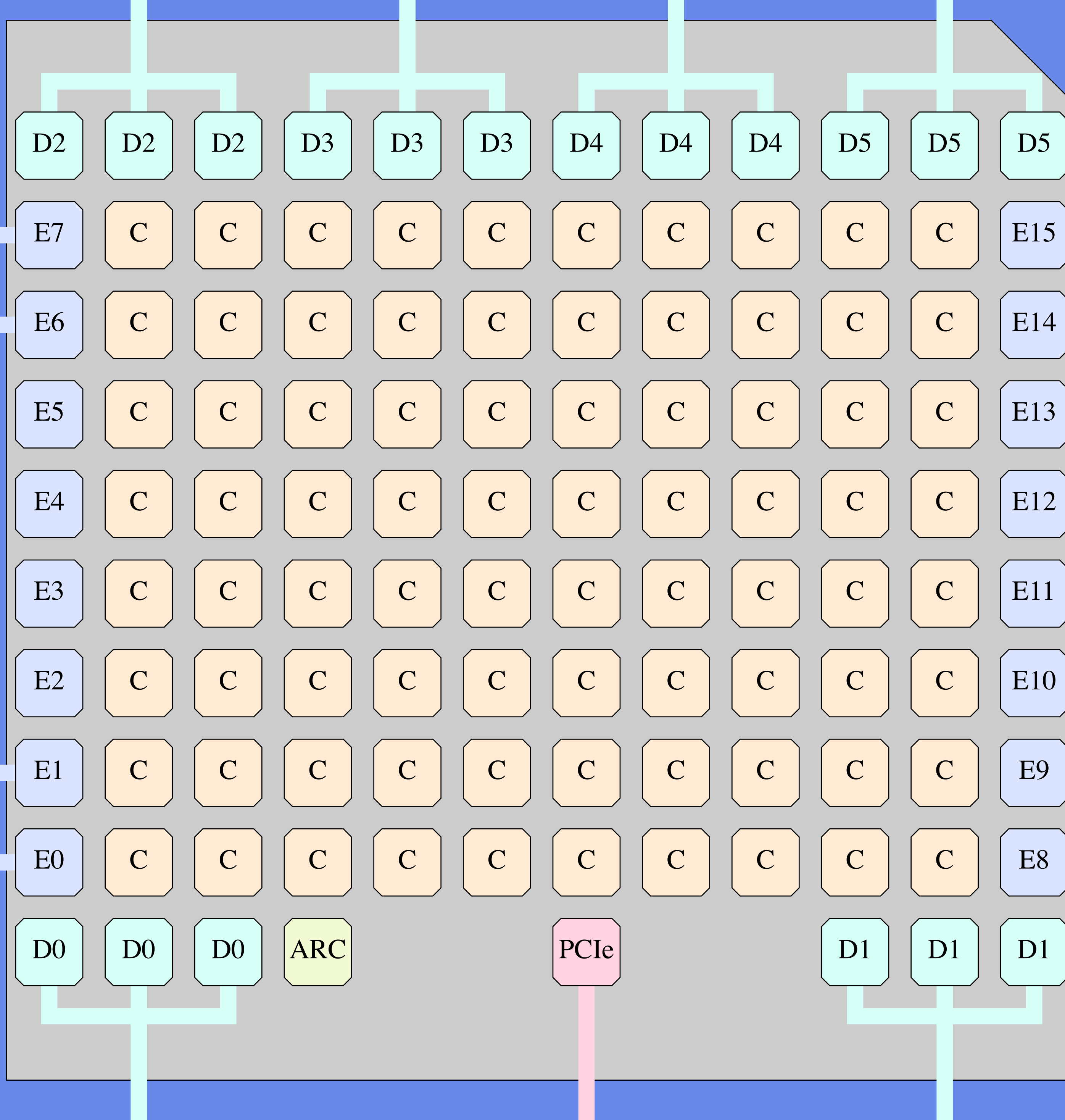
2 GiB
GDDR6

2 GiB
GDDR6



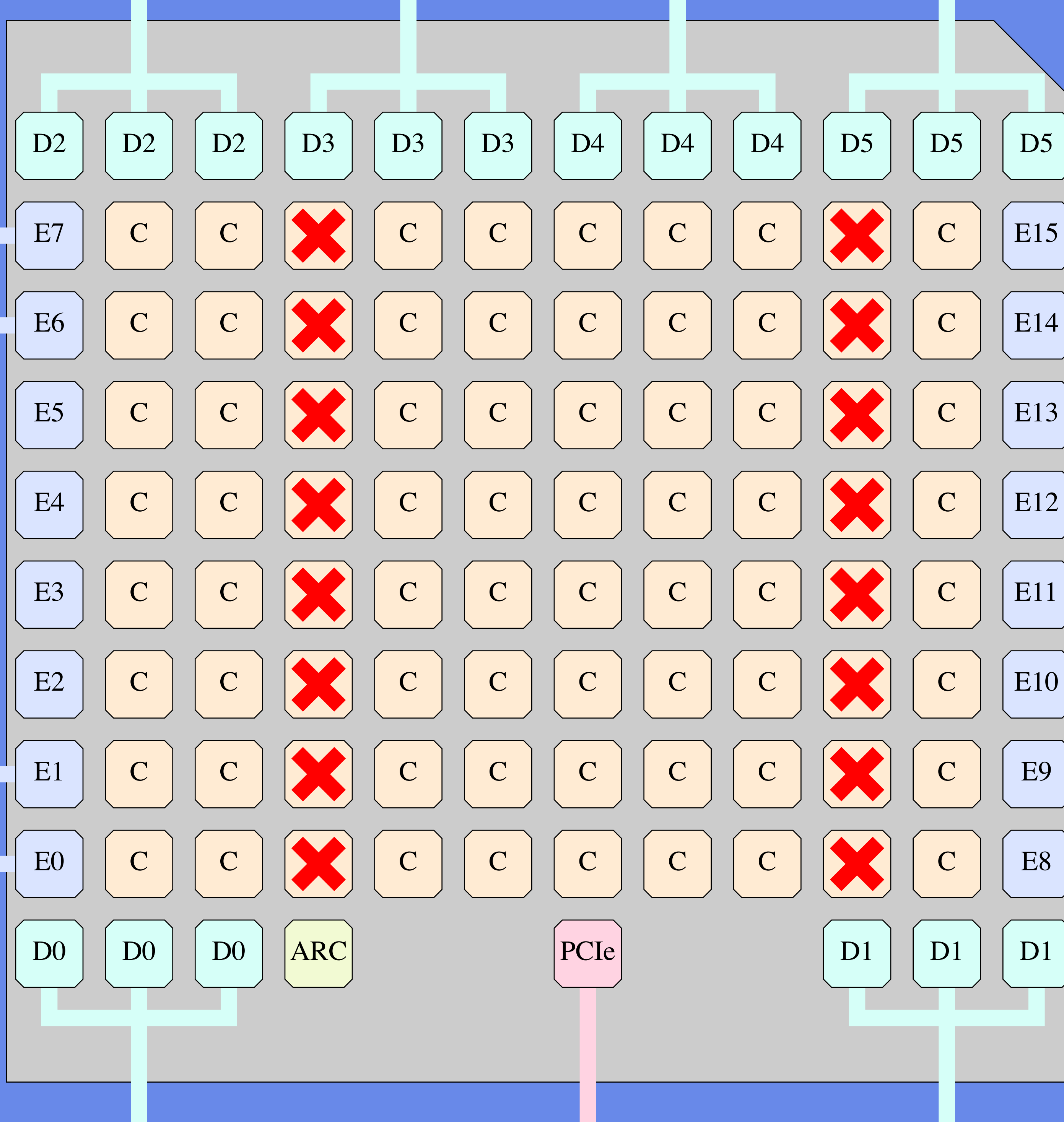
2 GiB

2 GiB



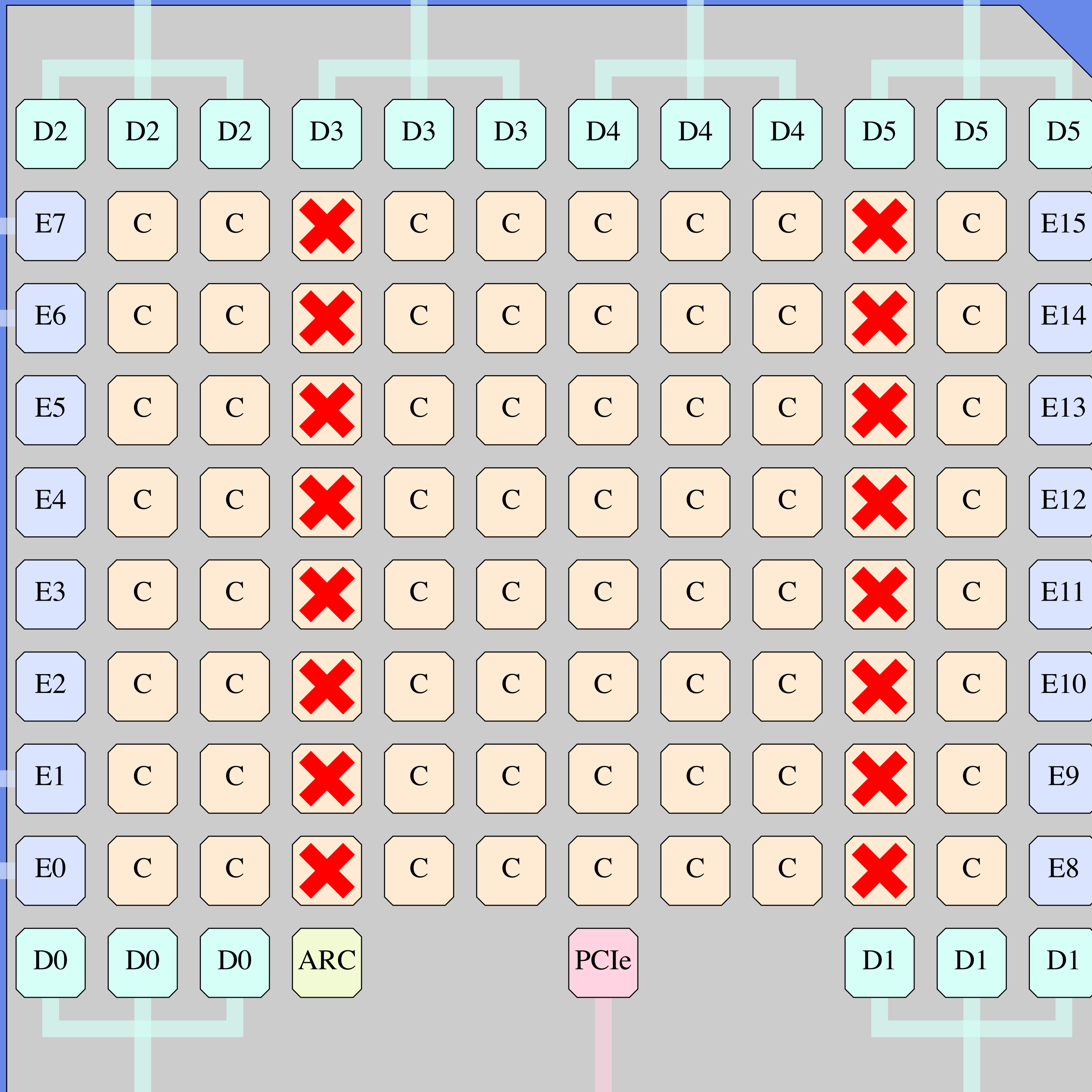
Five types of tile:

- ARC (1)
- C: Compute (80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



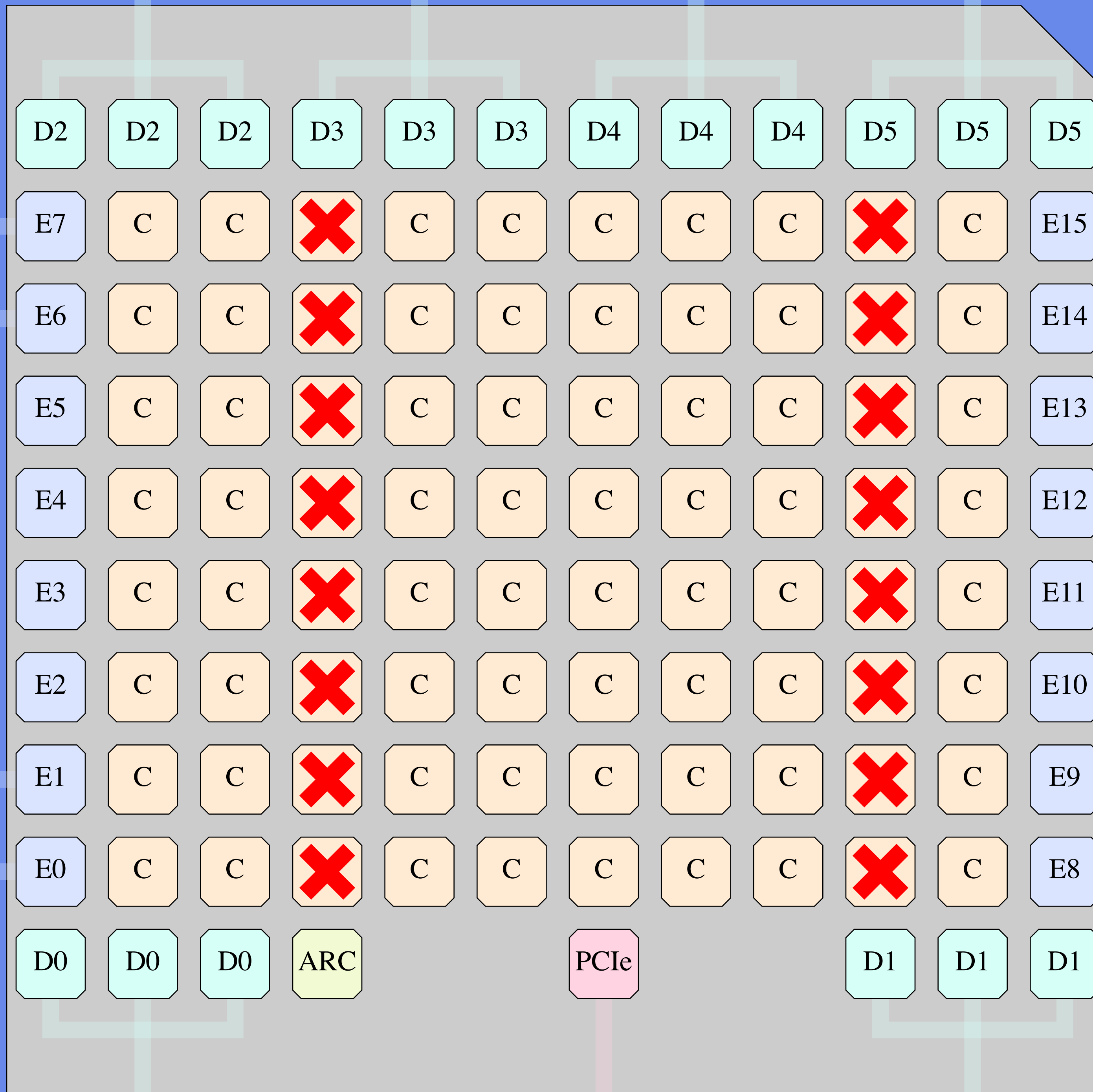
Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



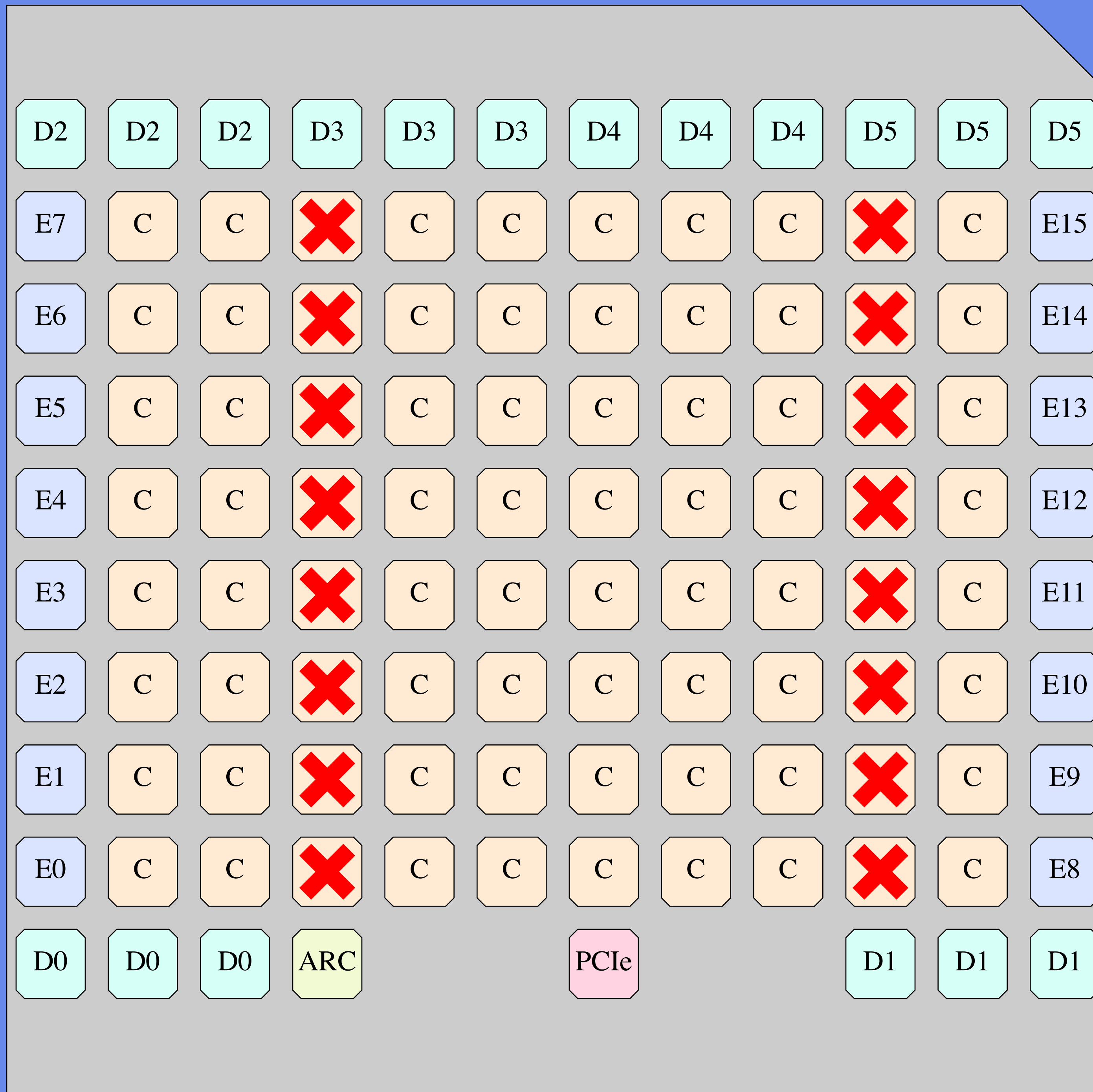
Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



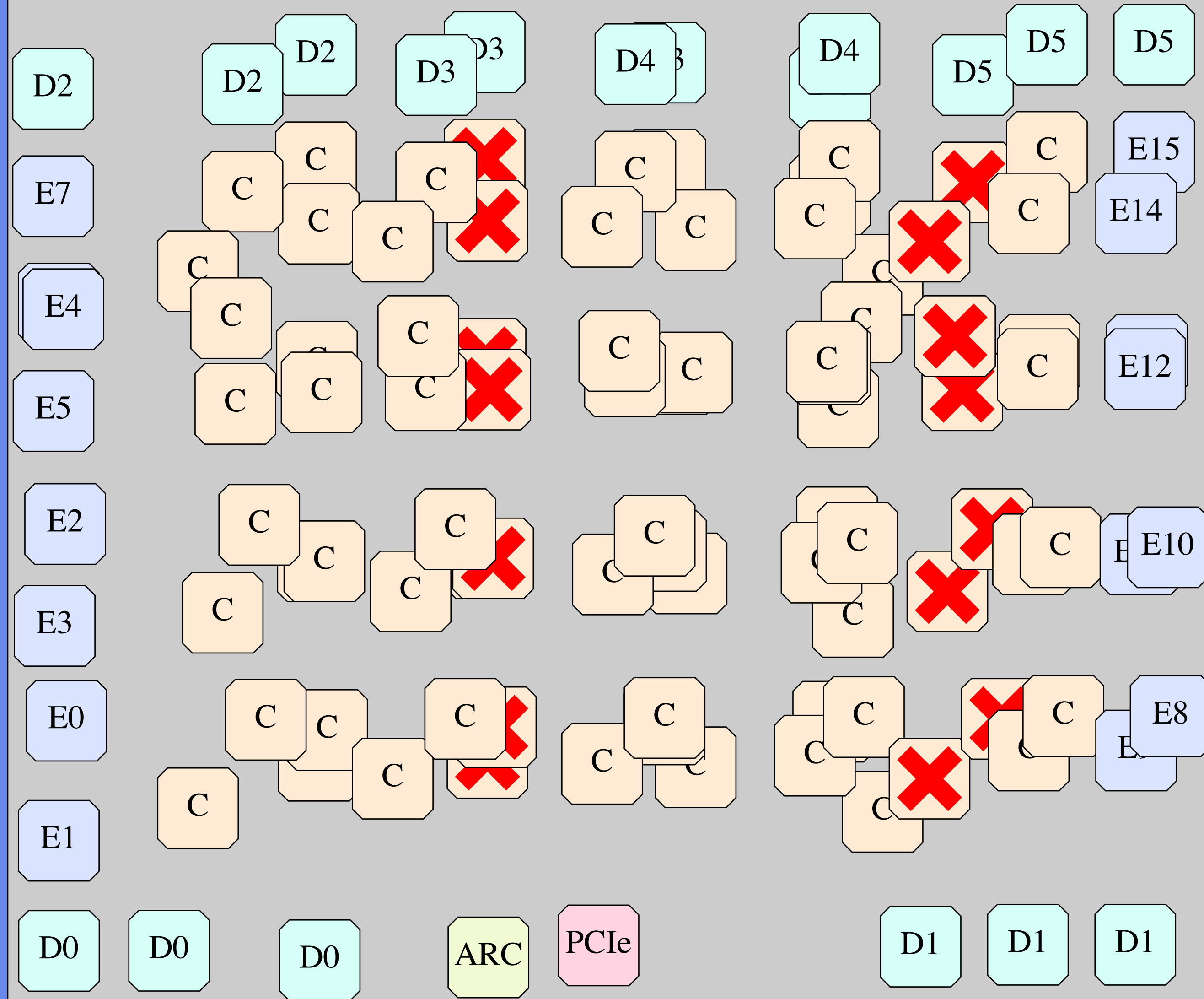
Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



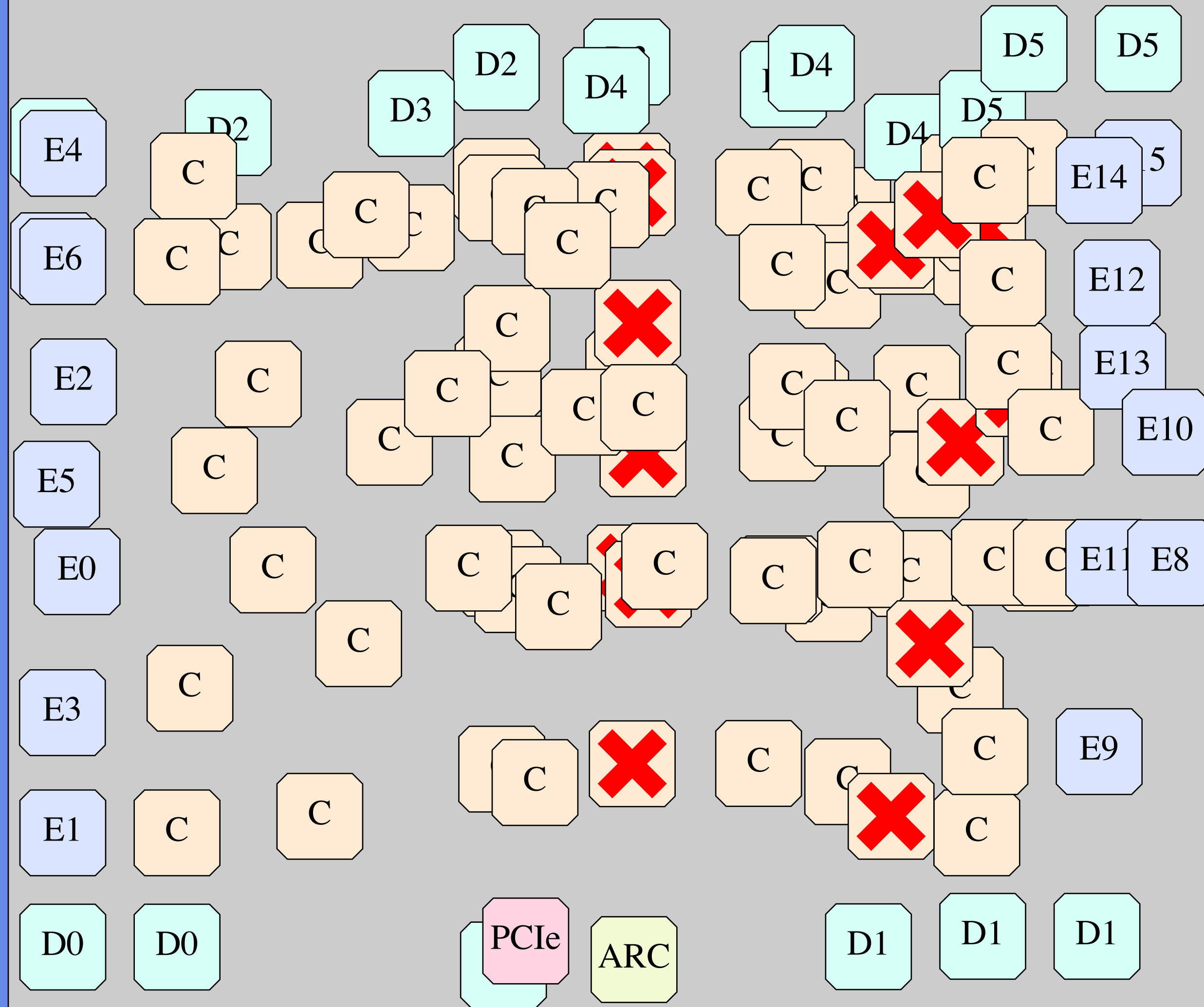
Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



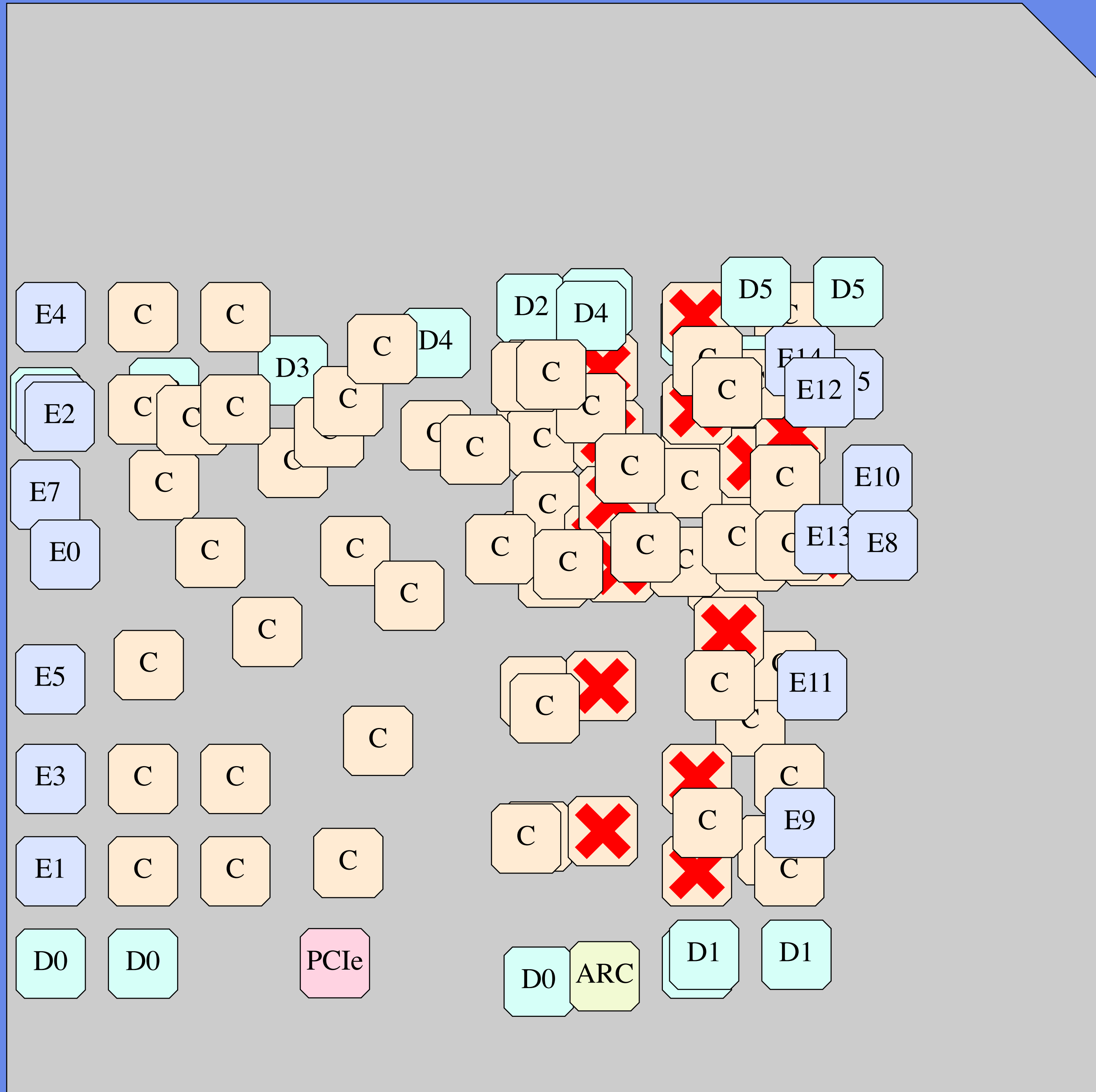
Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)

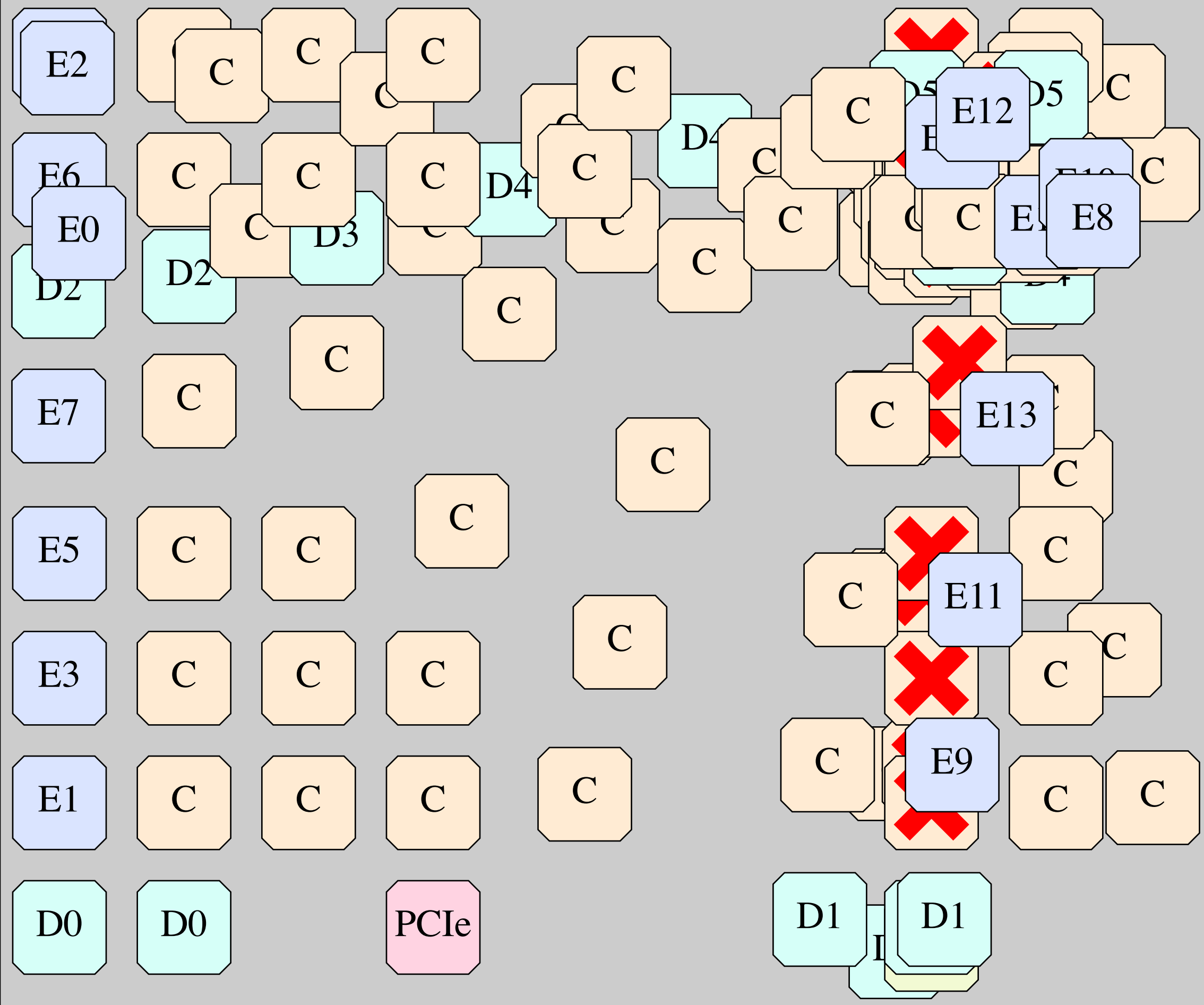


Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)

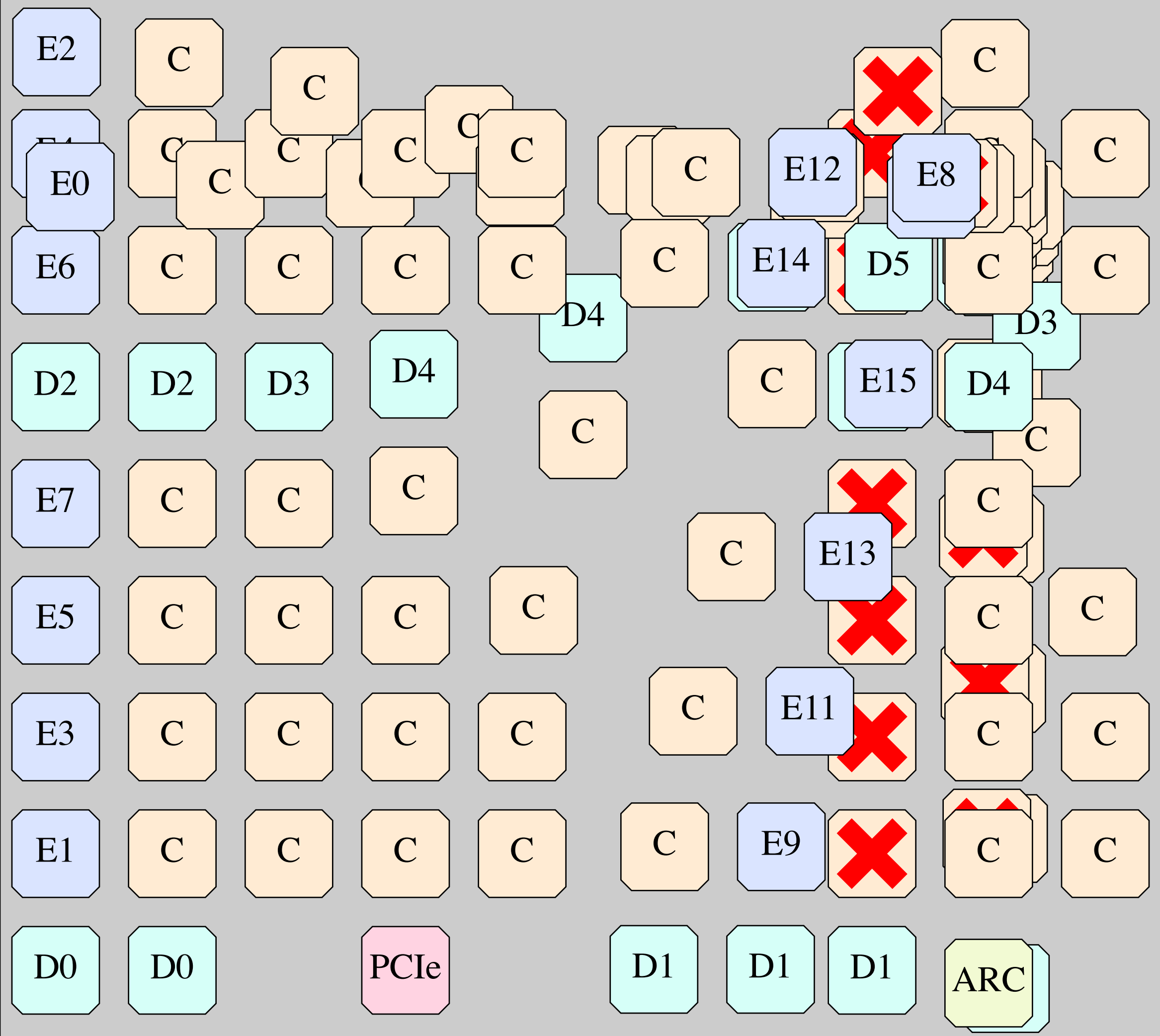
Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



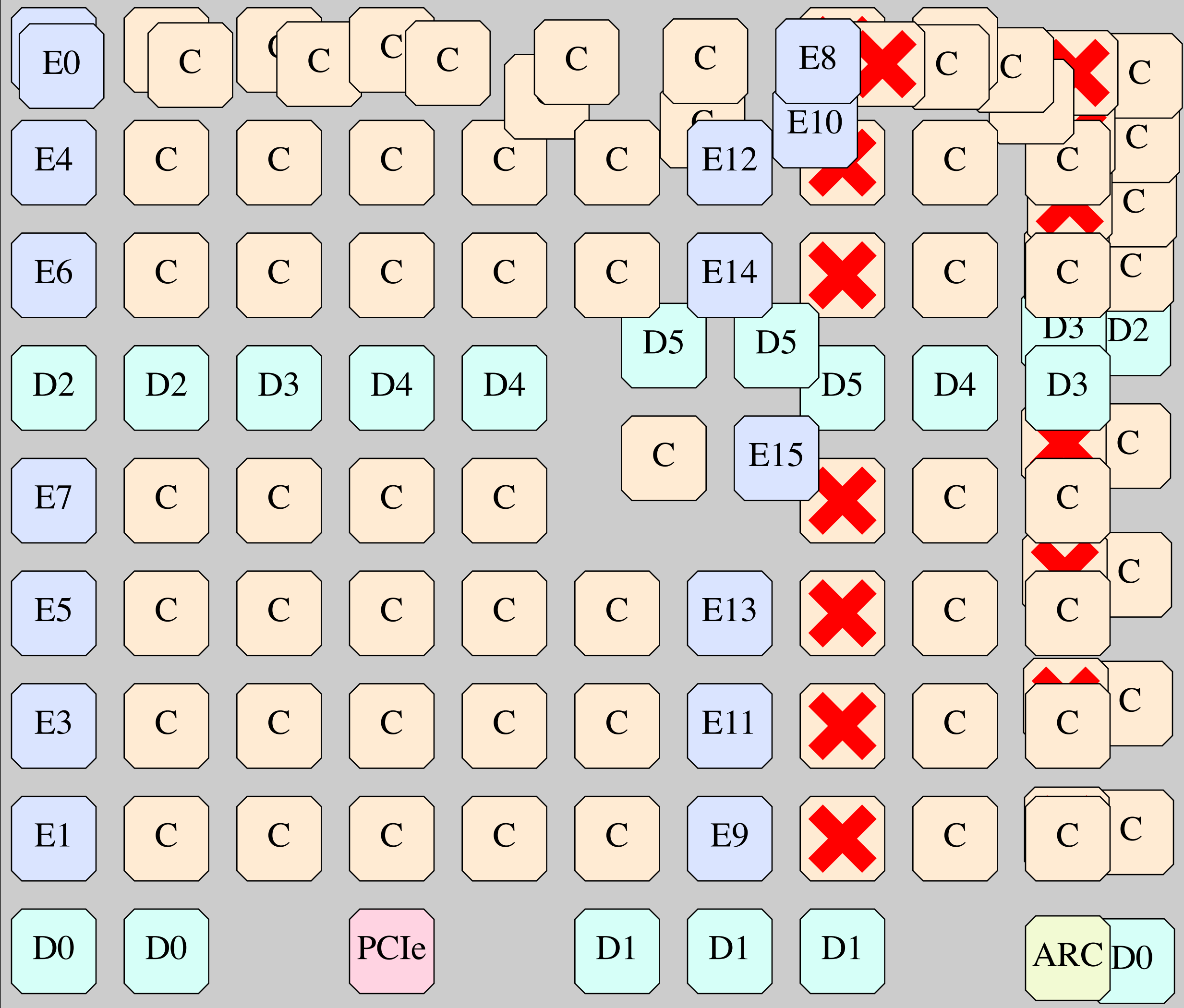
Five types of tile:

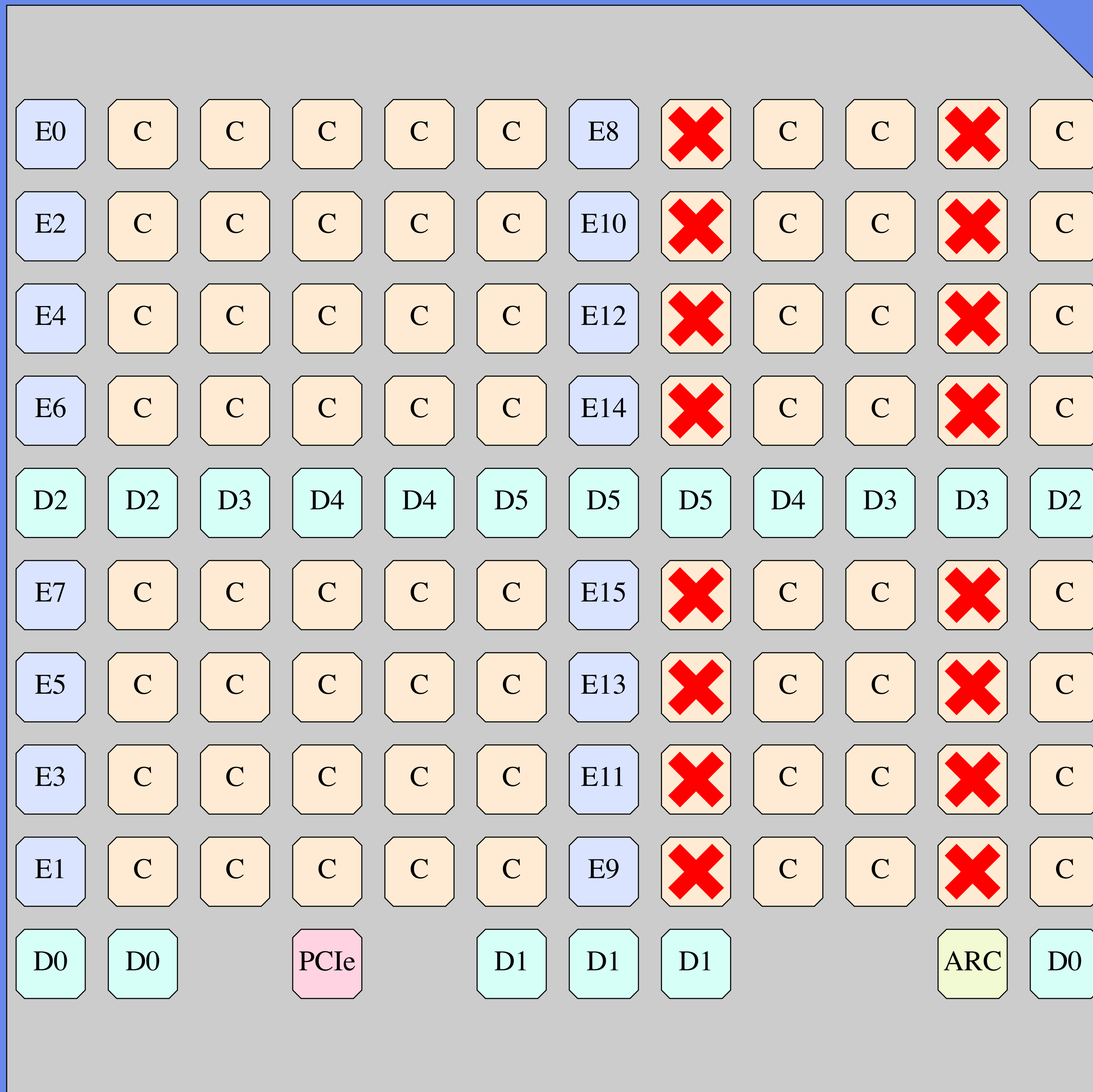
- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



Five types of tile:

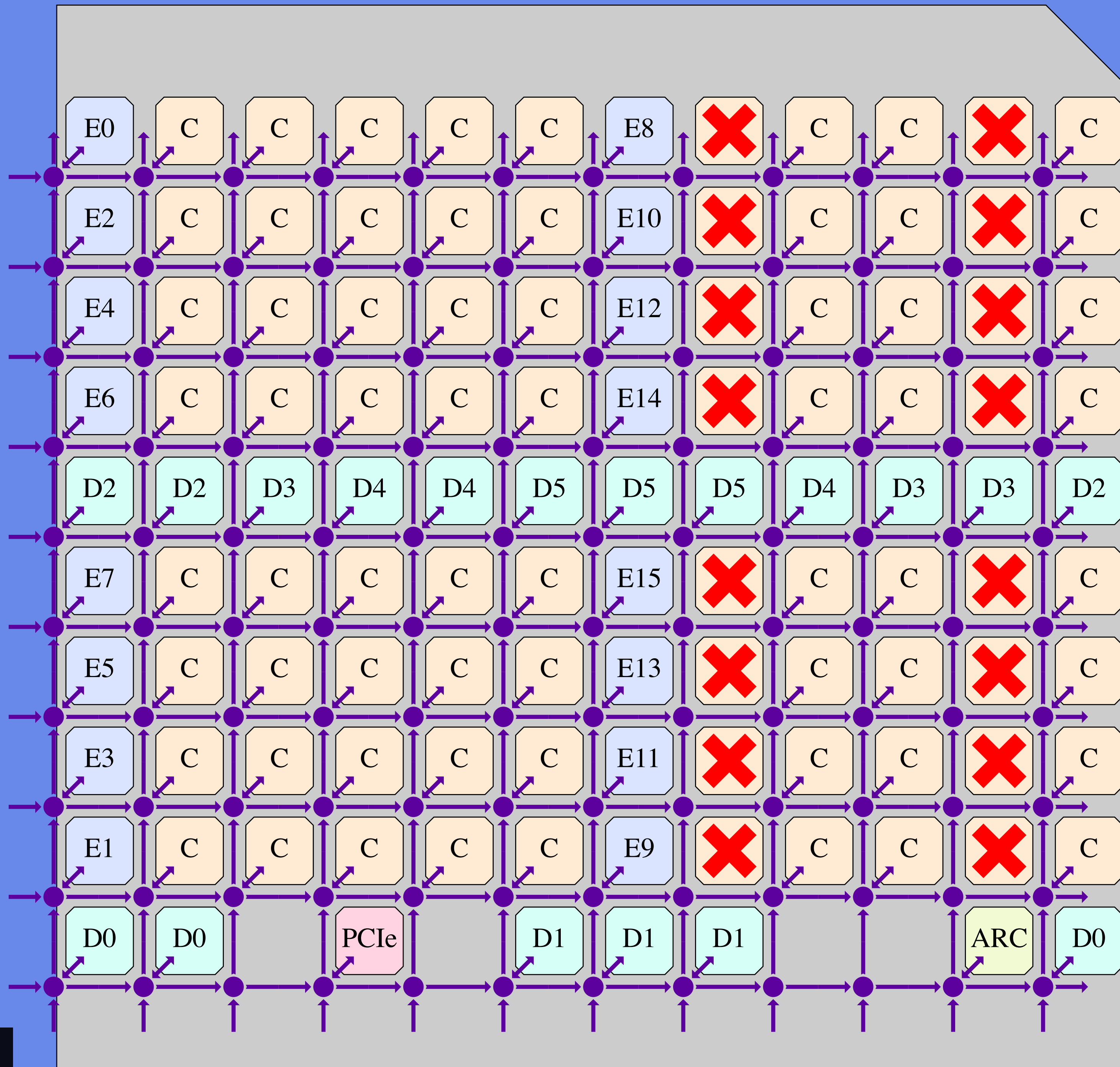
- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)





Five types of tile:

- ARC (1)
- C: Compute (64 or 72 or 80)
- D: DRAM (18)
- E: Ethernet (16)
- PCIe (1)



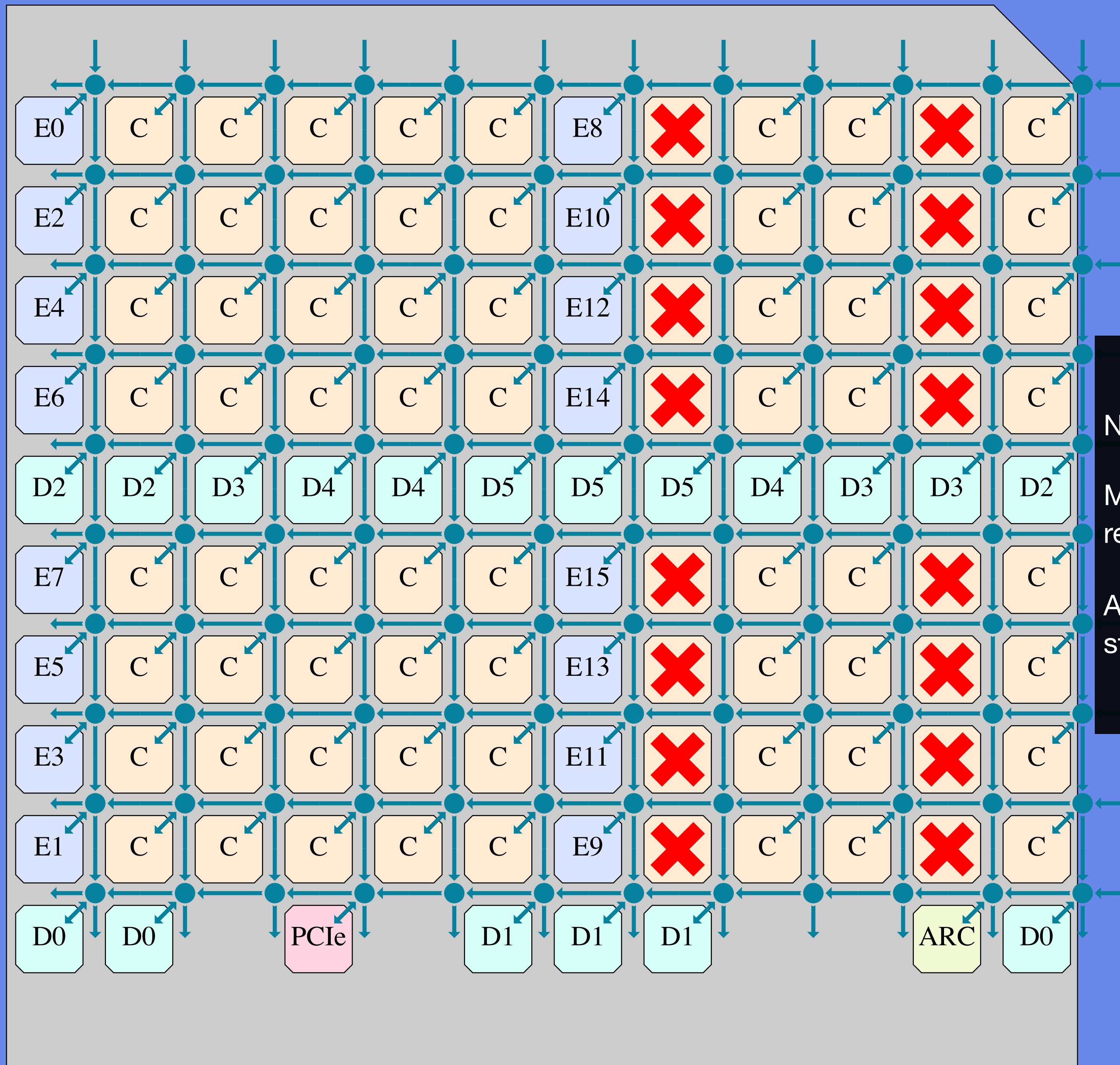
NoC does asynchronous memcopy:

```

u4 dst_x, u4 dst_y,
u32 dst_addr,
u4 src_x, u4 src_y,
u32 src_addr,
u14 num_bytes

```

Each arrow 32 bytes wide

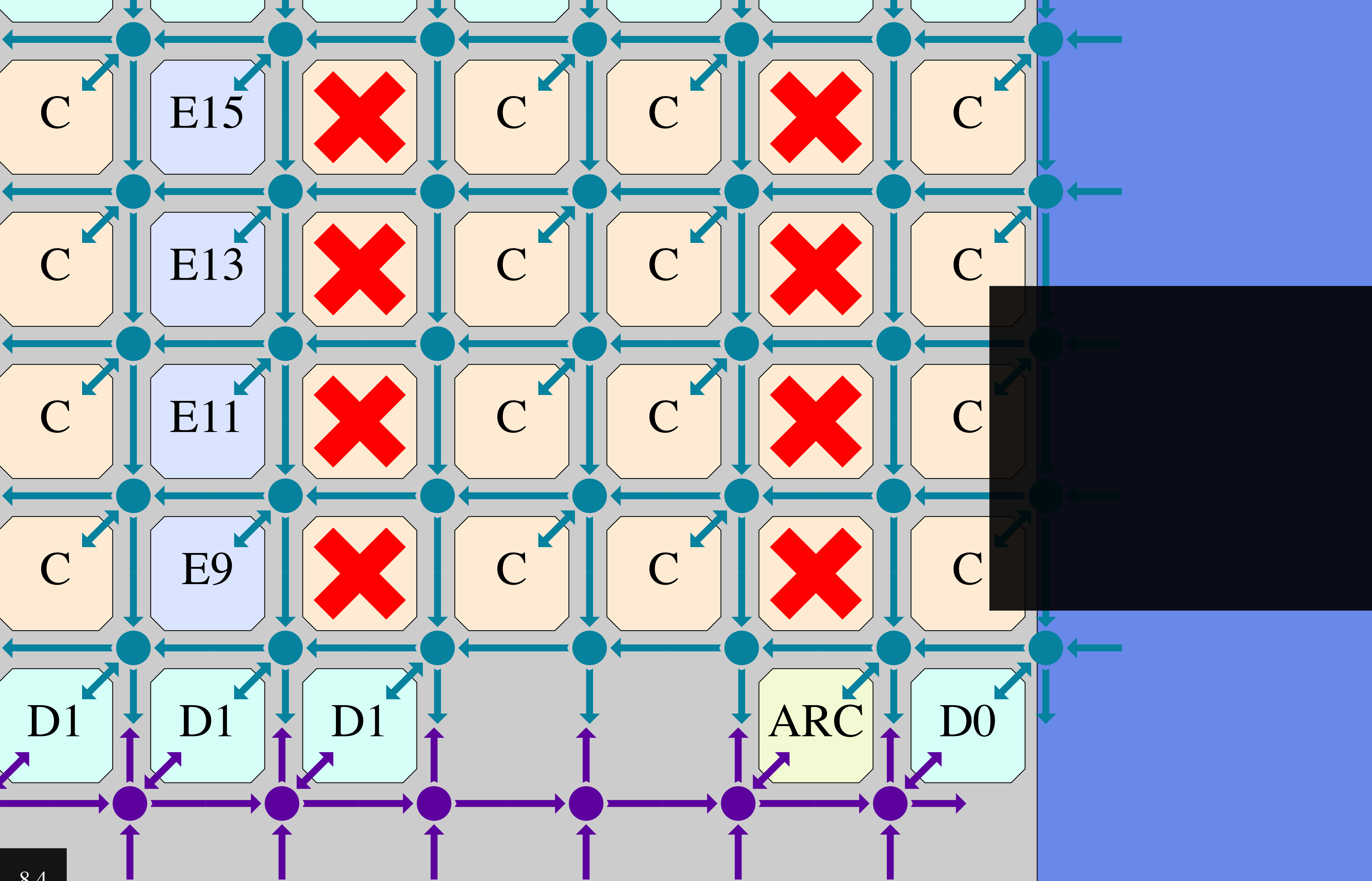


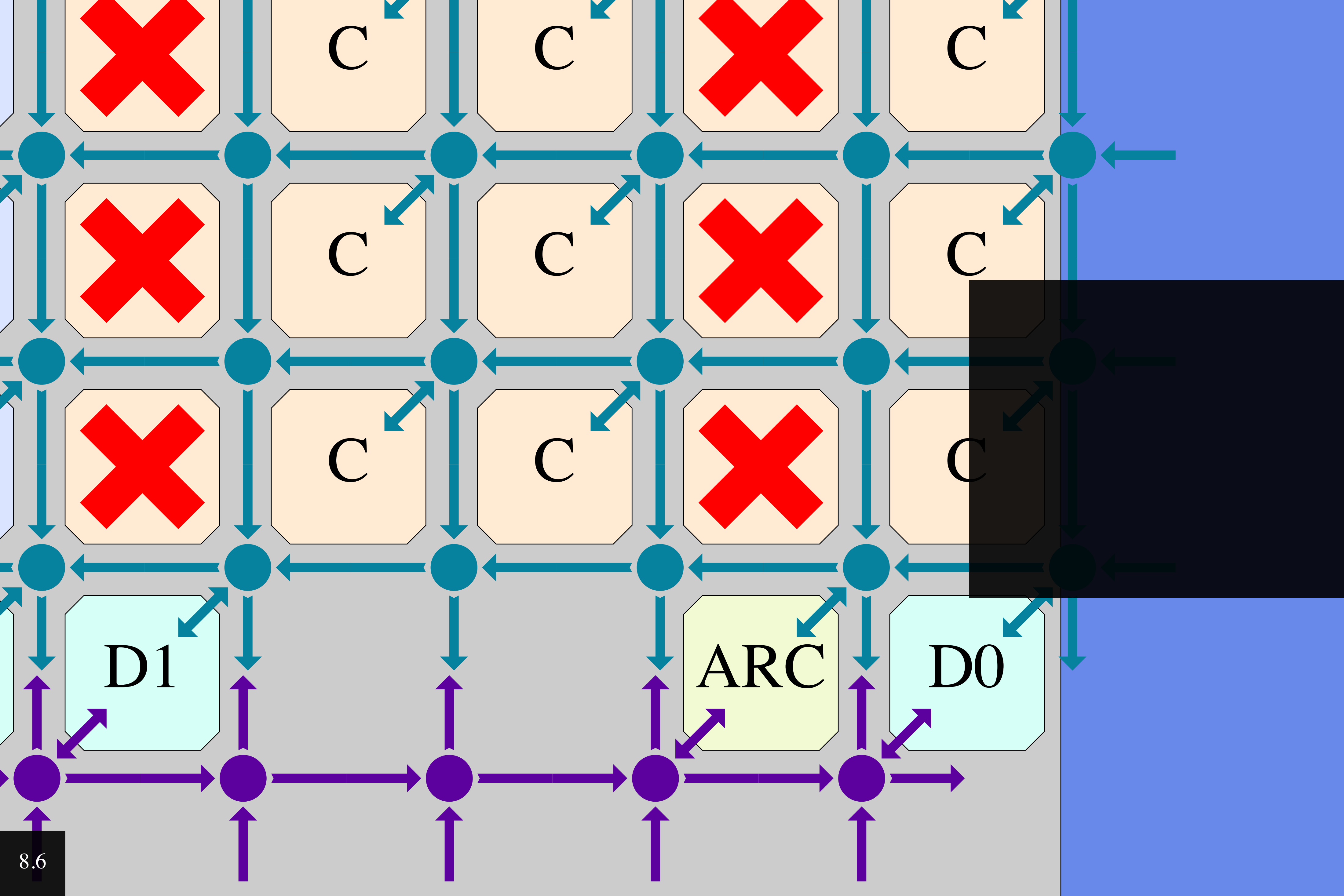
NoC also does:

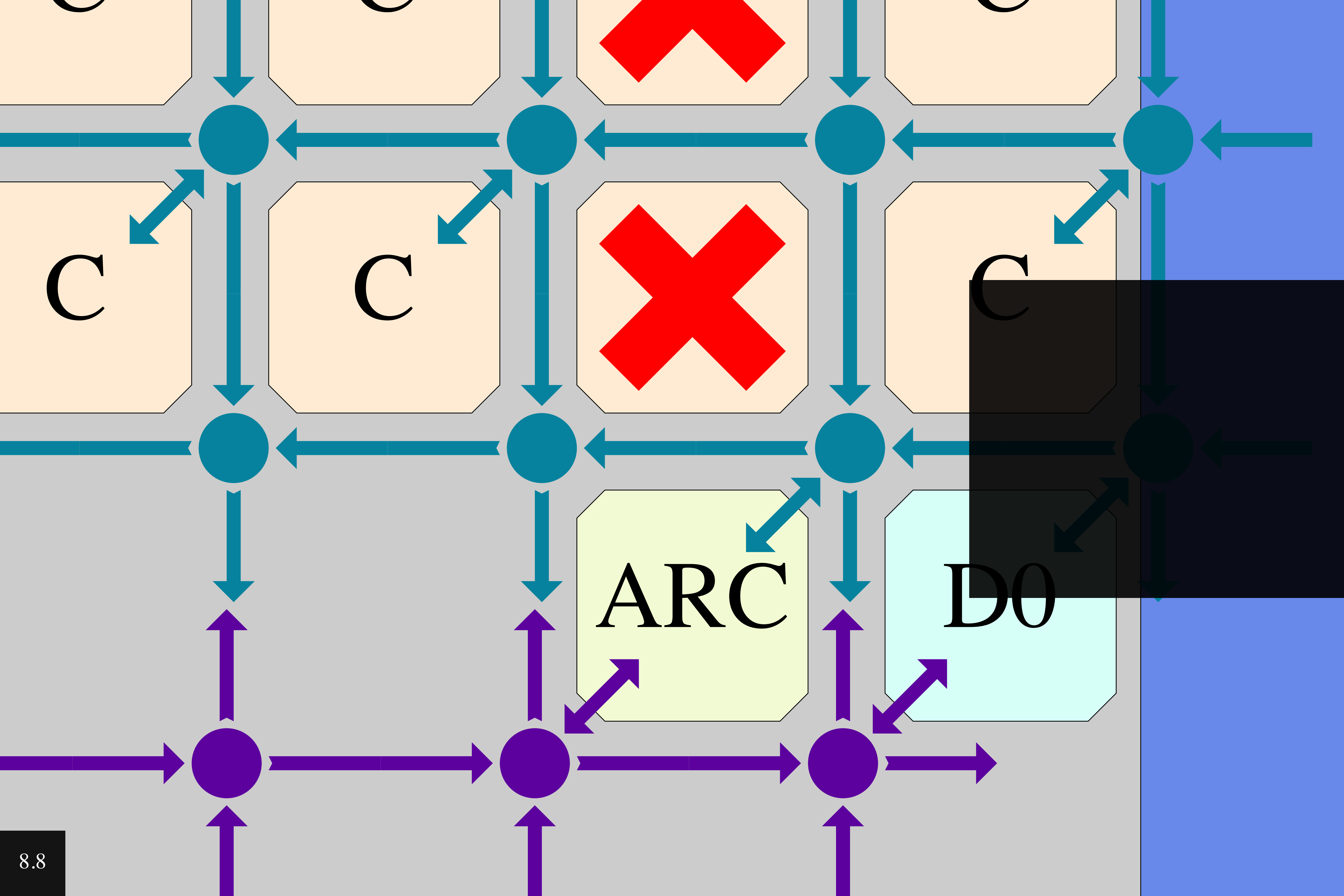
Multicast (dst as an x/y/w/h rectangle rather than x/y point)

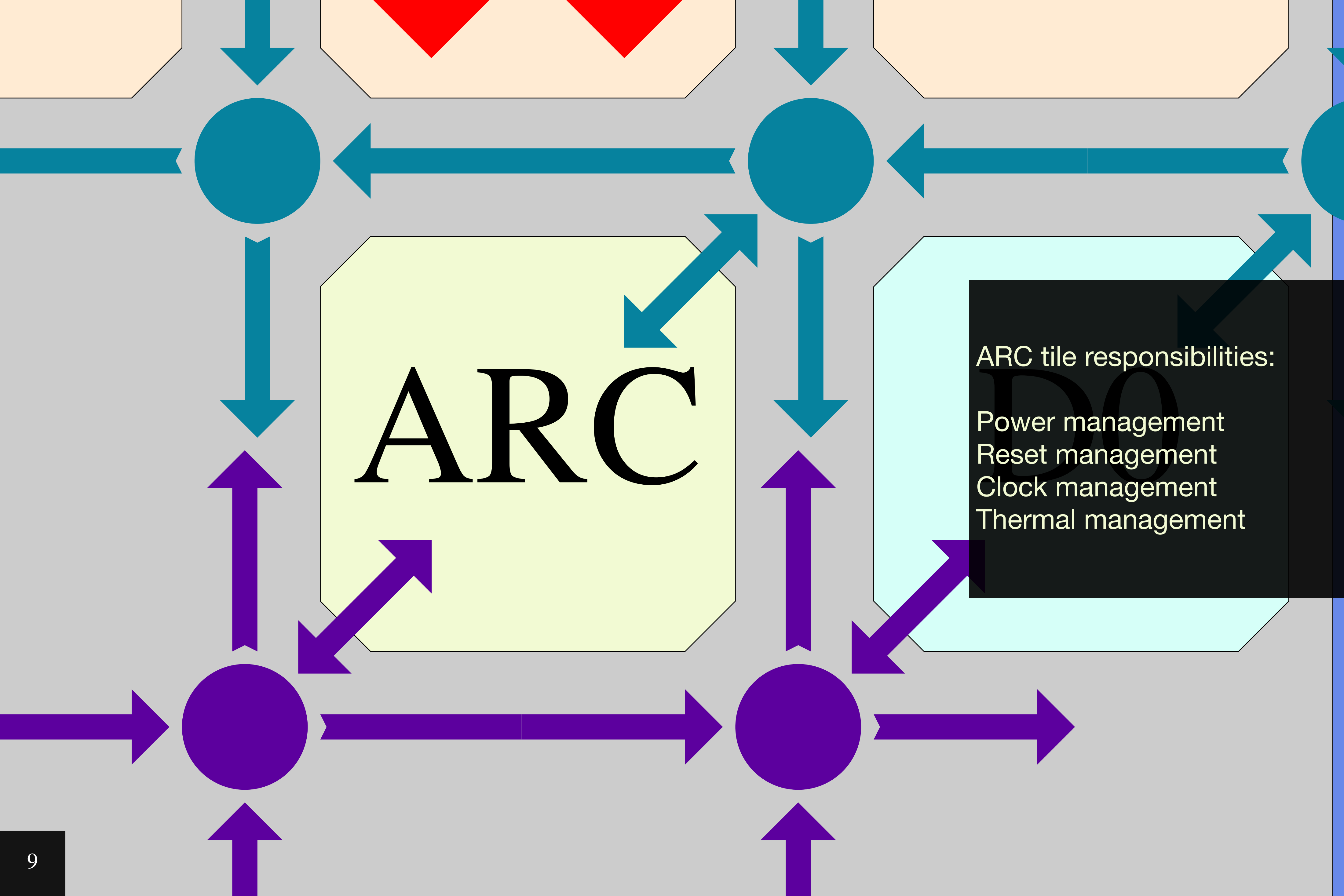
Atomics (e.g. increment at src, store original value to dst)







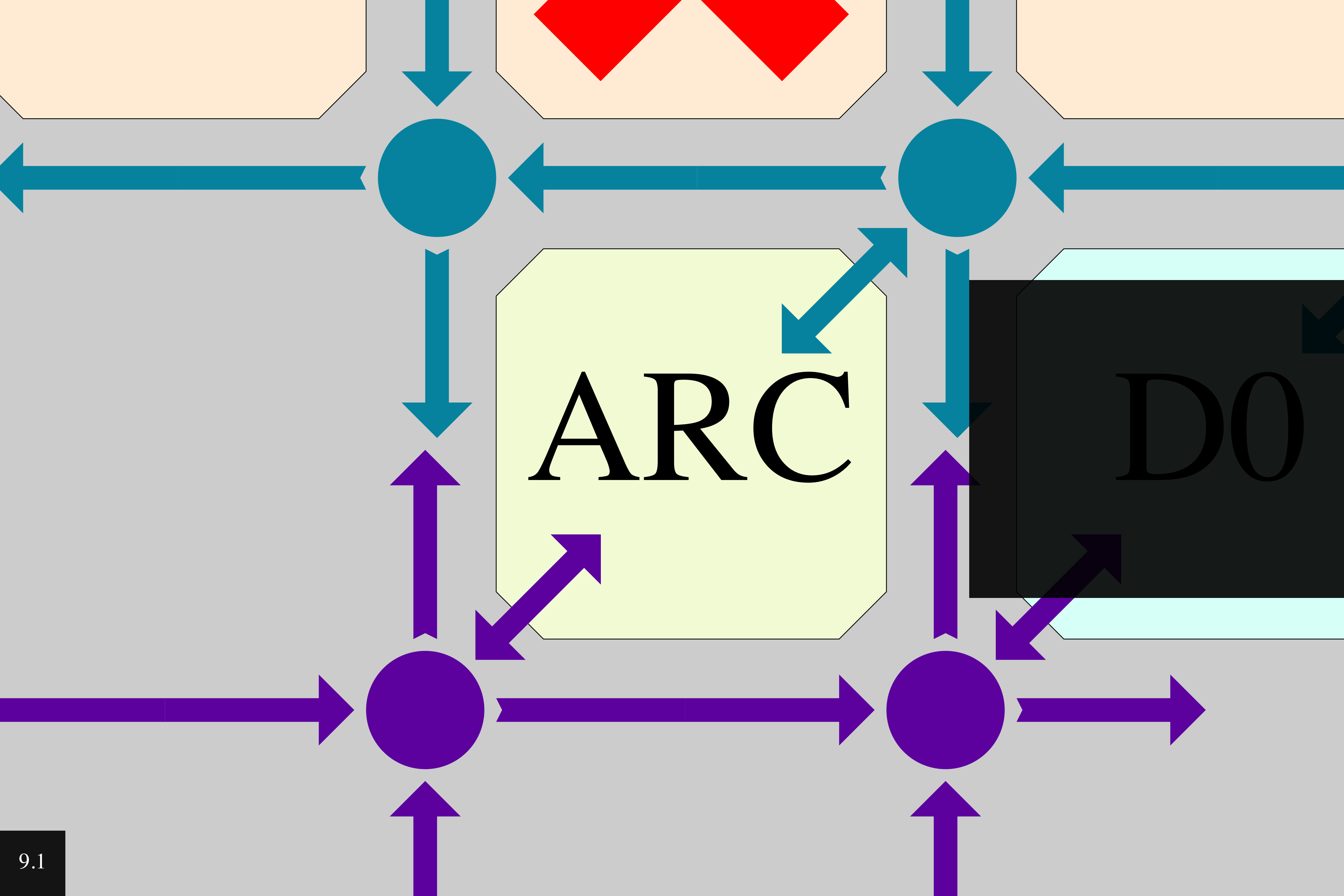


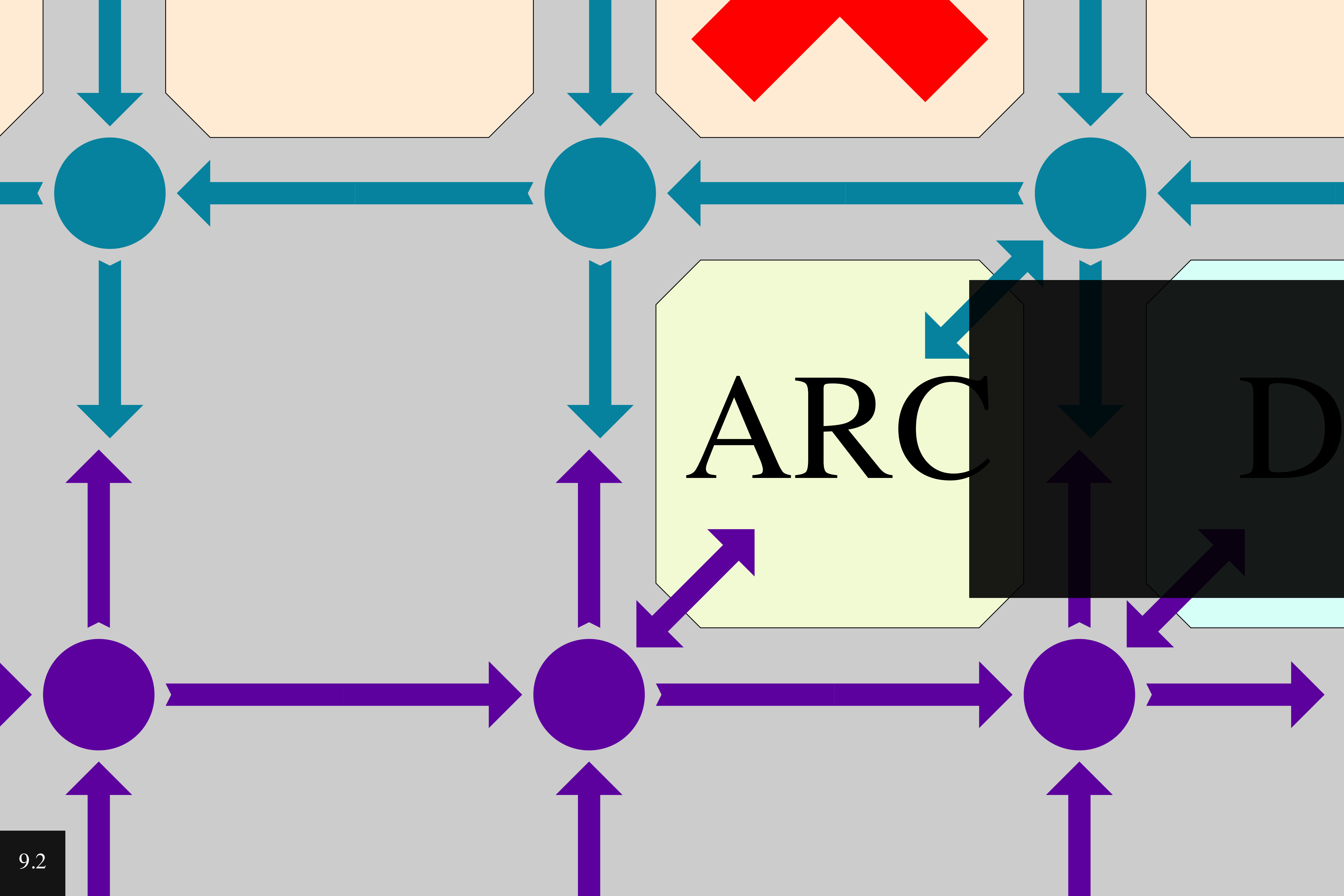


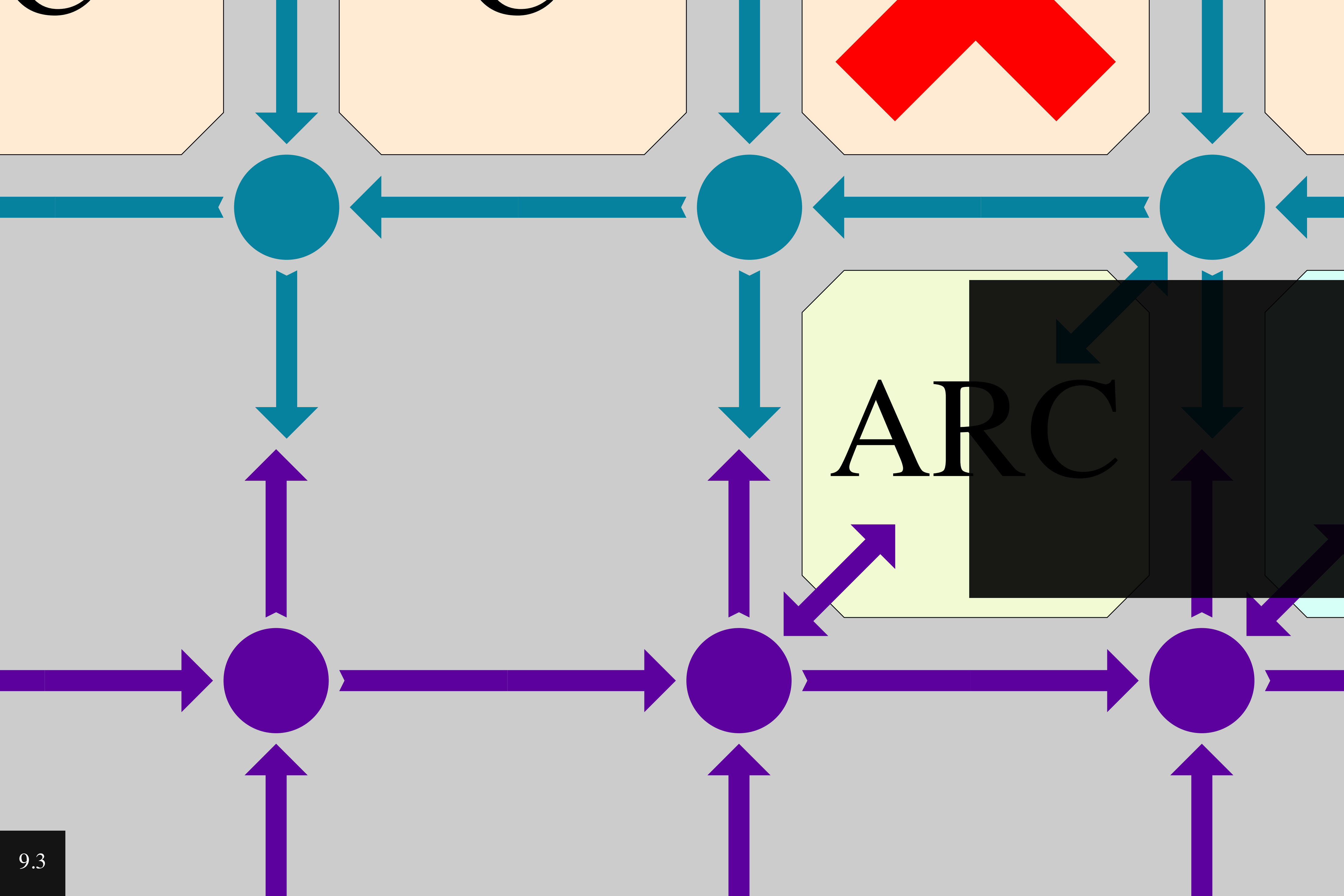
ARC

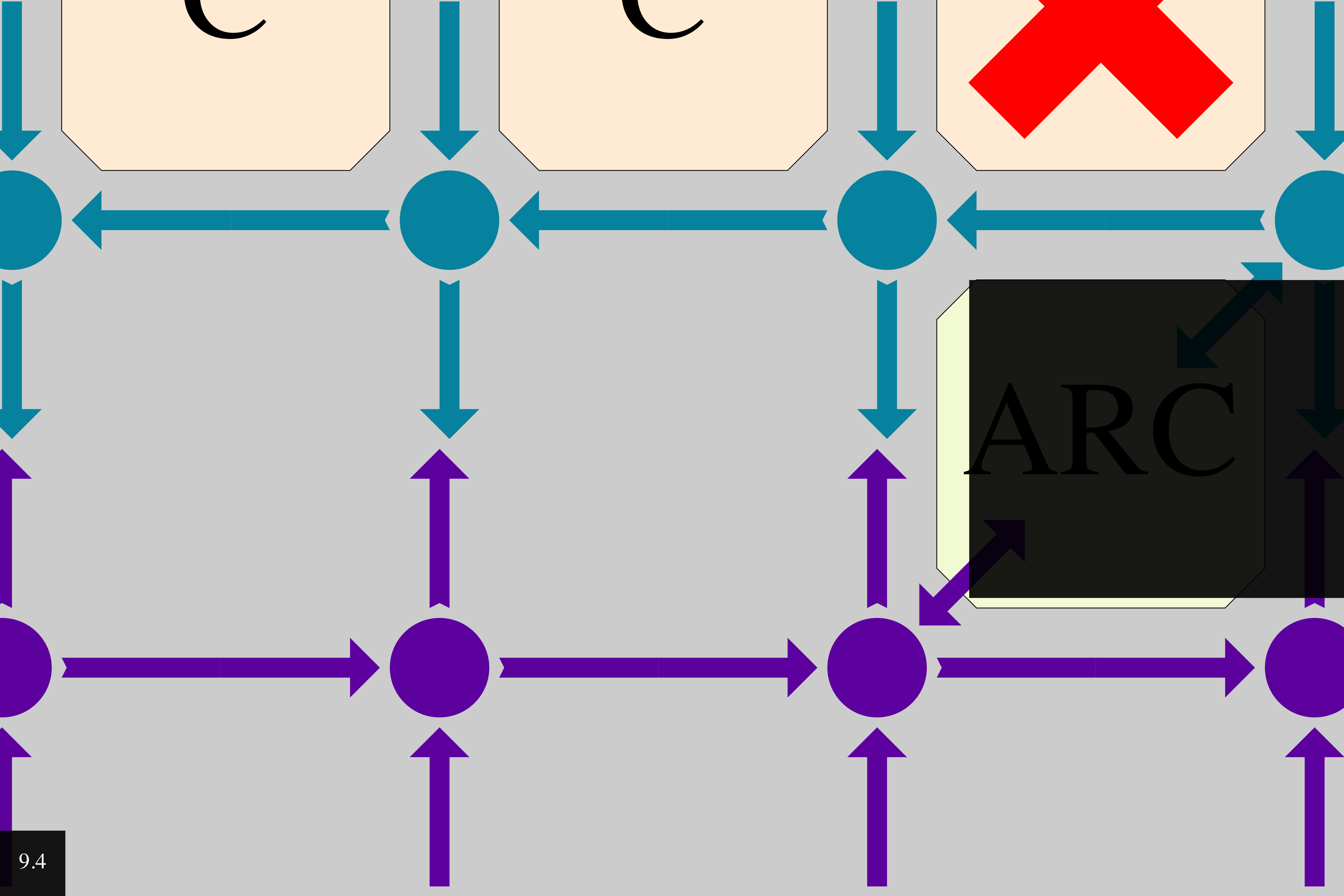
ARC tile responsibilities:

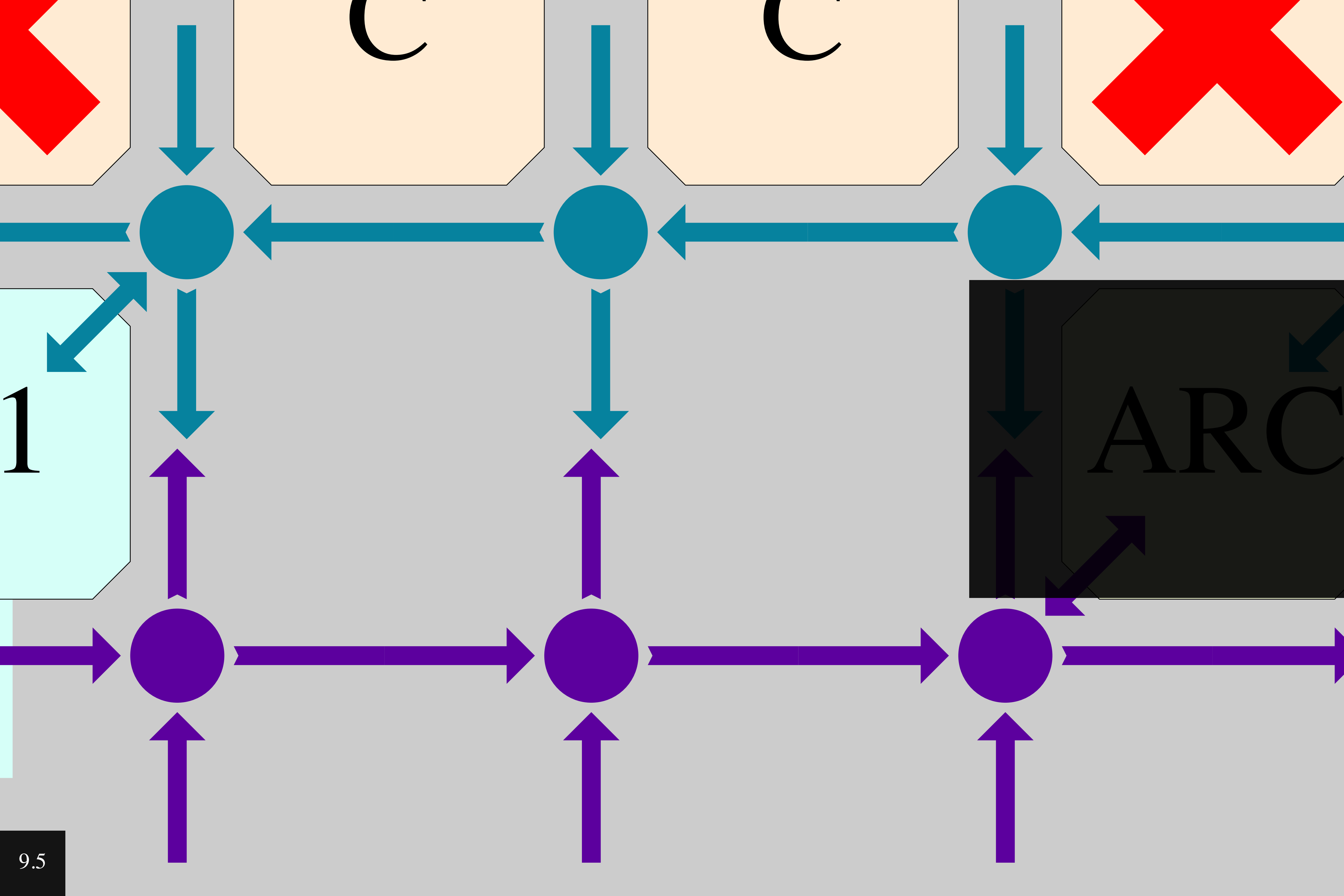
- Power management
- Reset management
- Clock management
- Thermal management





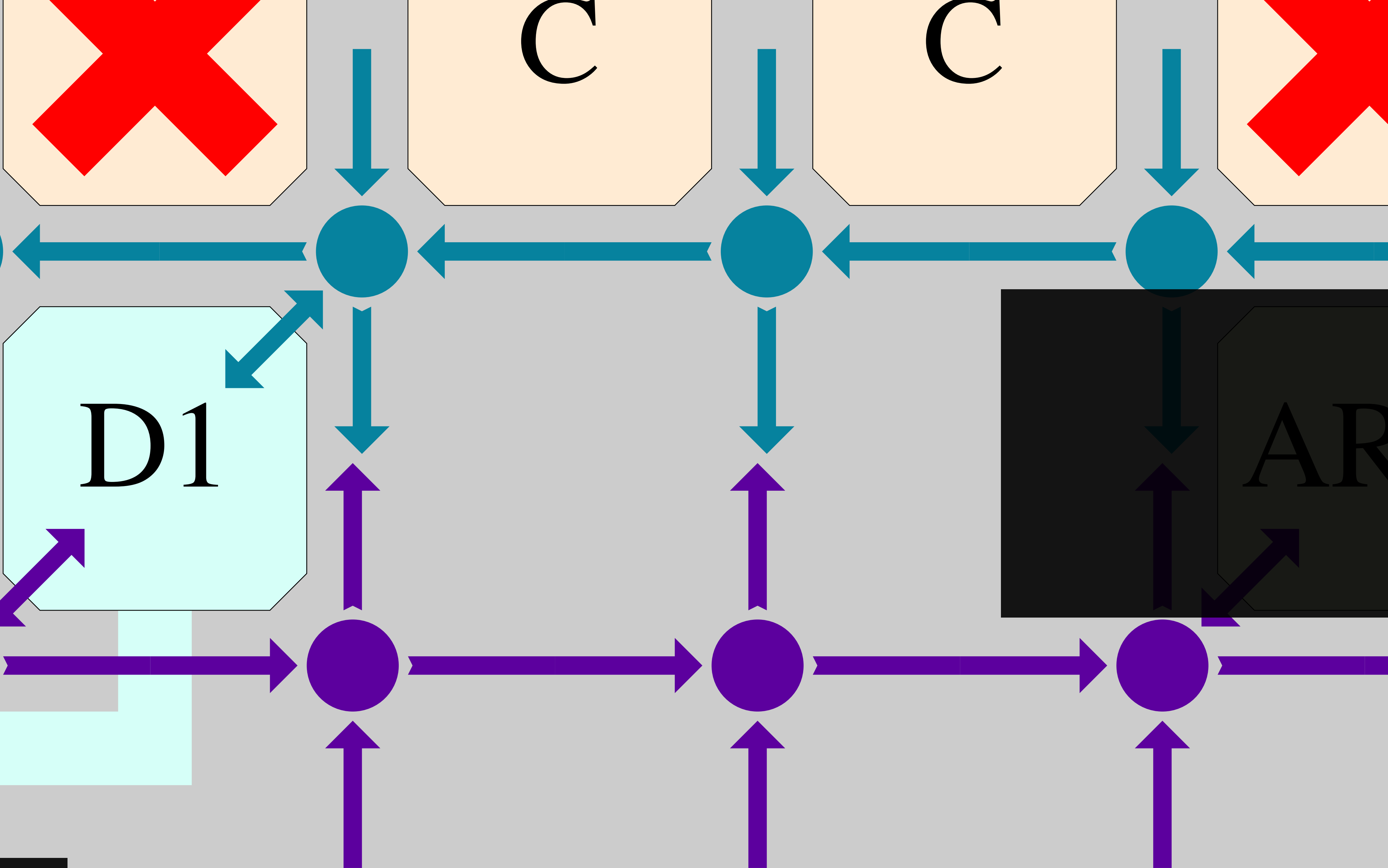


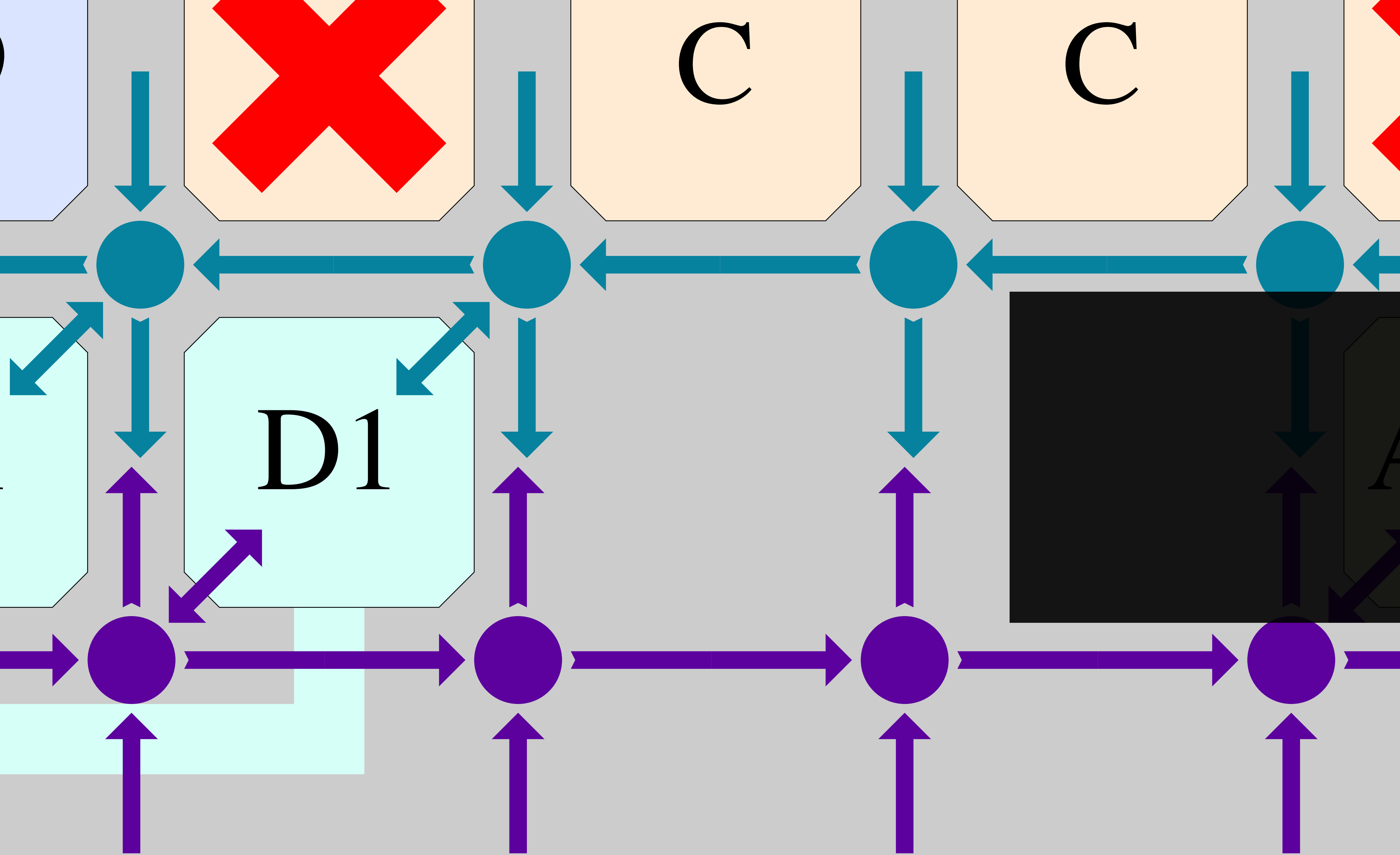


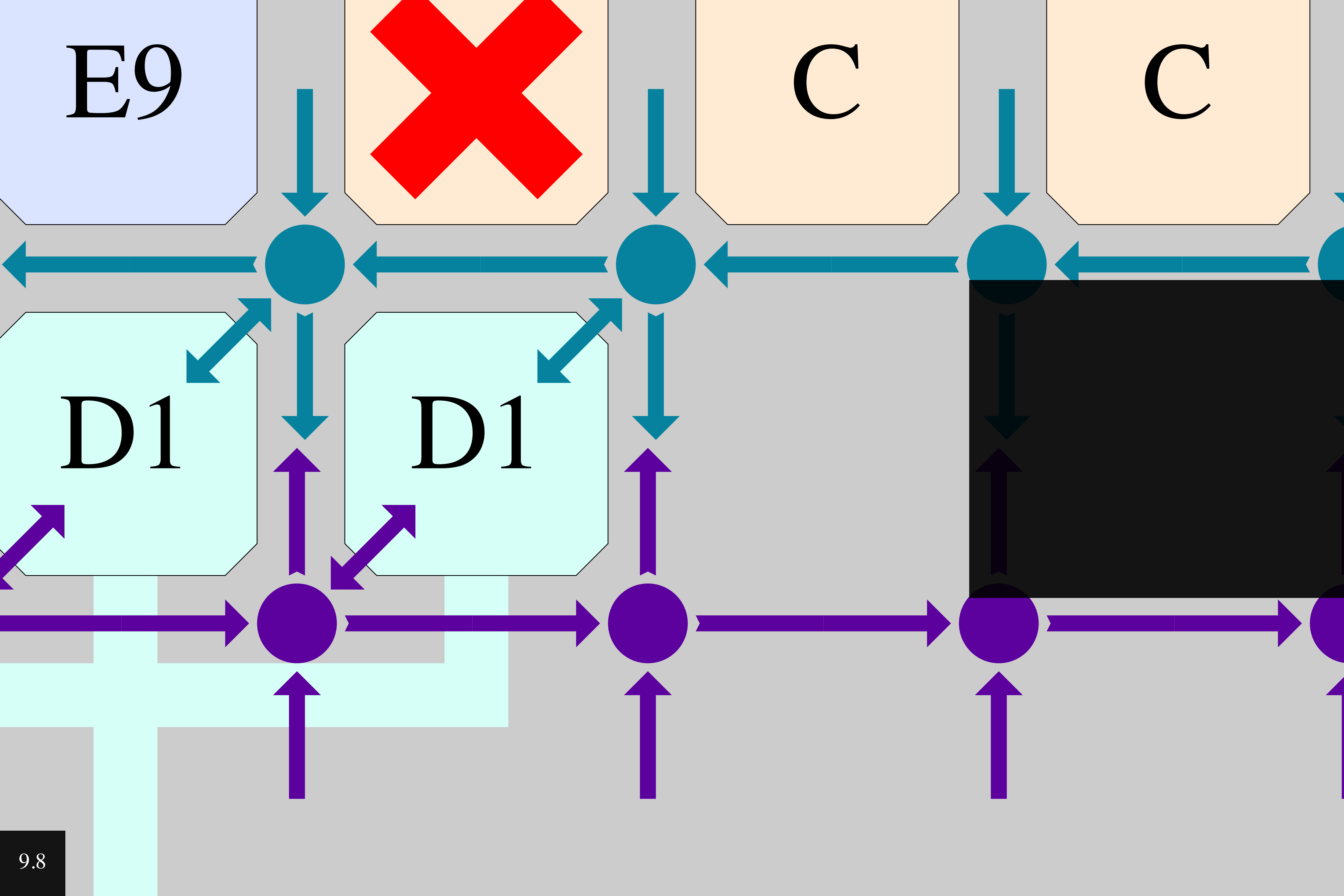


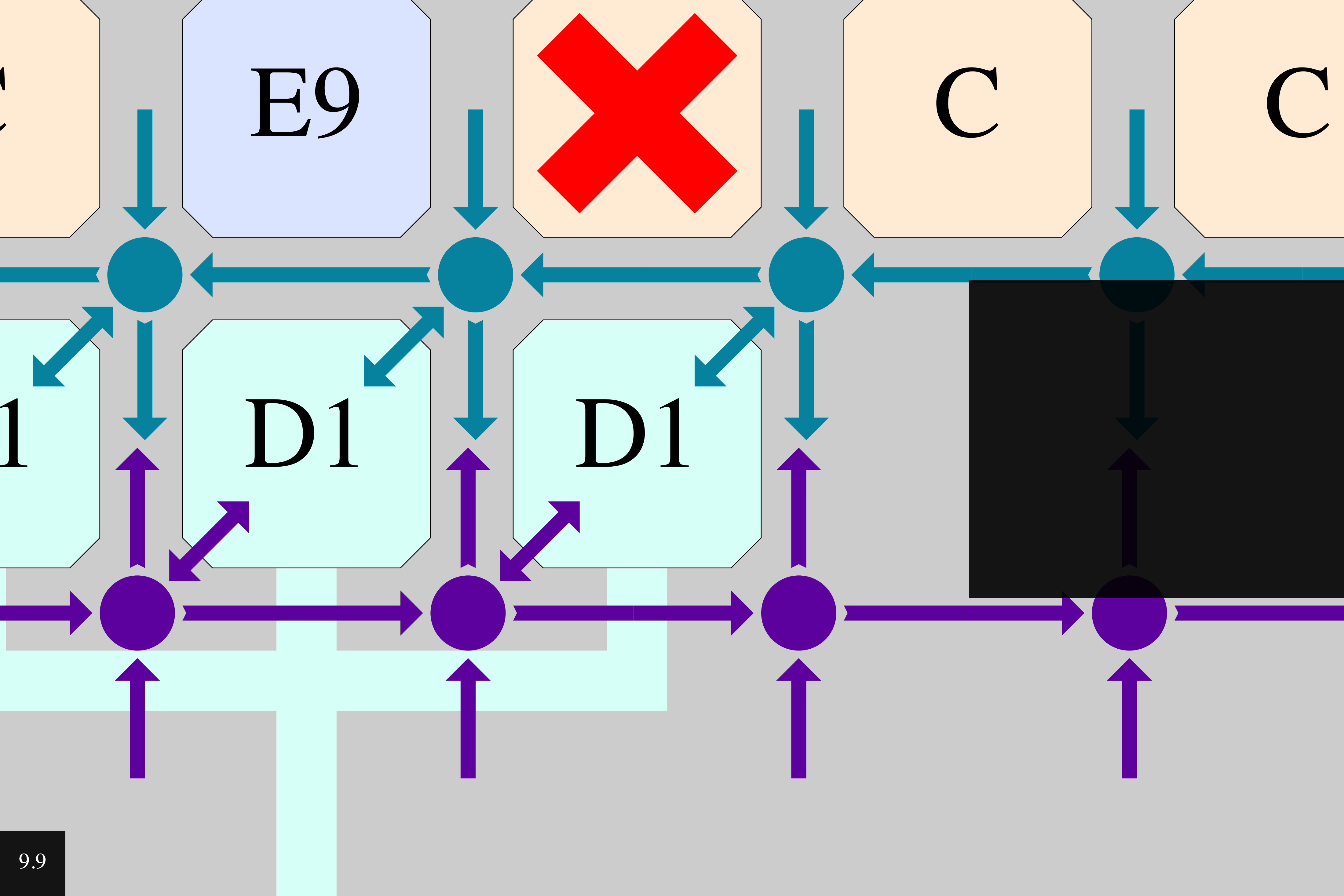
1

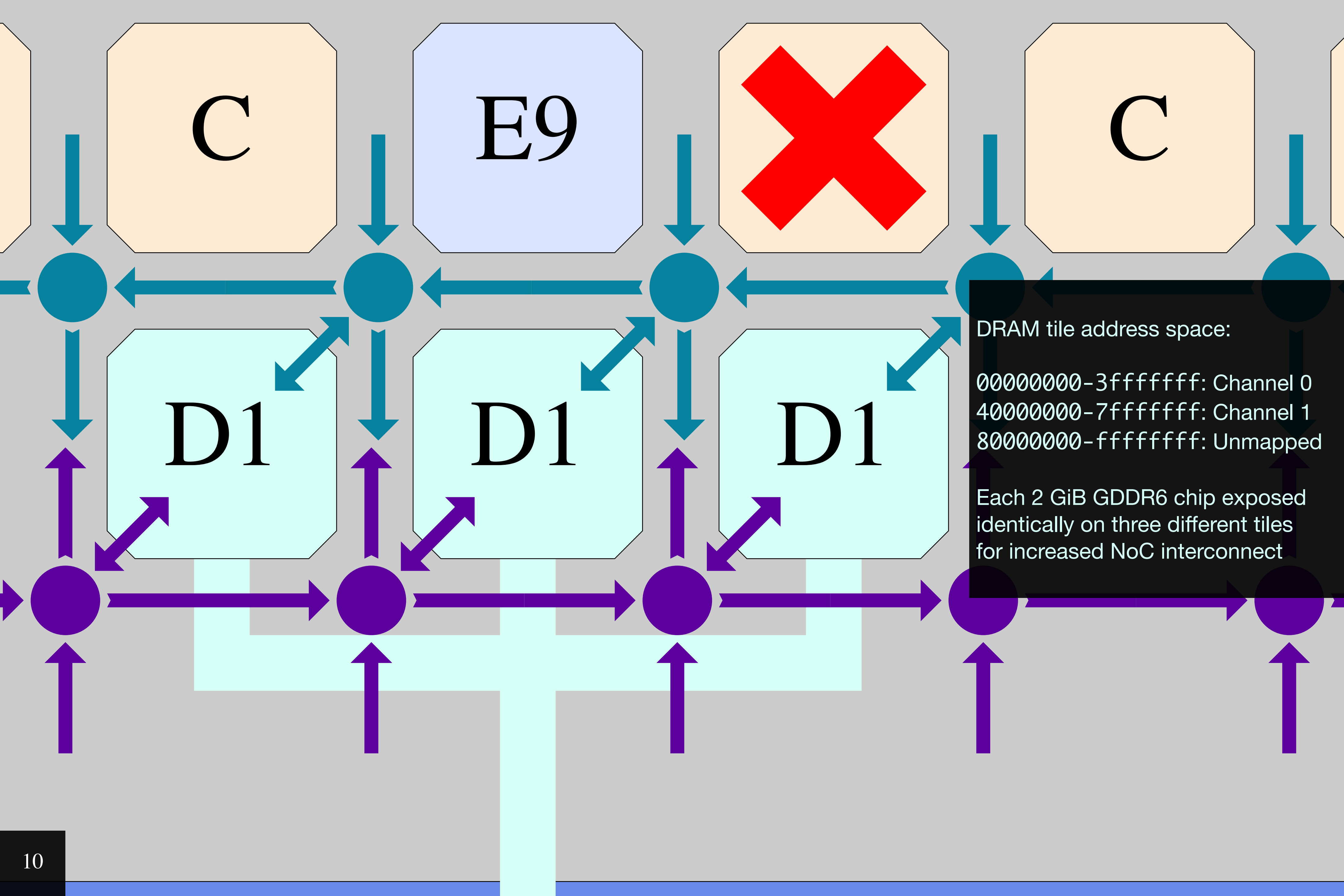
ARC

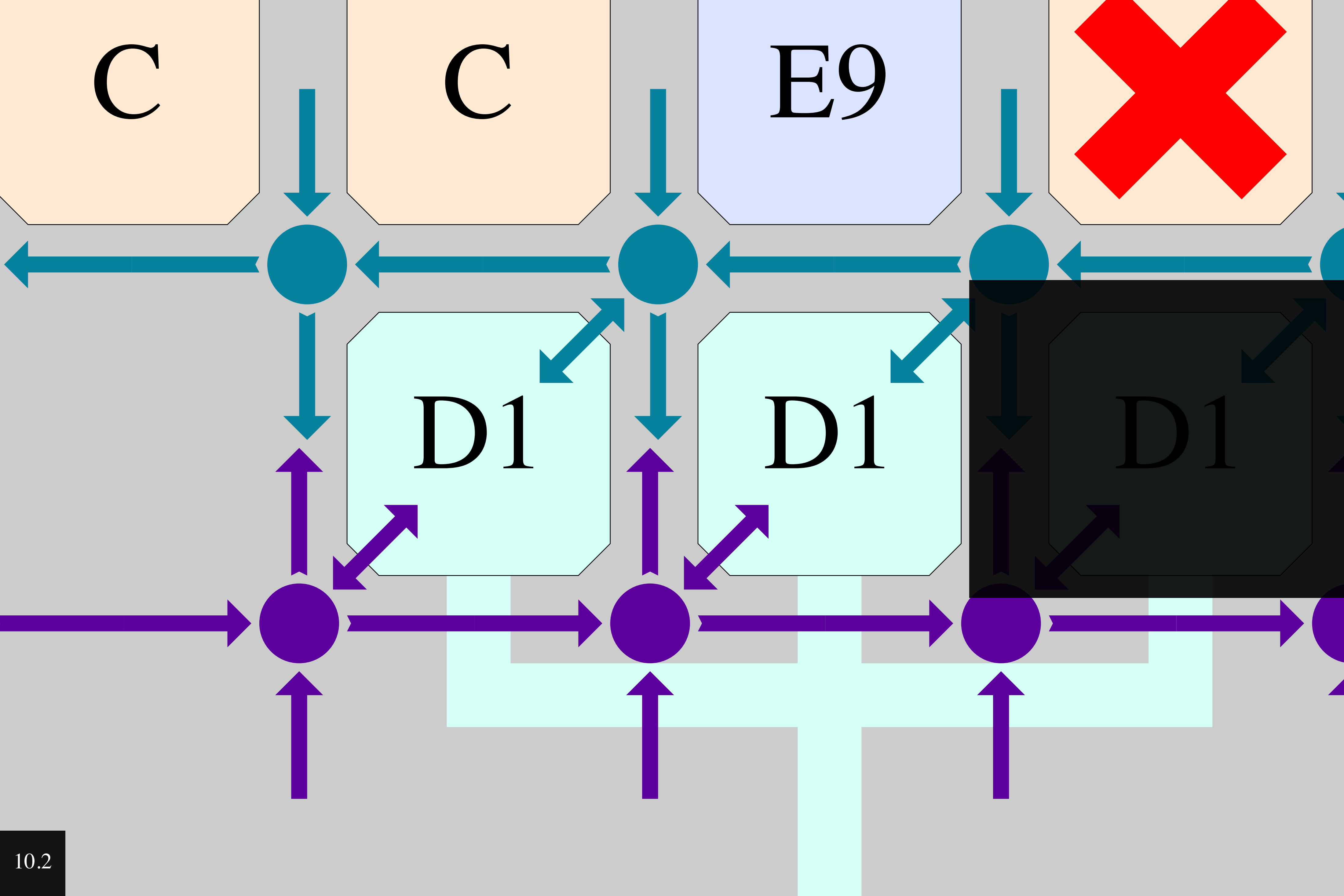


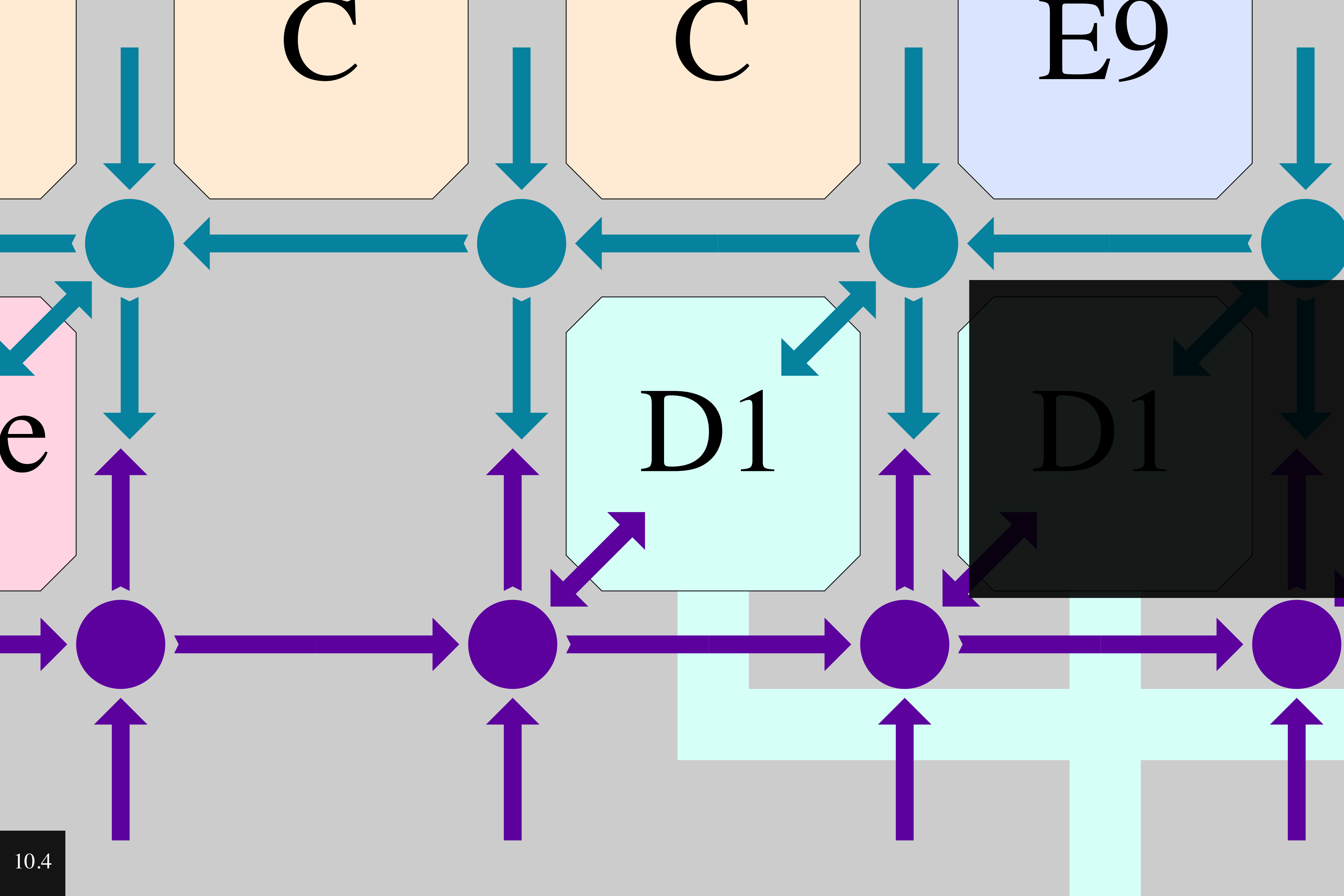


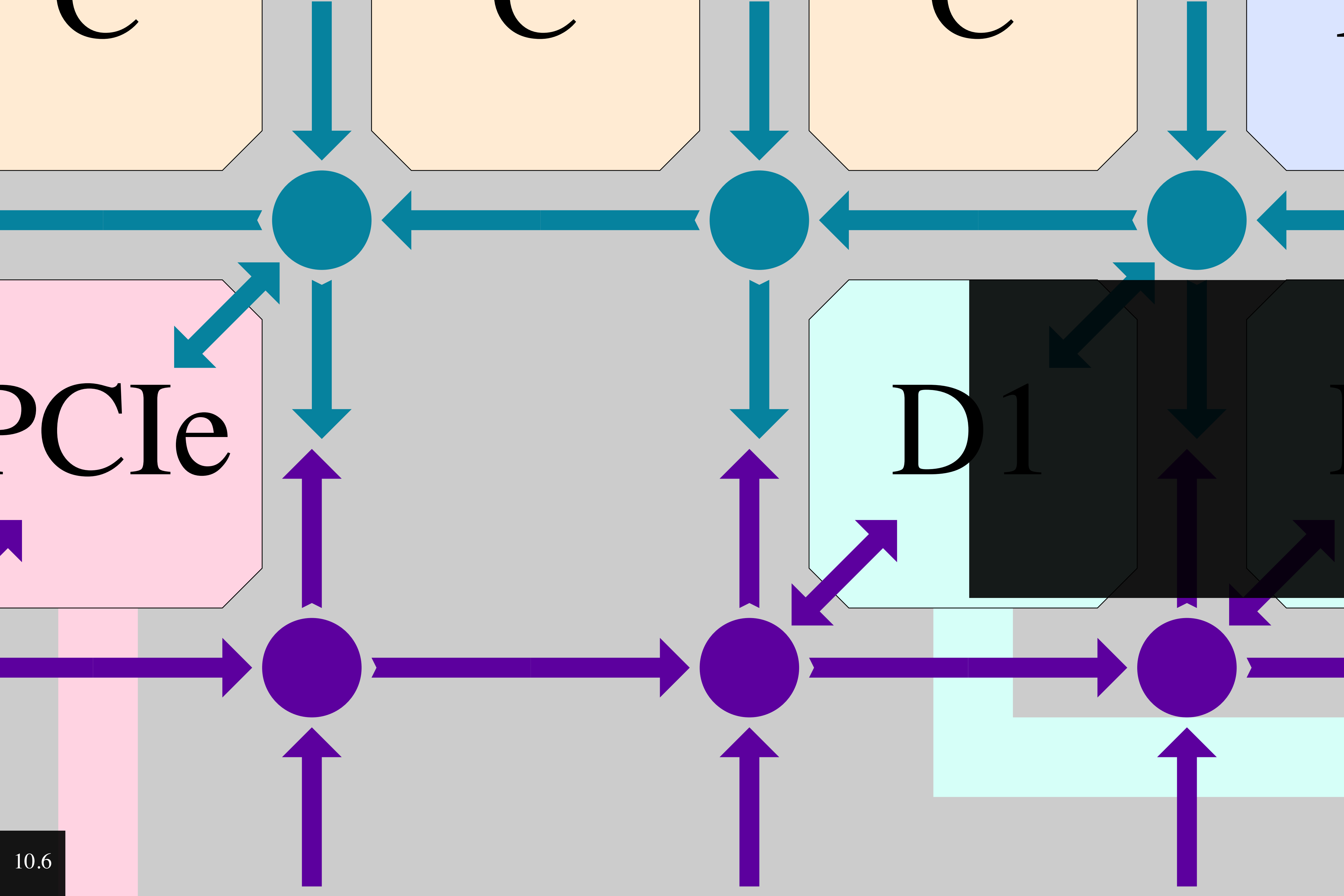






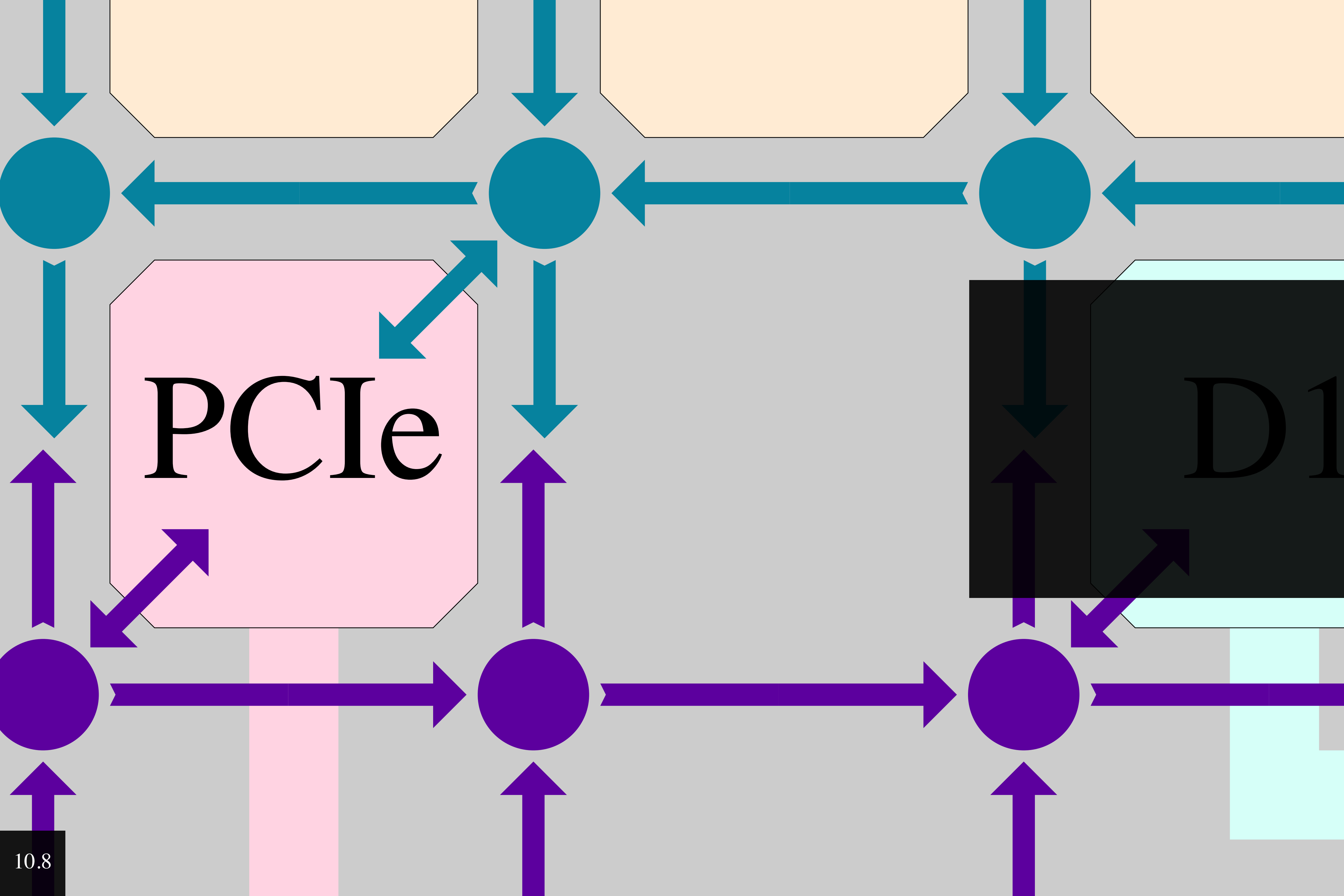






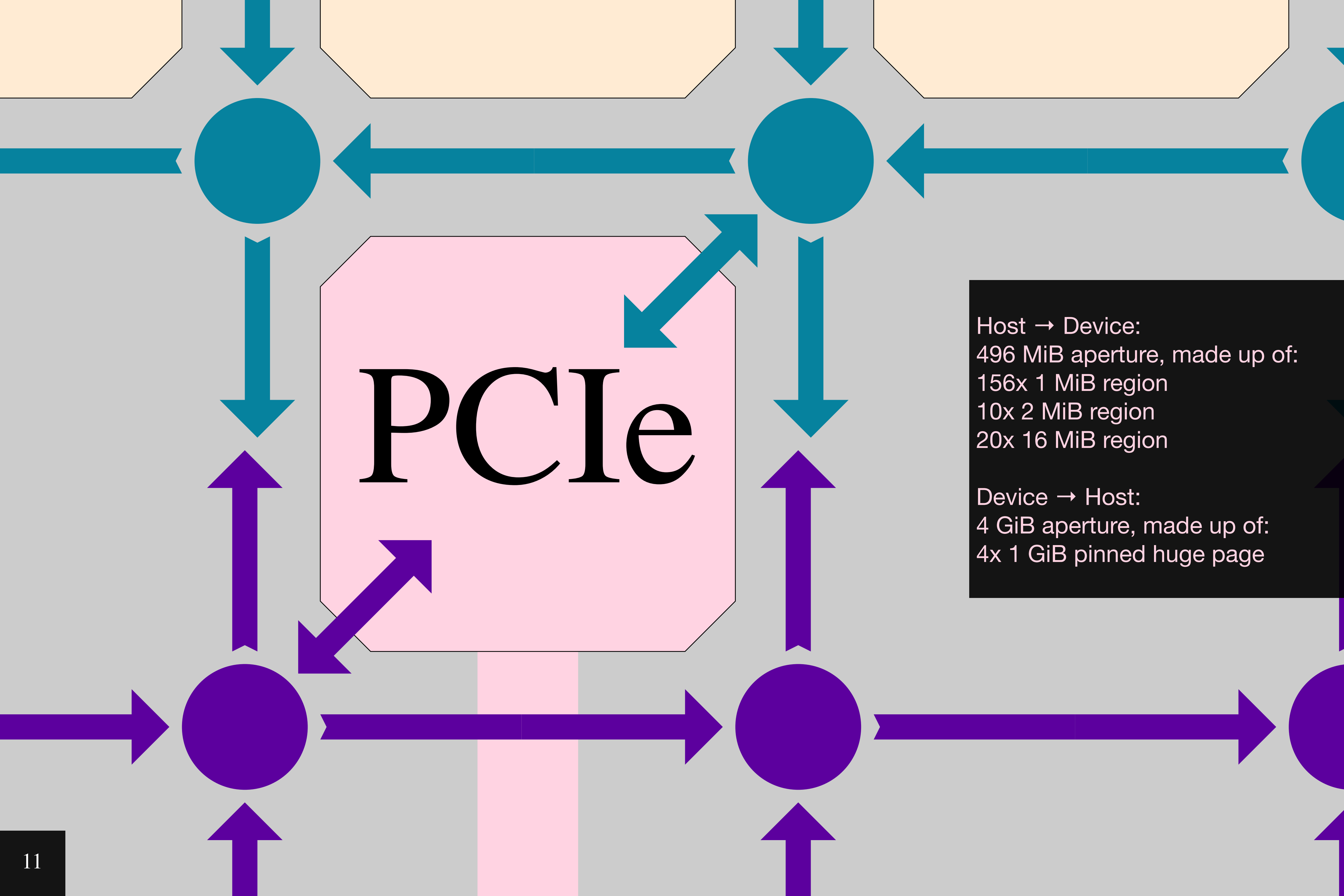
PCIE

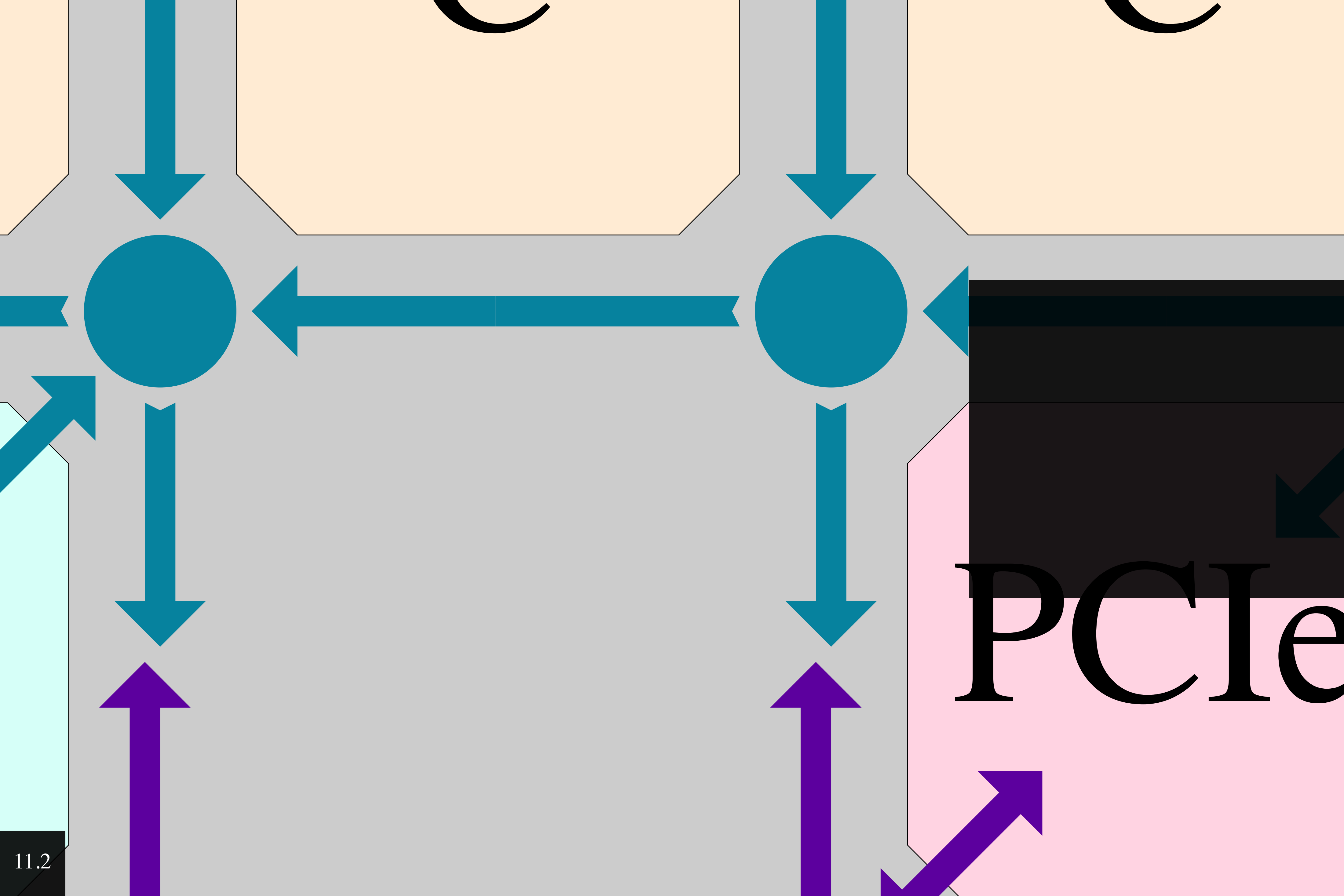
D



PCIe

D1

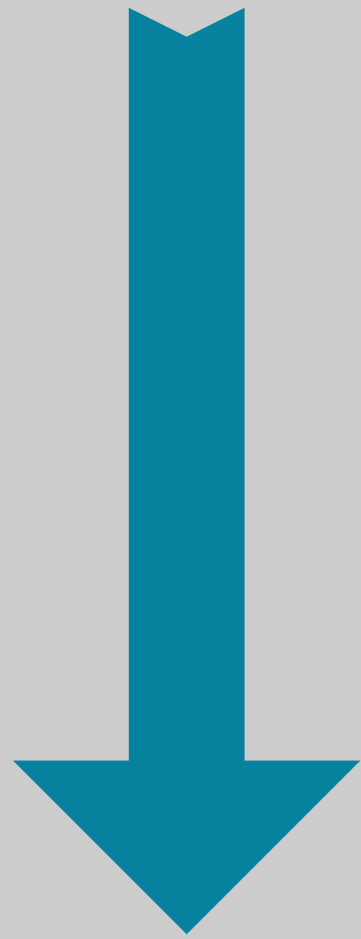
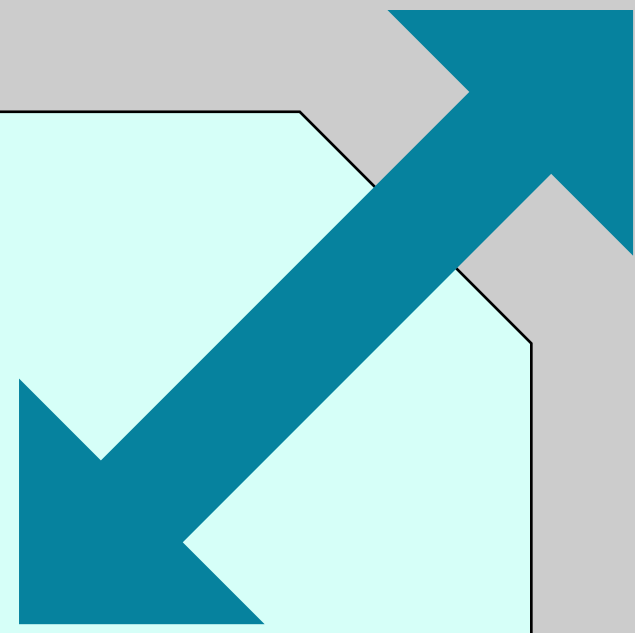
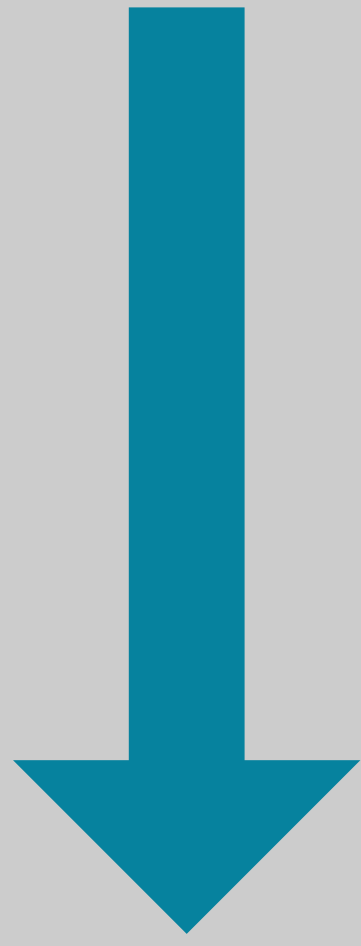




PCIe

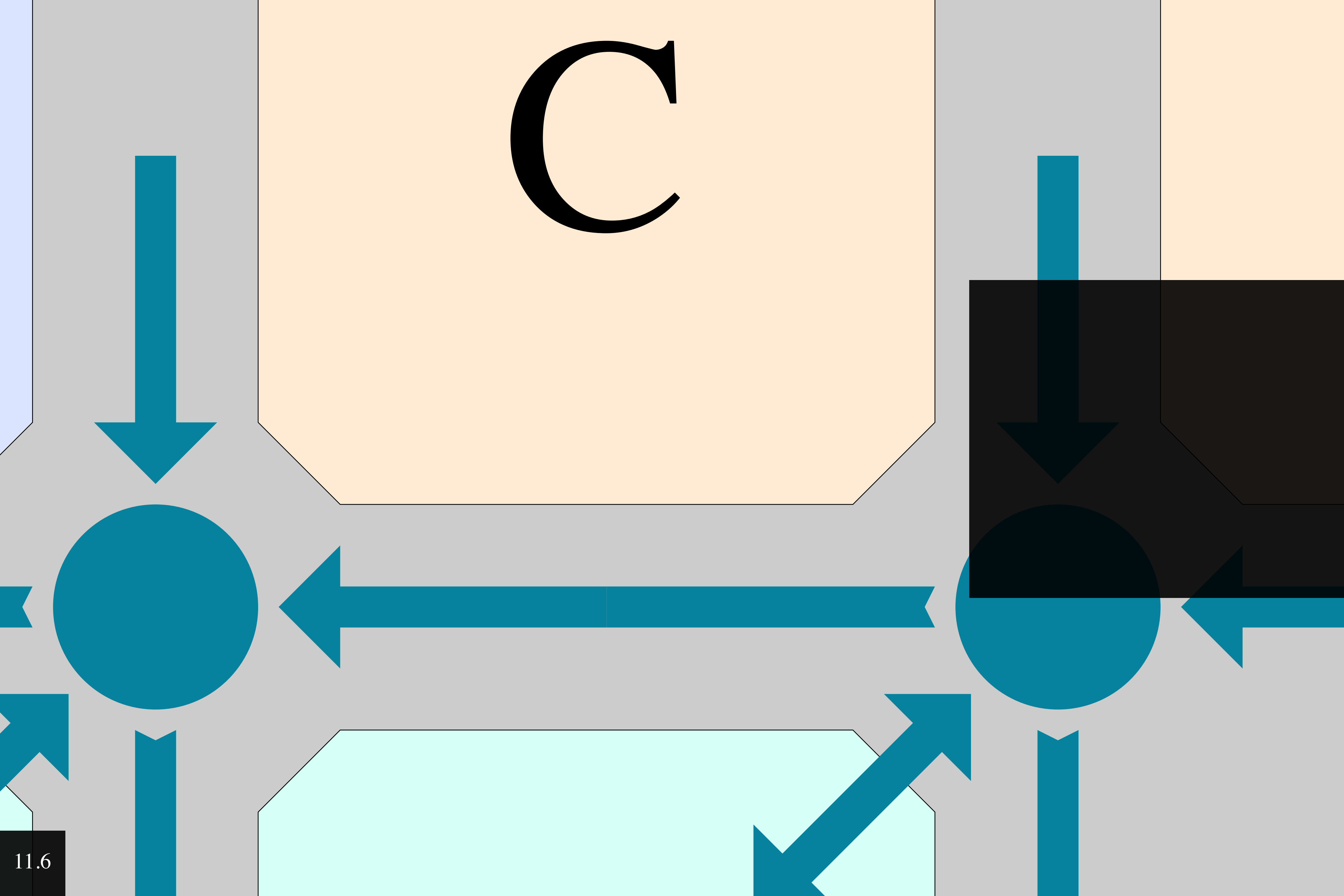
C

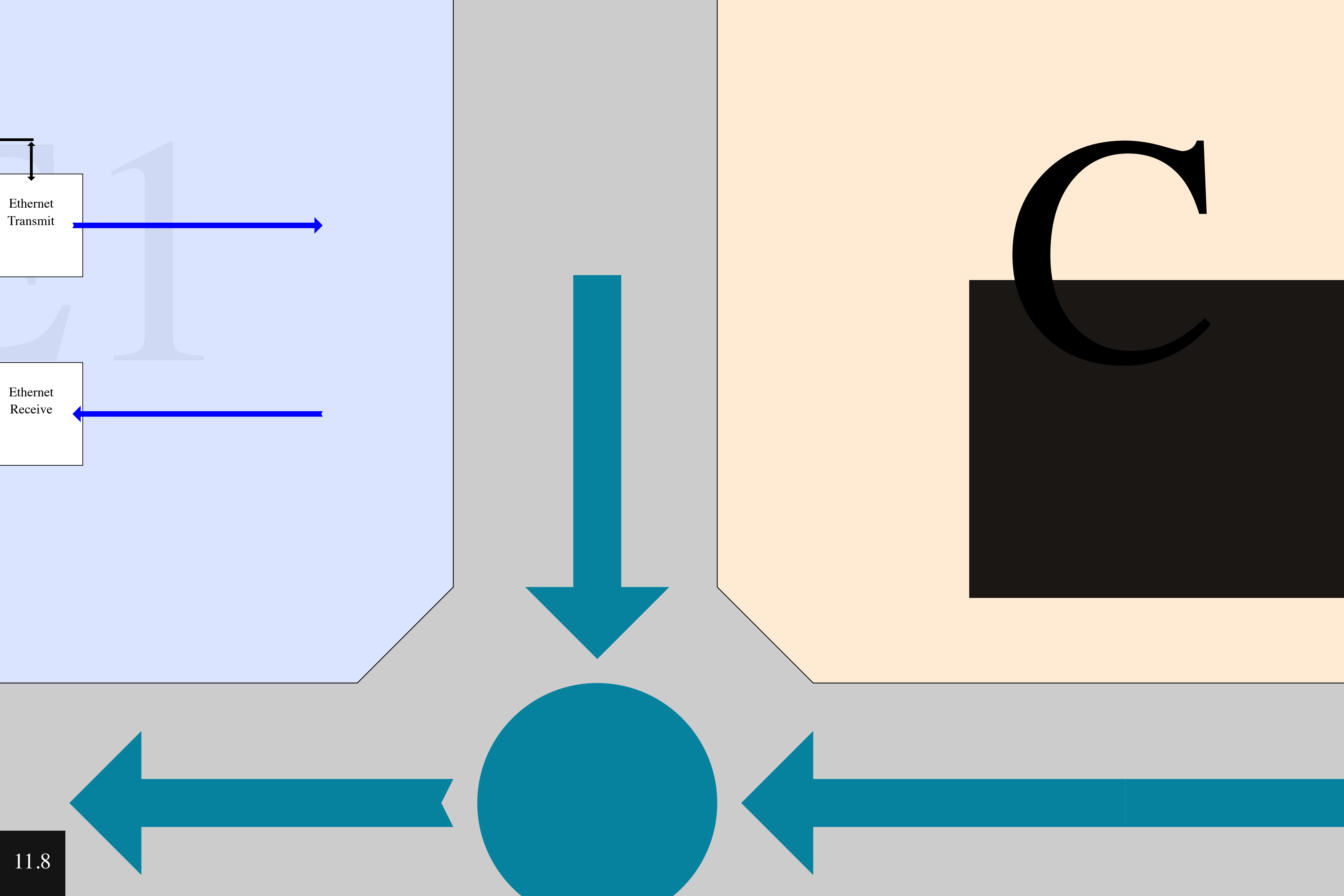
C

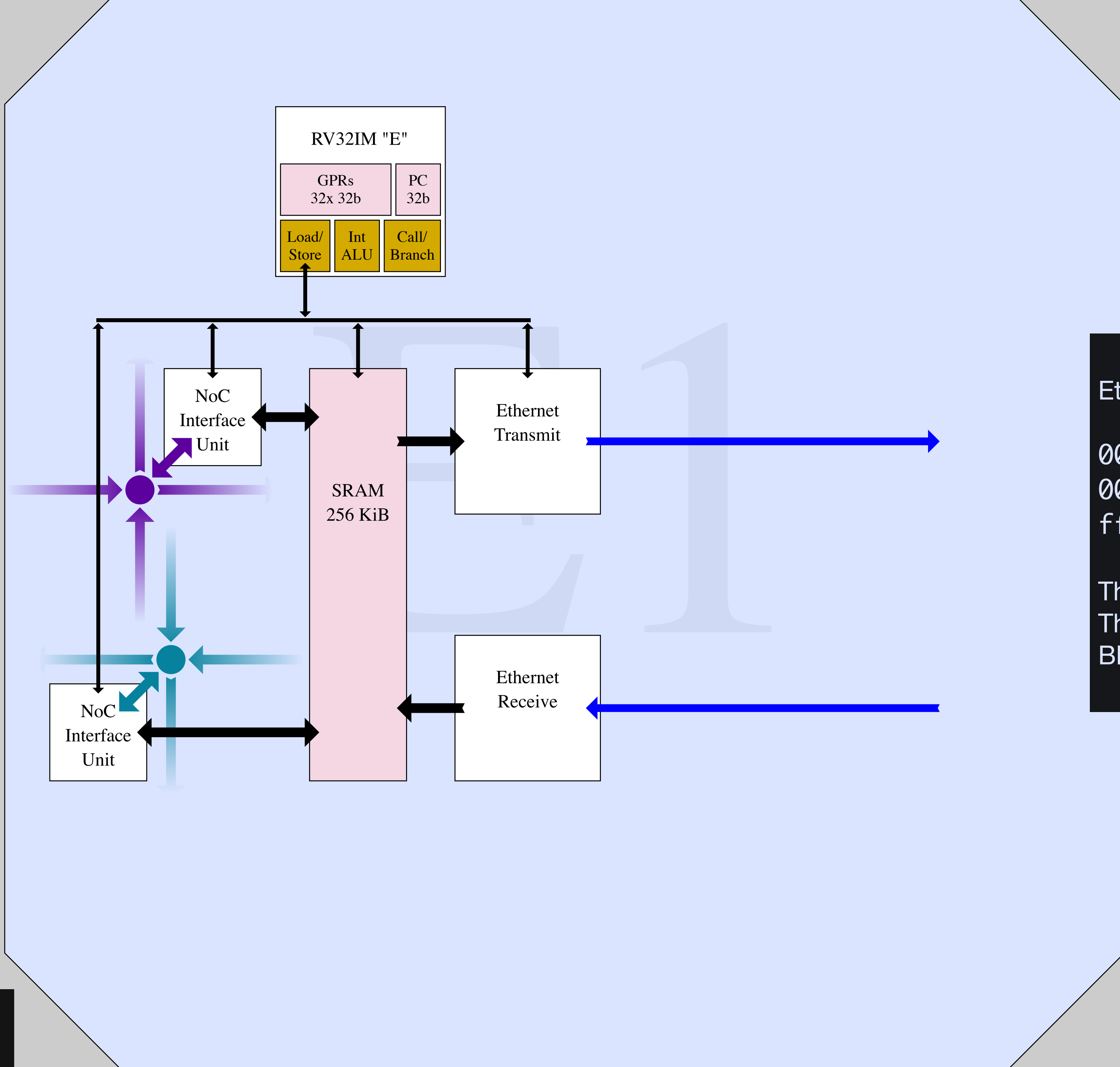


DO

C



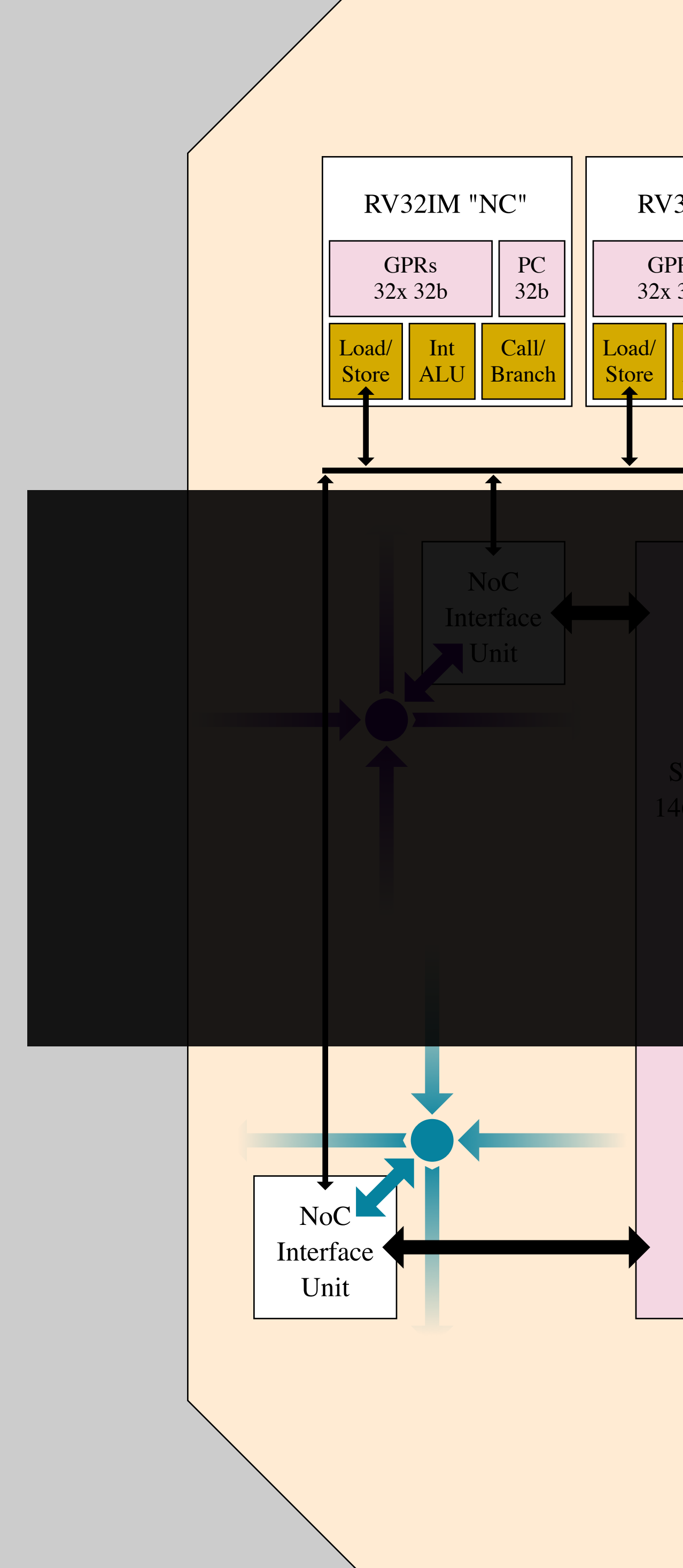
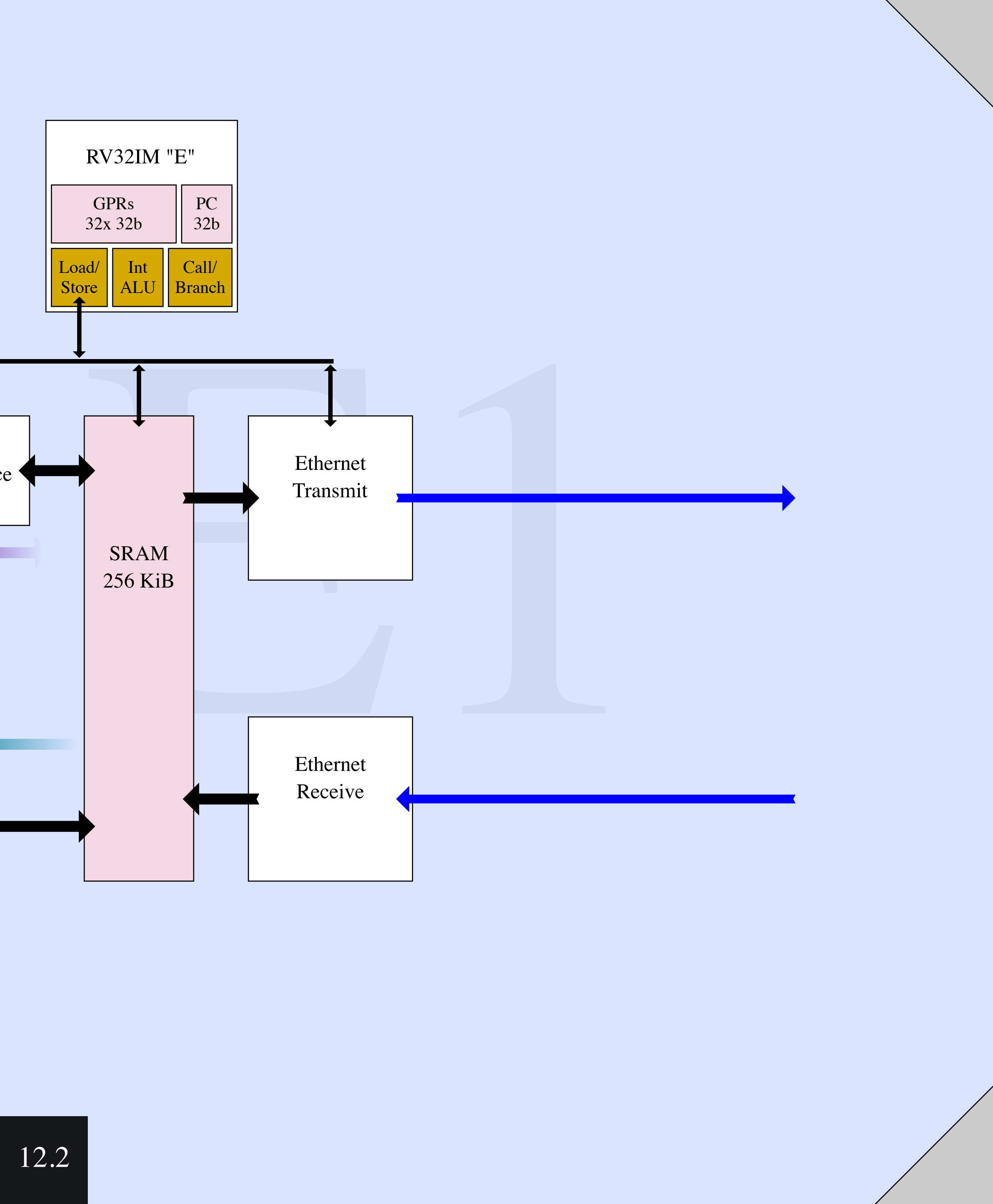


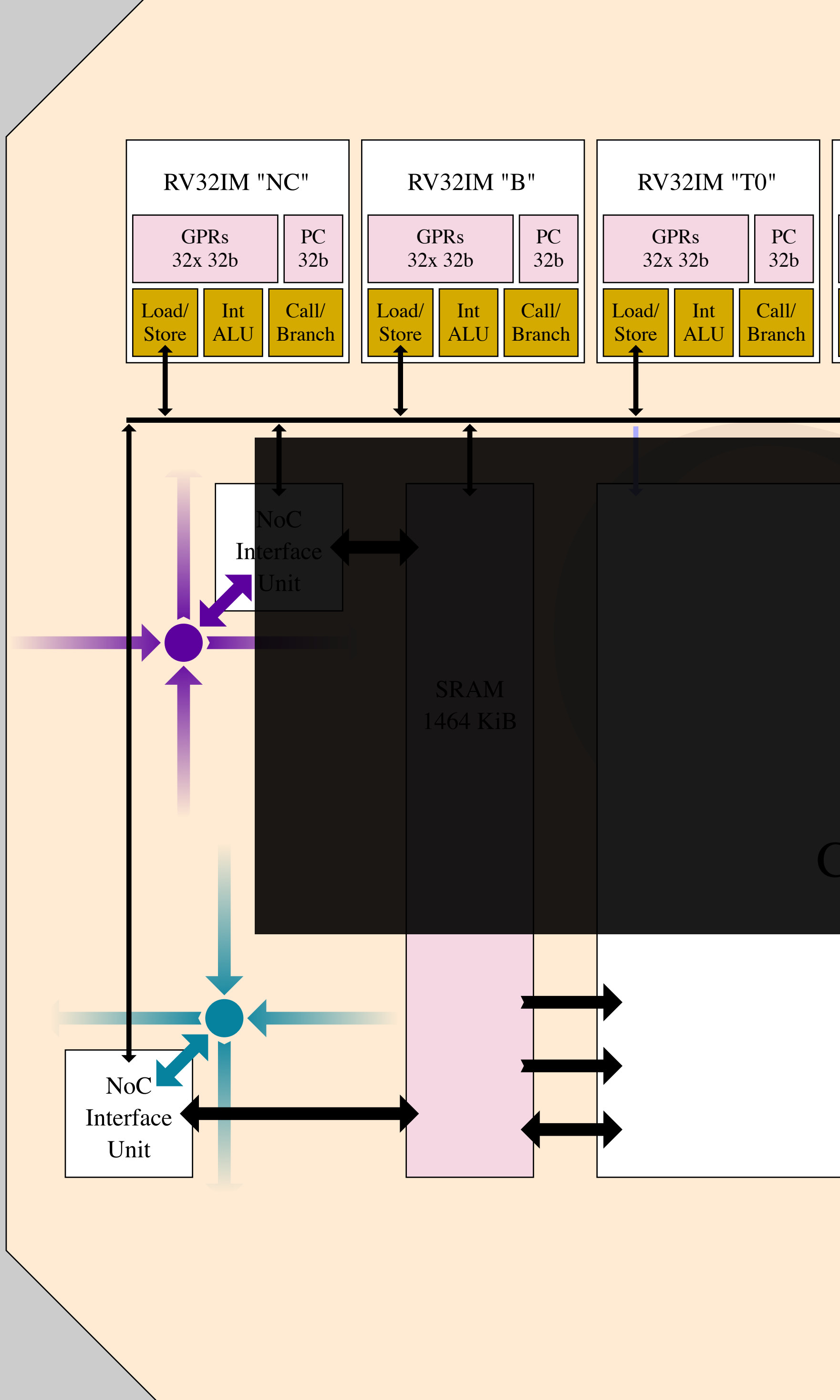
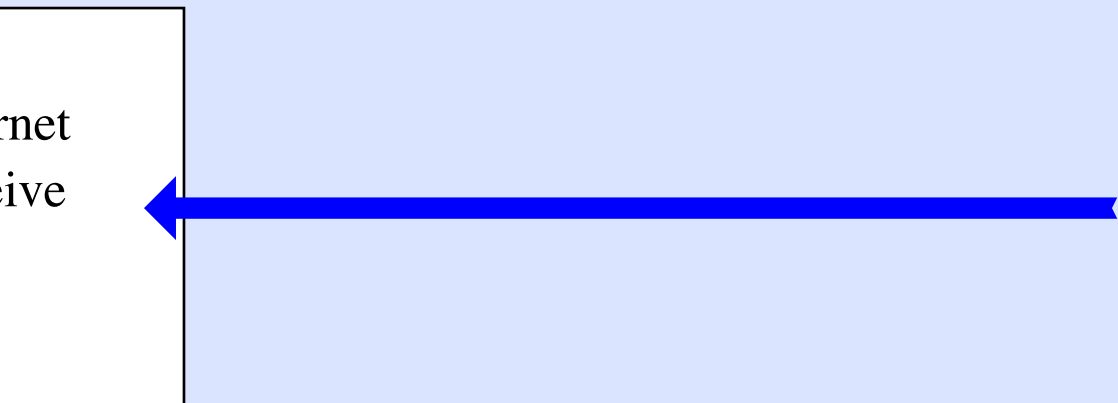
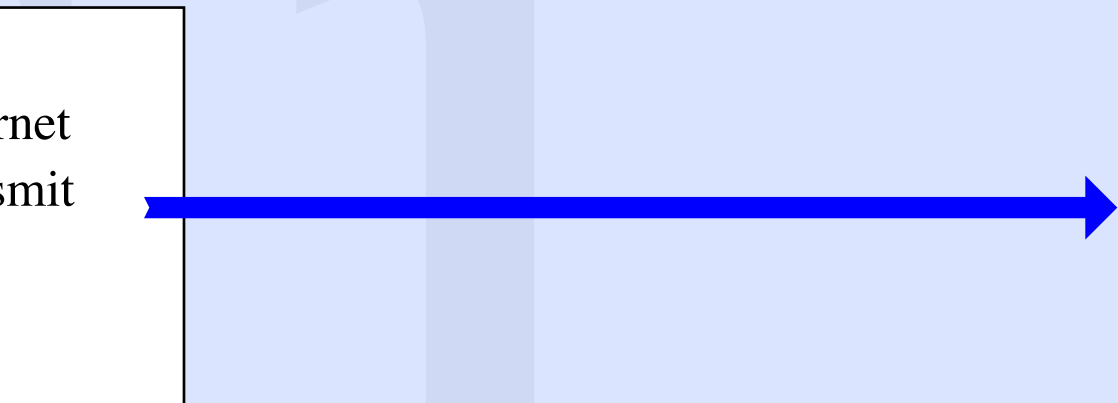


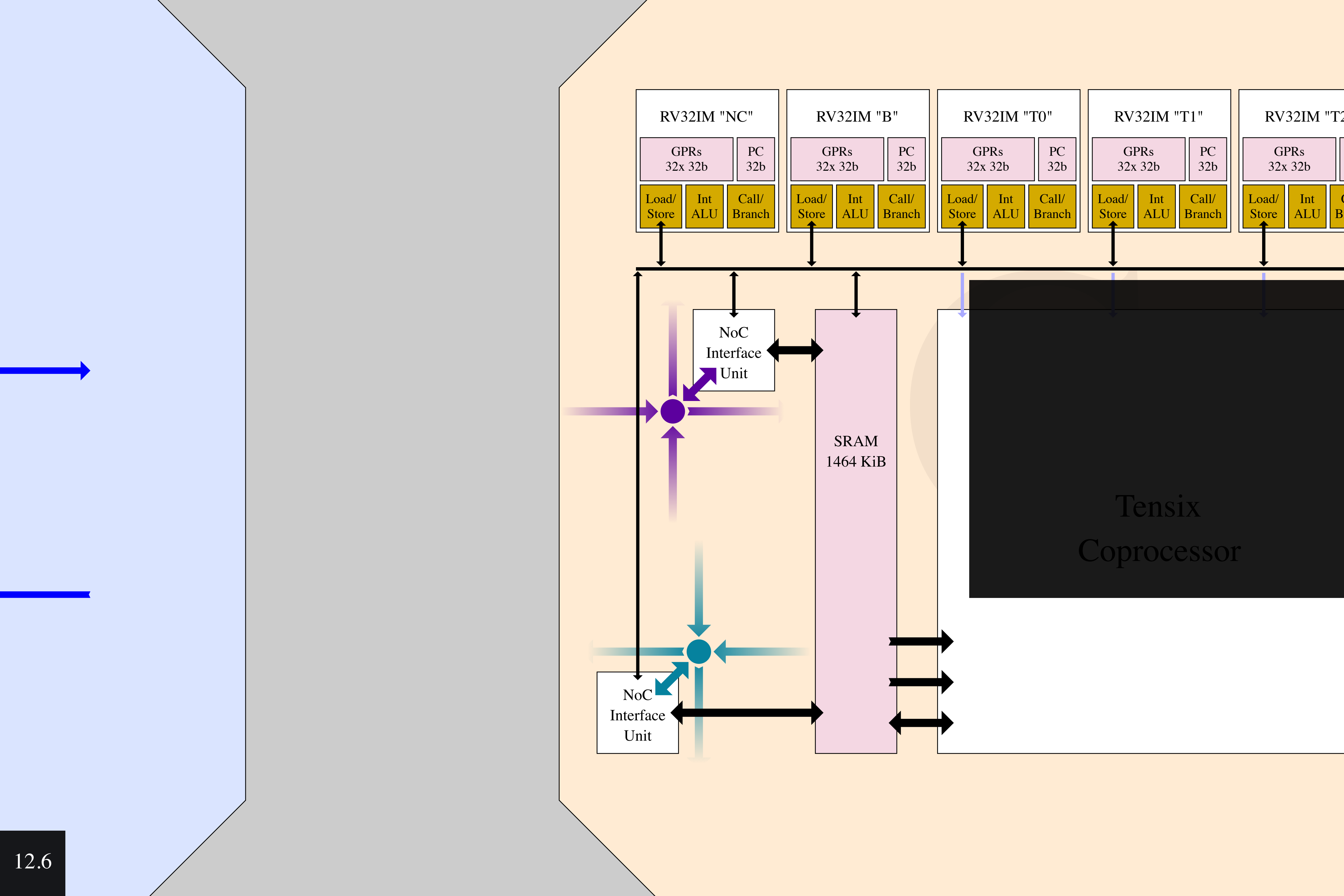
Ethernet tile address space:

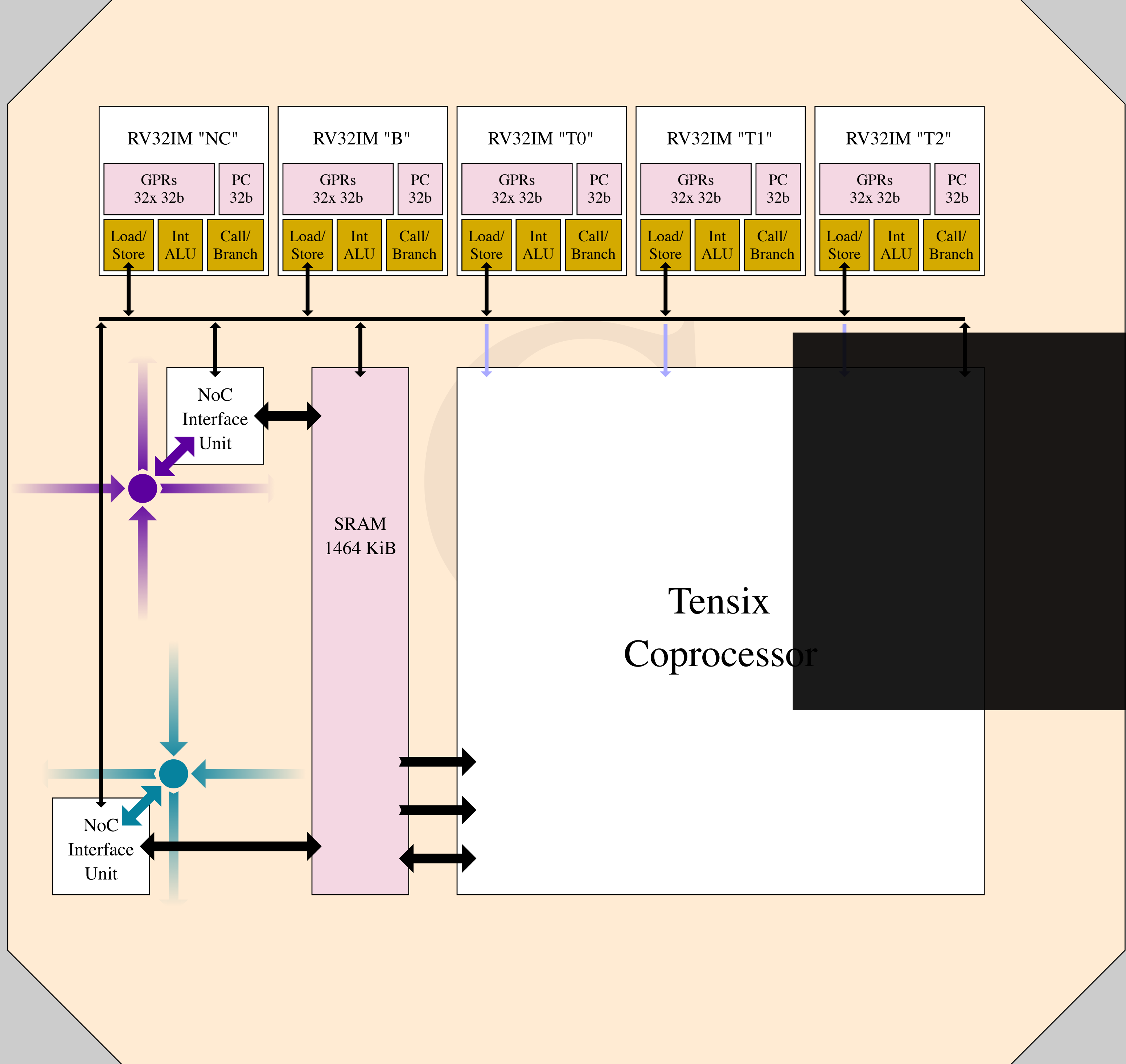
- 00000000-0003ffff: SRAM
- 00040000-feffffff: Unmapped
- ff000000-ffffffff: Peripherals

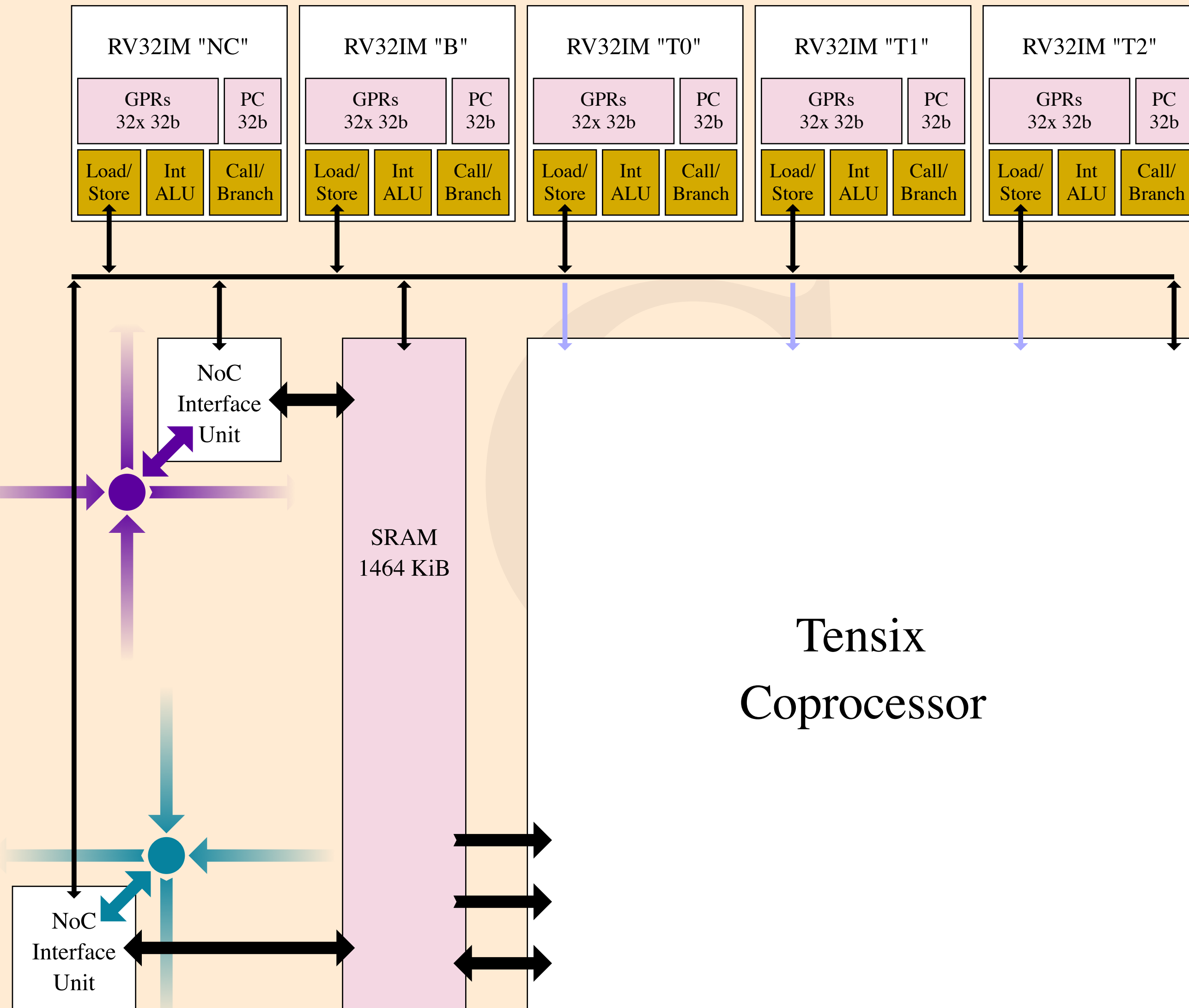
Thin arrows 4 bytes wide
 Thick arrows 32 bytes wide
 Blue arrows 100GbE







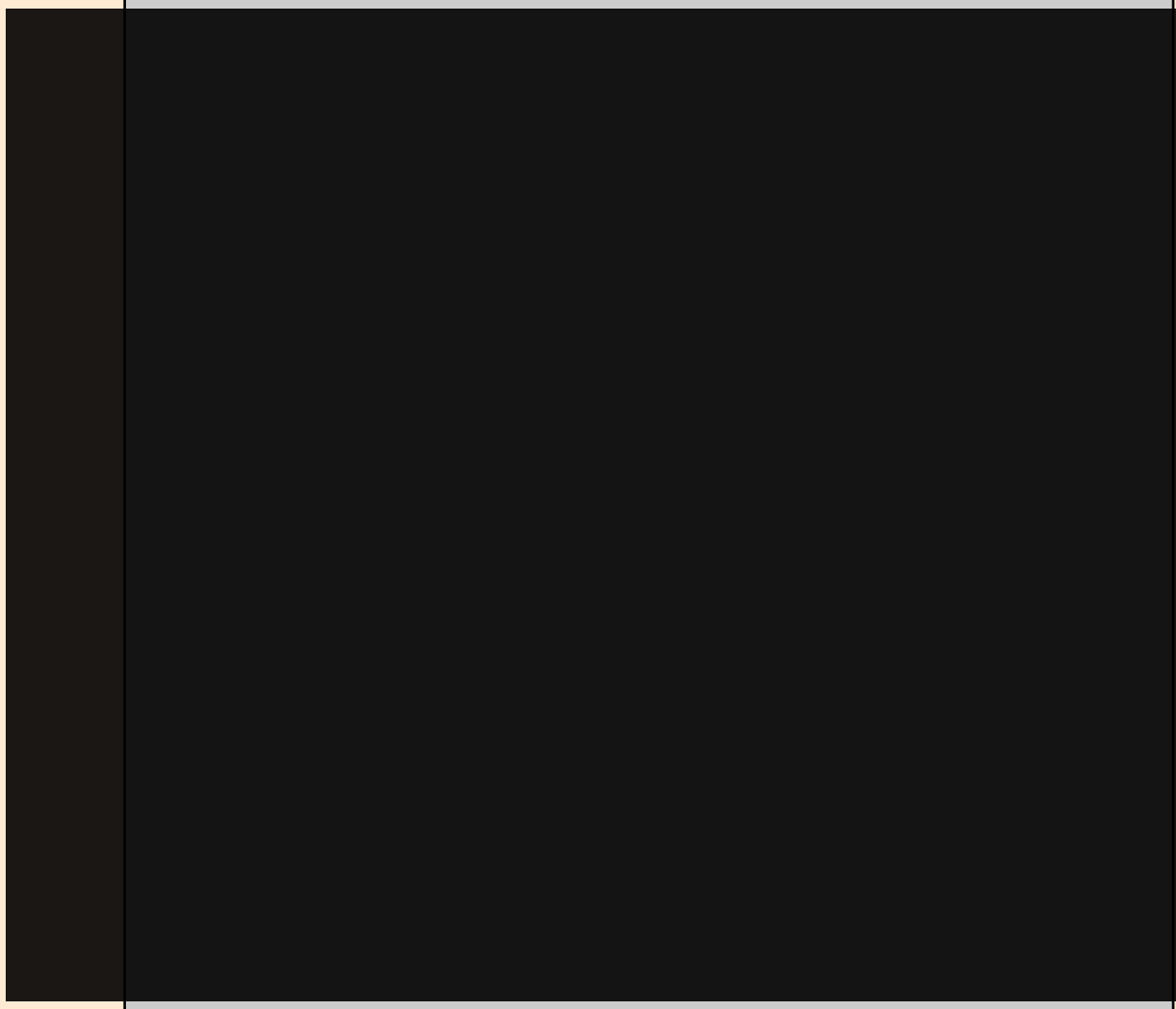
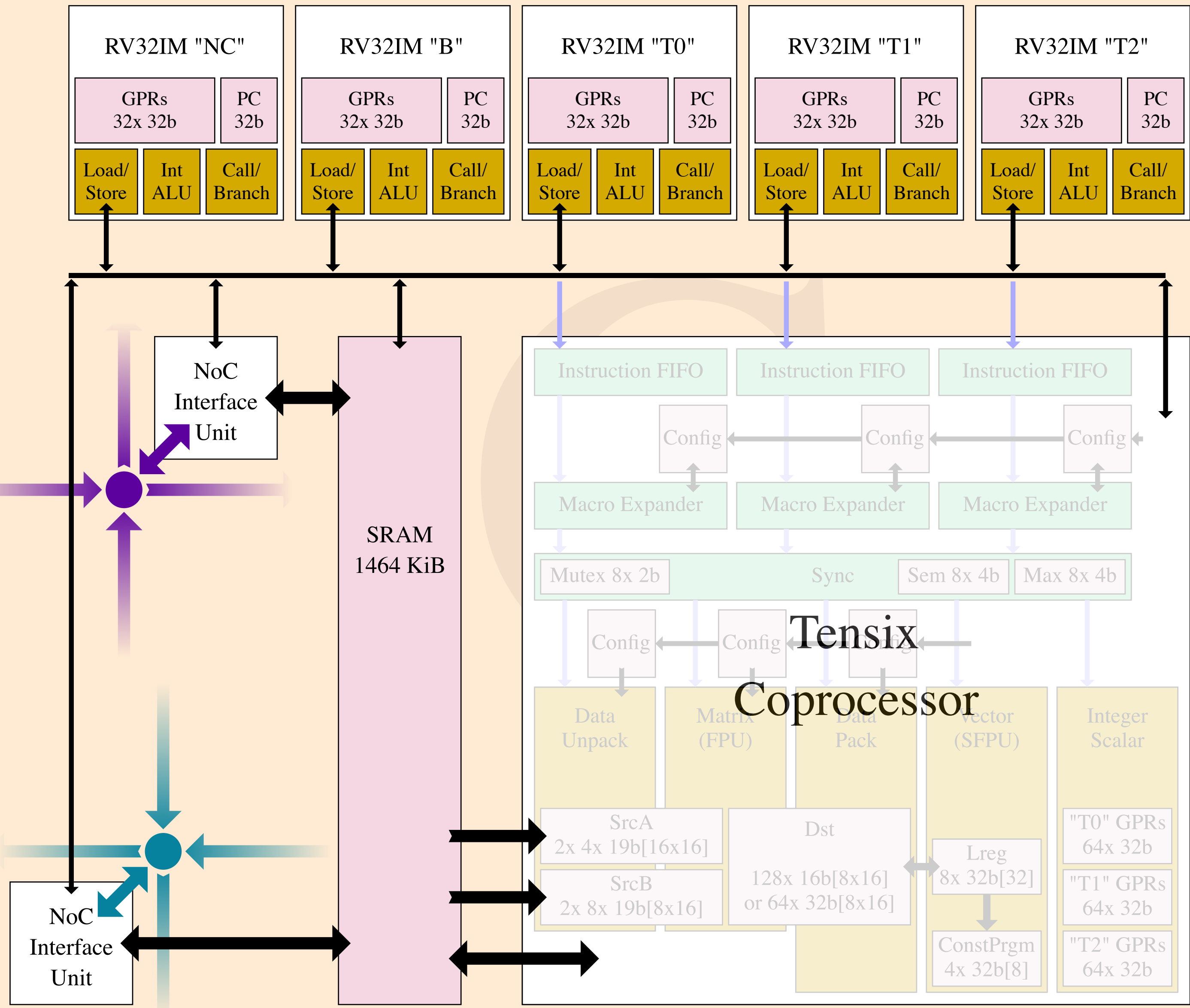


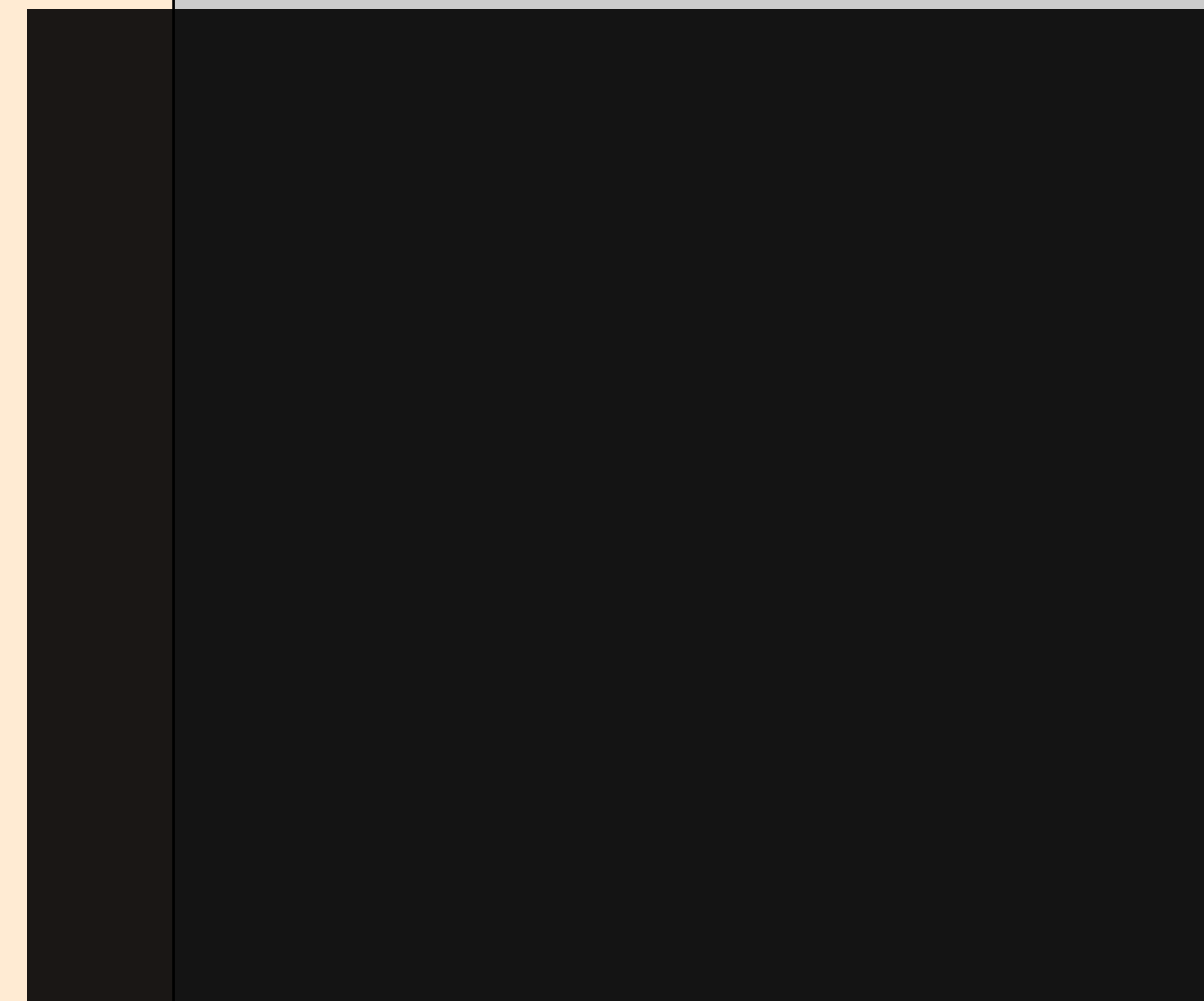
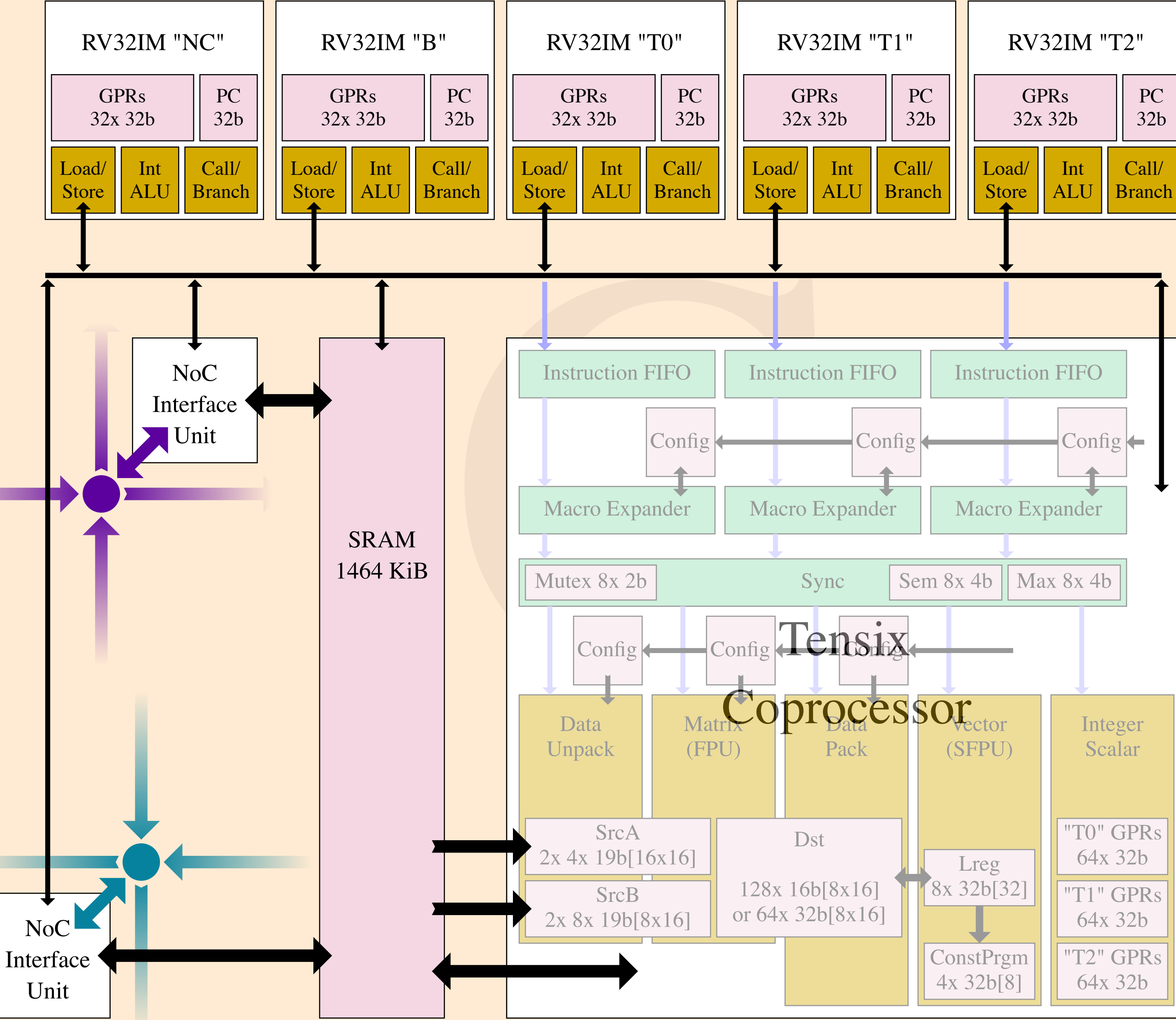


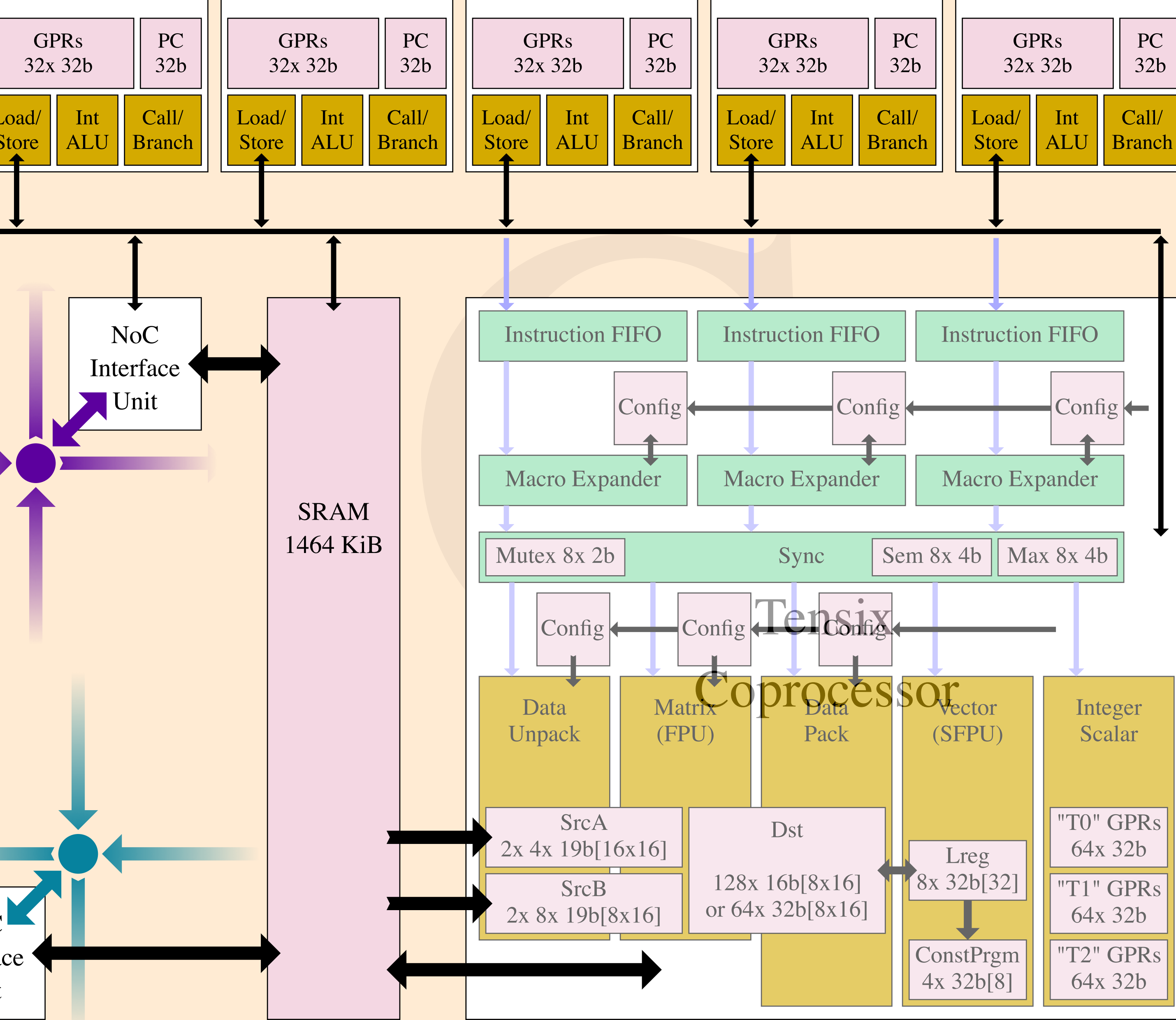
Compute tile address space:

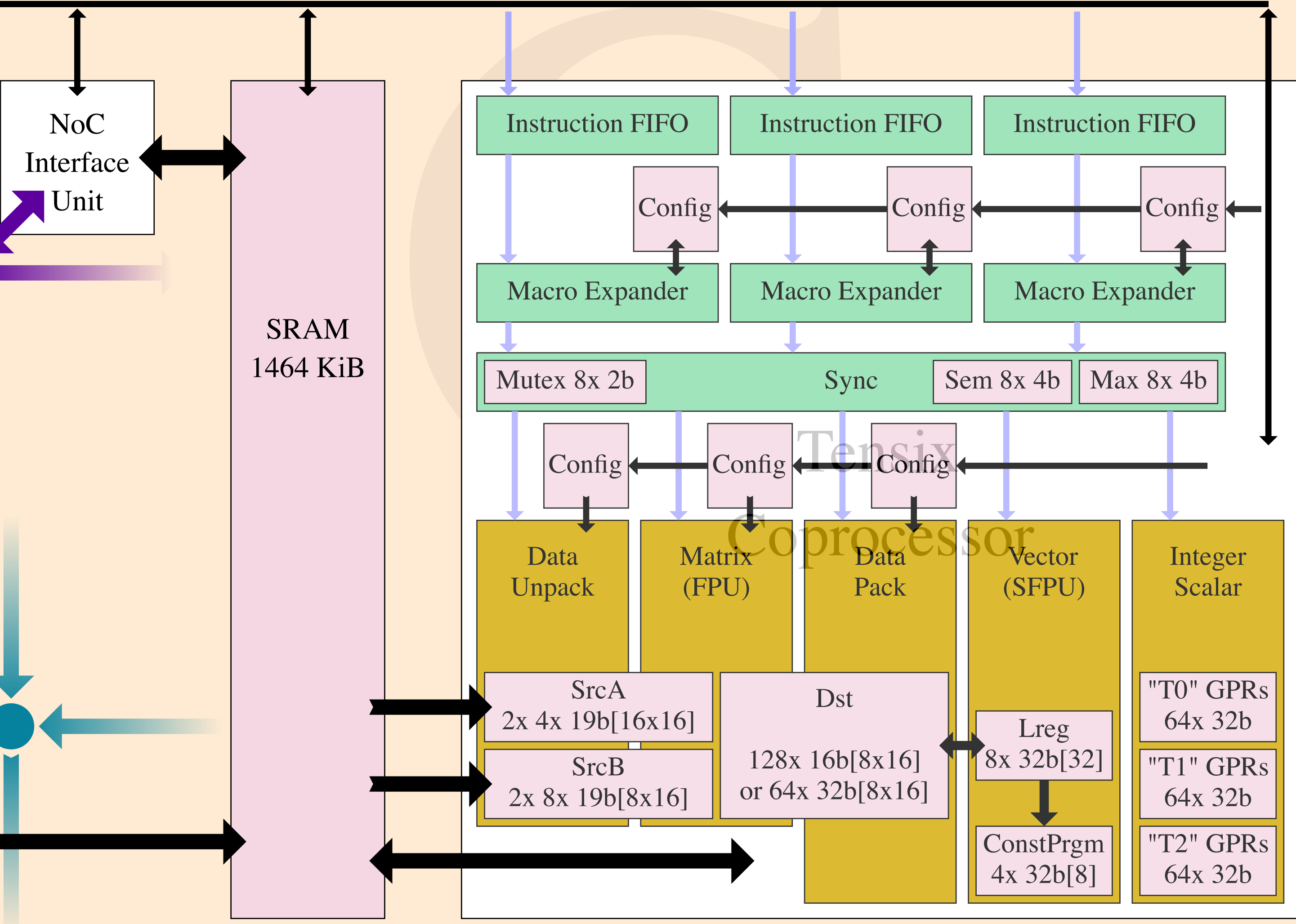
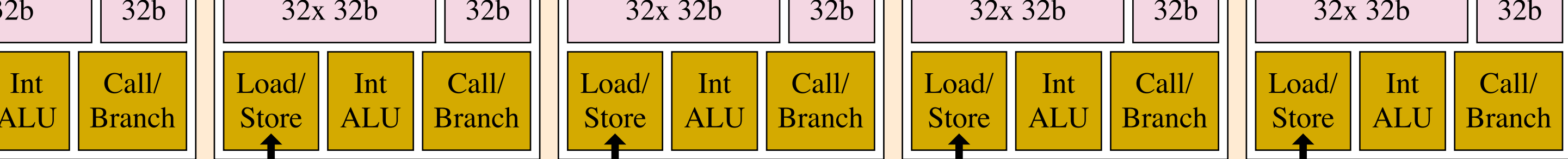
- 00000000-0016dfff: SRAM
- 0016e000-feffffff: Unmapped
- ff000000-ffffffff: Peripherals

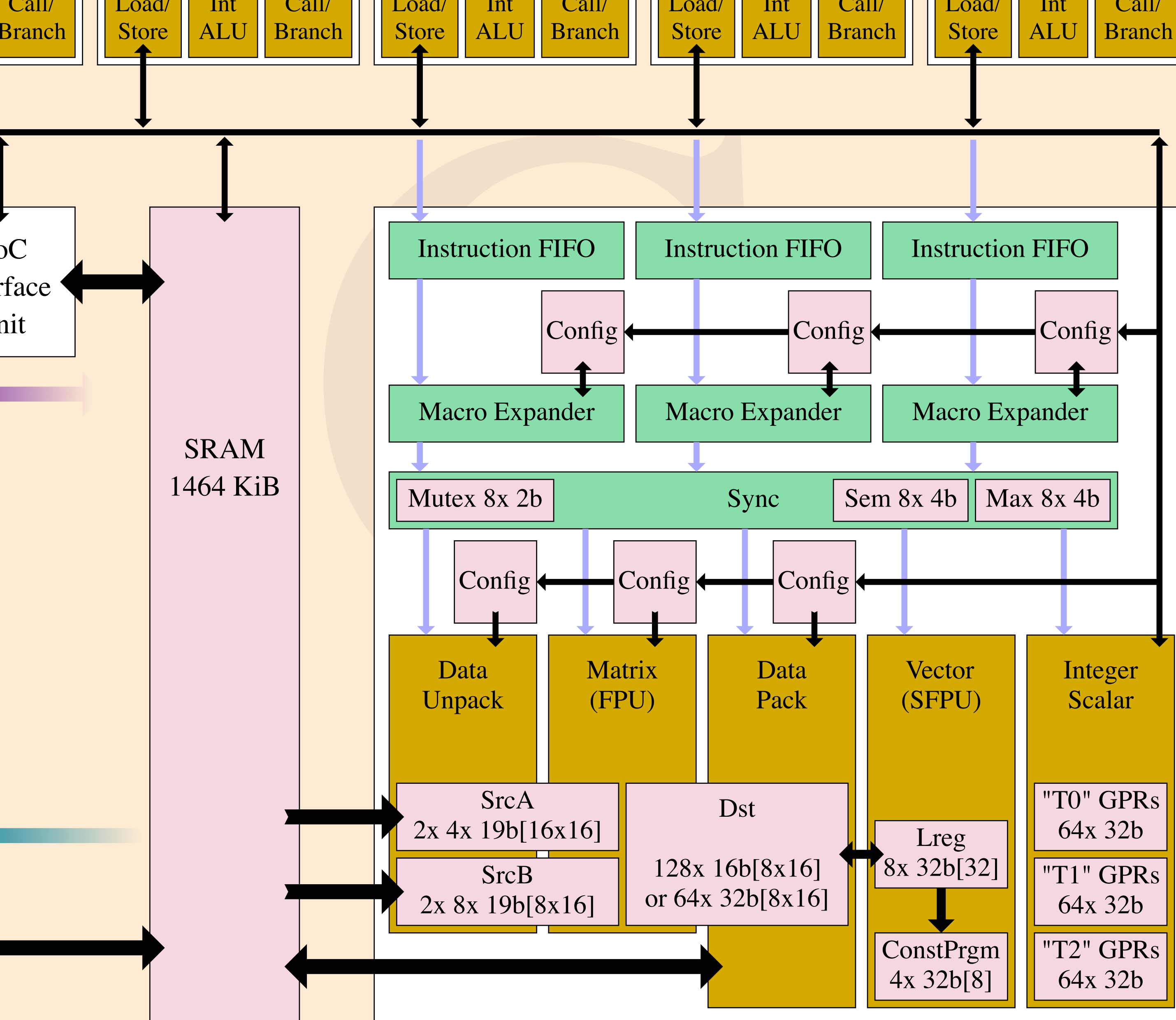
Thin arrows 4 bytes wide
Thick arrows 32 bytes wide











Data Unpack: SRAM → SrcA
 Data Unpack: SRAM → SrcB
 Matrix: Dst += SrcB @ SrcA
 Matrix 8x16 += 8x16 @ 16x16
 Matrix fp8, 1/2 rate bf16, 1/4 rate fp16
 Data Pack: Dst → SRAM
 Data Pack: or SRAM += Dst
 Vector: 32 lanes of SIMD
 Vector: fp32 or int32

E2

C

C



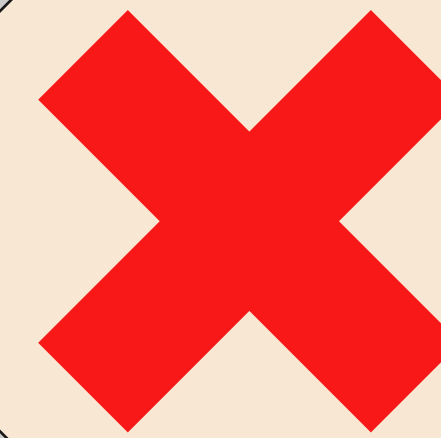
C

C

E1

C

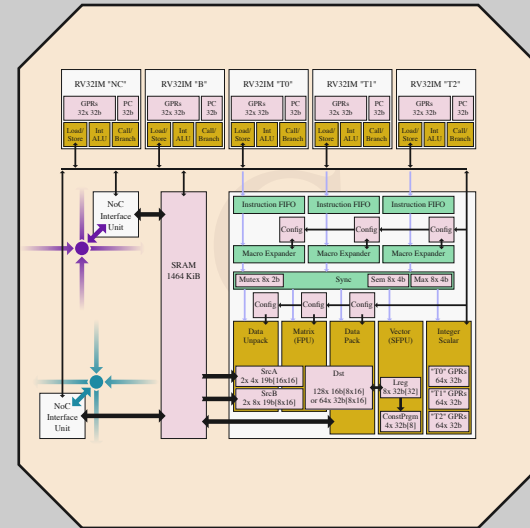
C



C

C

E0



C



C

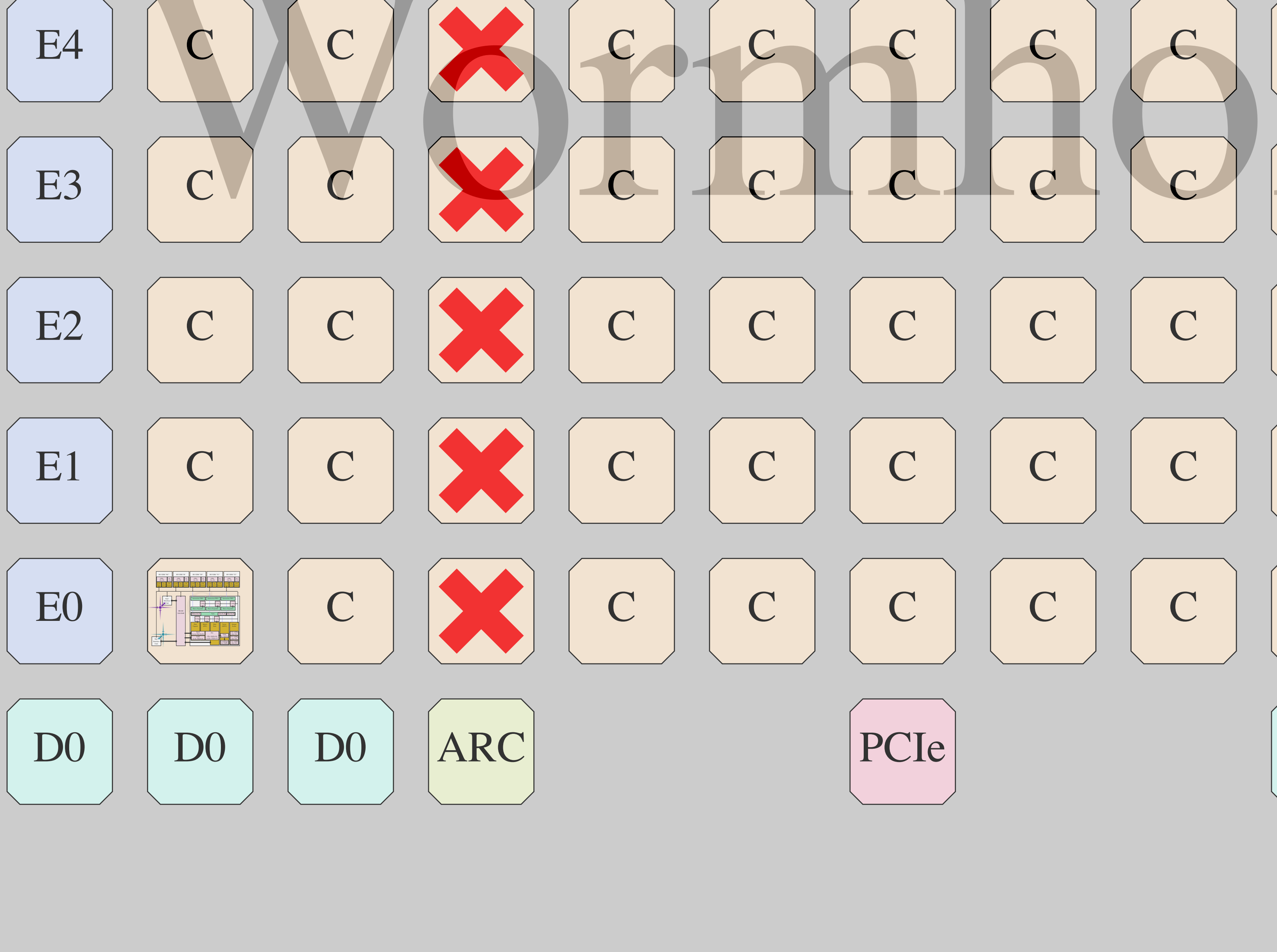
C

D0

D0

D0

ARC





2 GiB

2 GiB



2 GiB
GDDR6

2 GiB
GDDR6

Cage

GDDR6

GDDR6

GDDR6

GDDR6

GDDR6



Cage

2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6

P-DD Cage

2 GiB
GDDR6

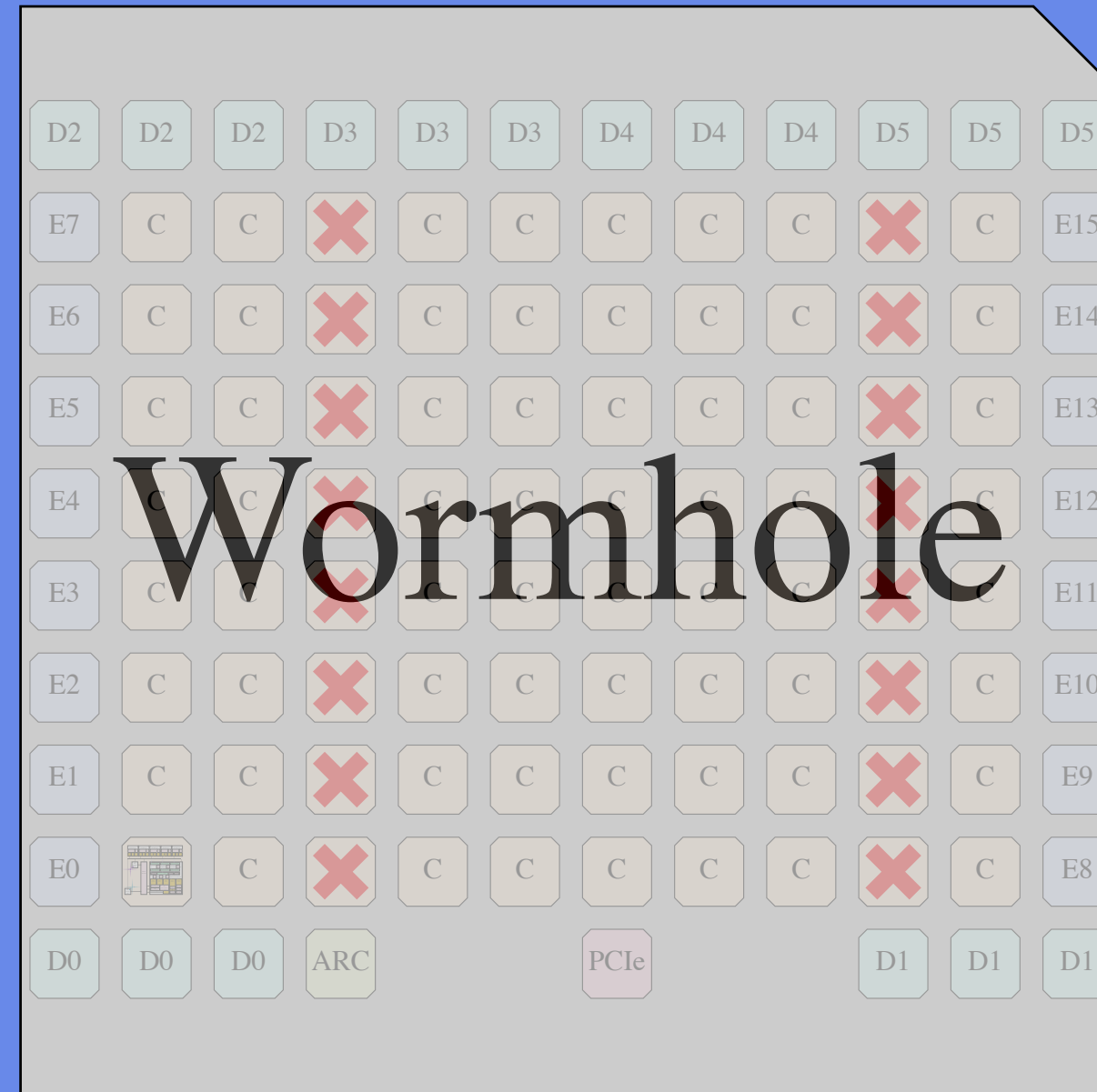
2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6



Wormhole

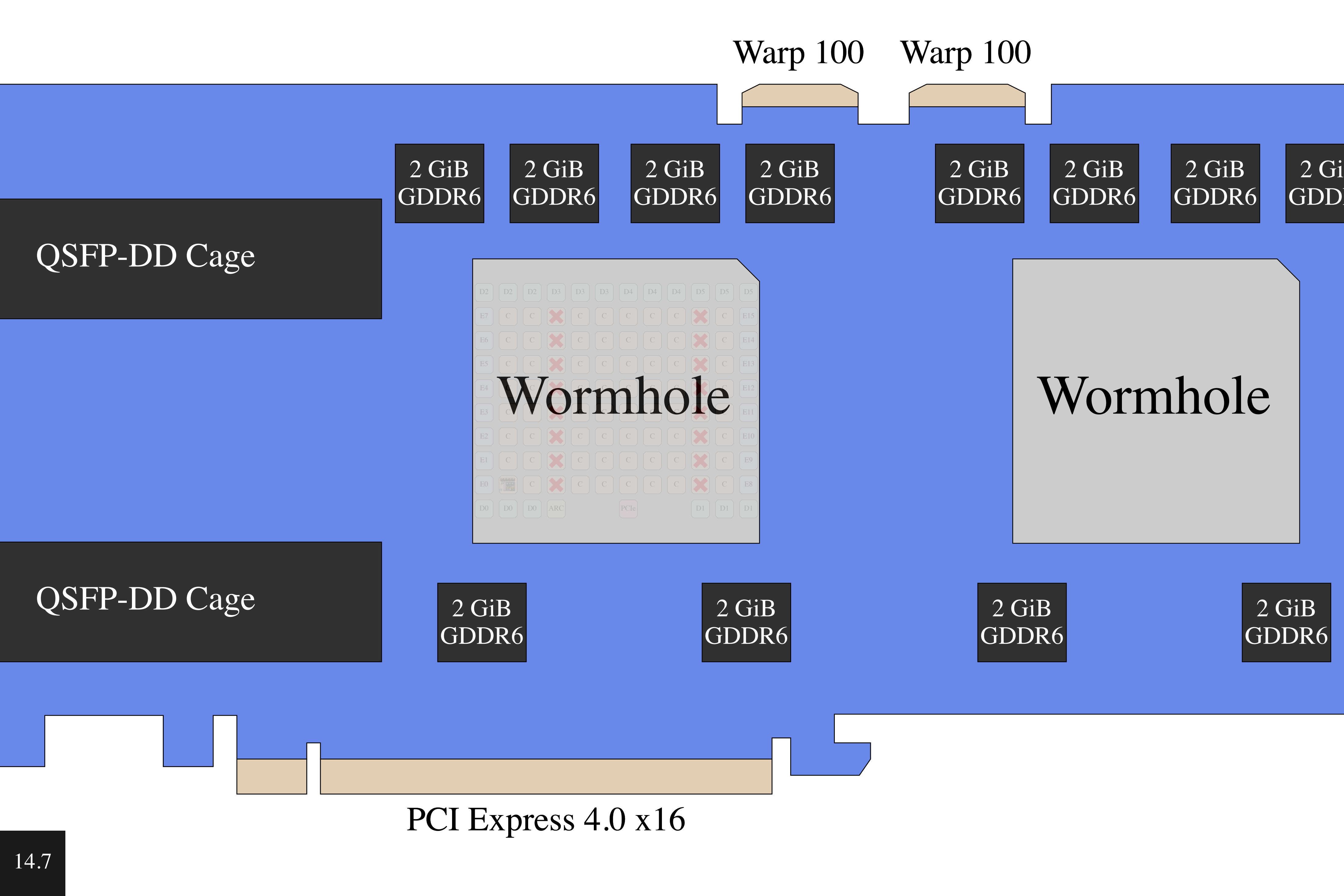
Worm

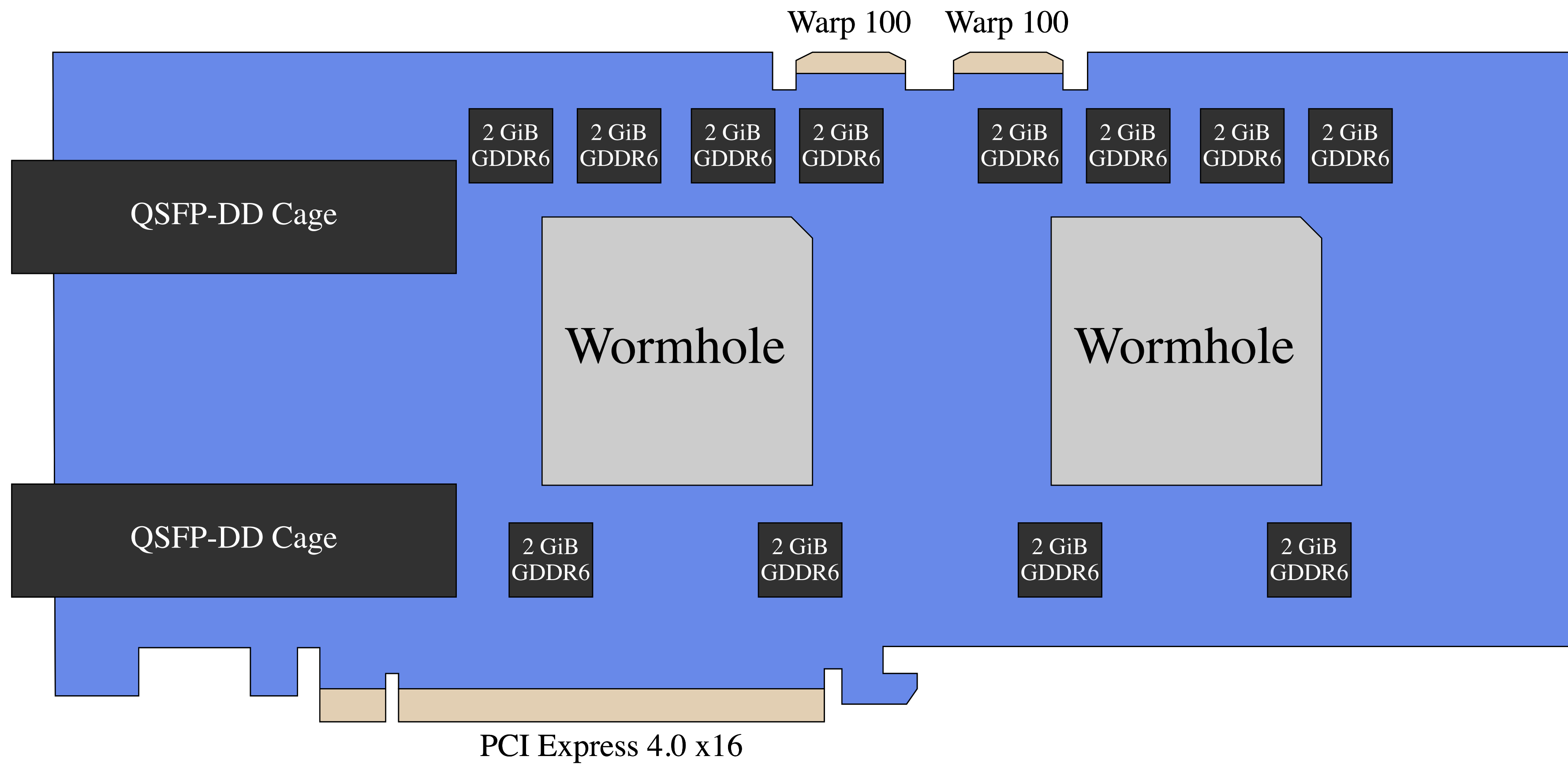
P-DD Cage

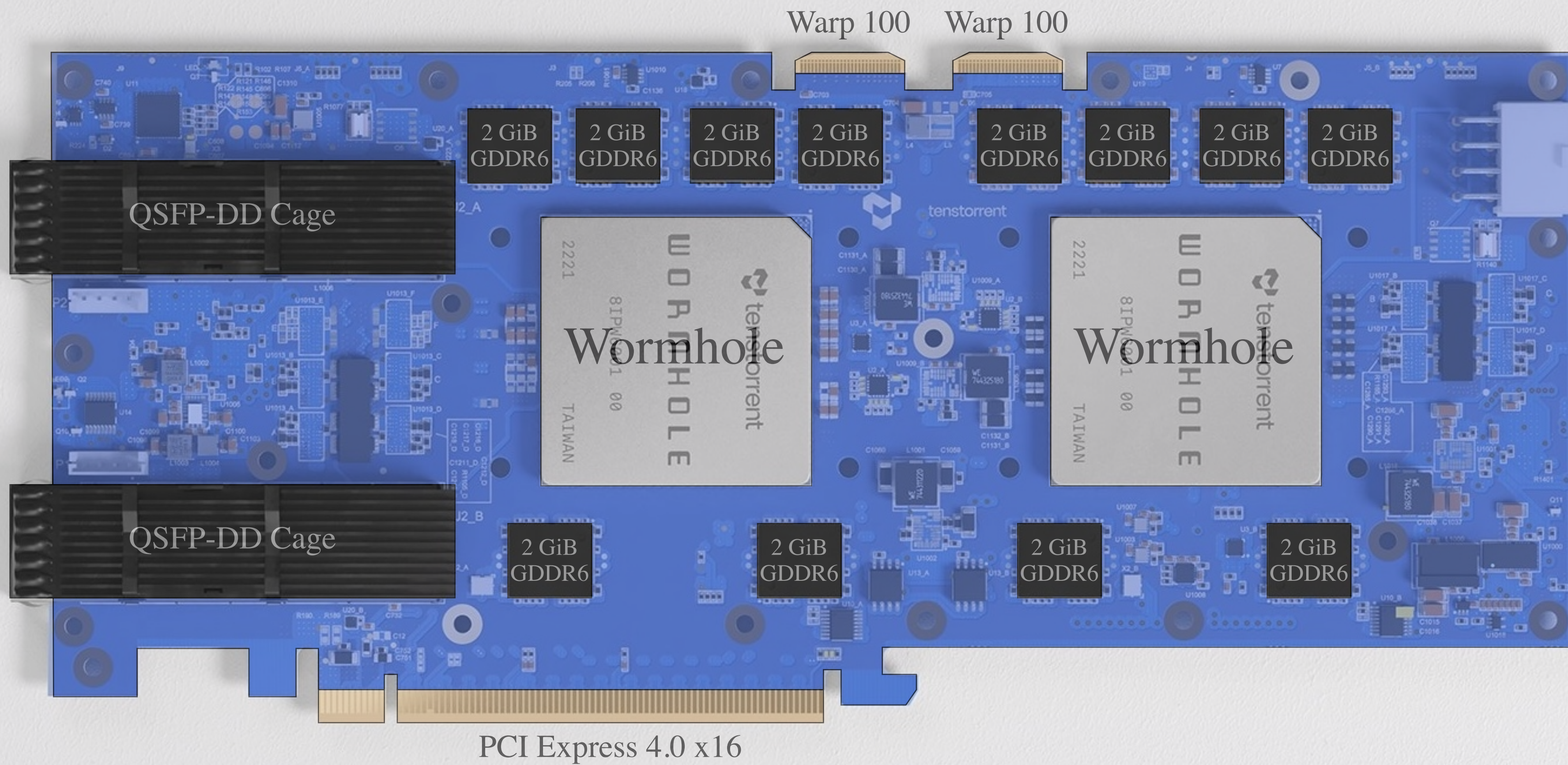
2 GiB
GDDR6

2 GiB
GDDR6

2 GiB
GDDR6







Warp 100 Warp 100

QSFP-DD Cage

2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6

Wormhole
tenstorrent
8IPW01 00
TAIWAN
2221

Wormhole
tenstorrent
8IPW01 00
TAIWAN
2221

QSFP-DD Cage

2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6 2 GiB GDDR6

PCI Express 4.0 x16

