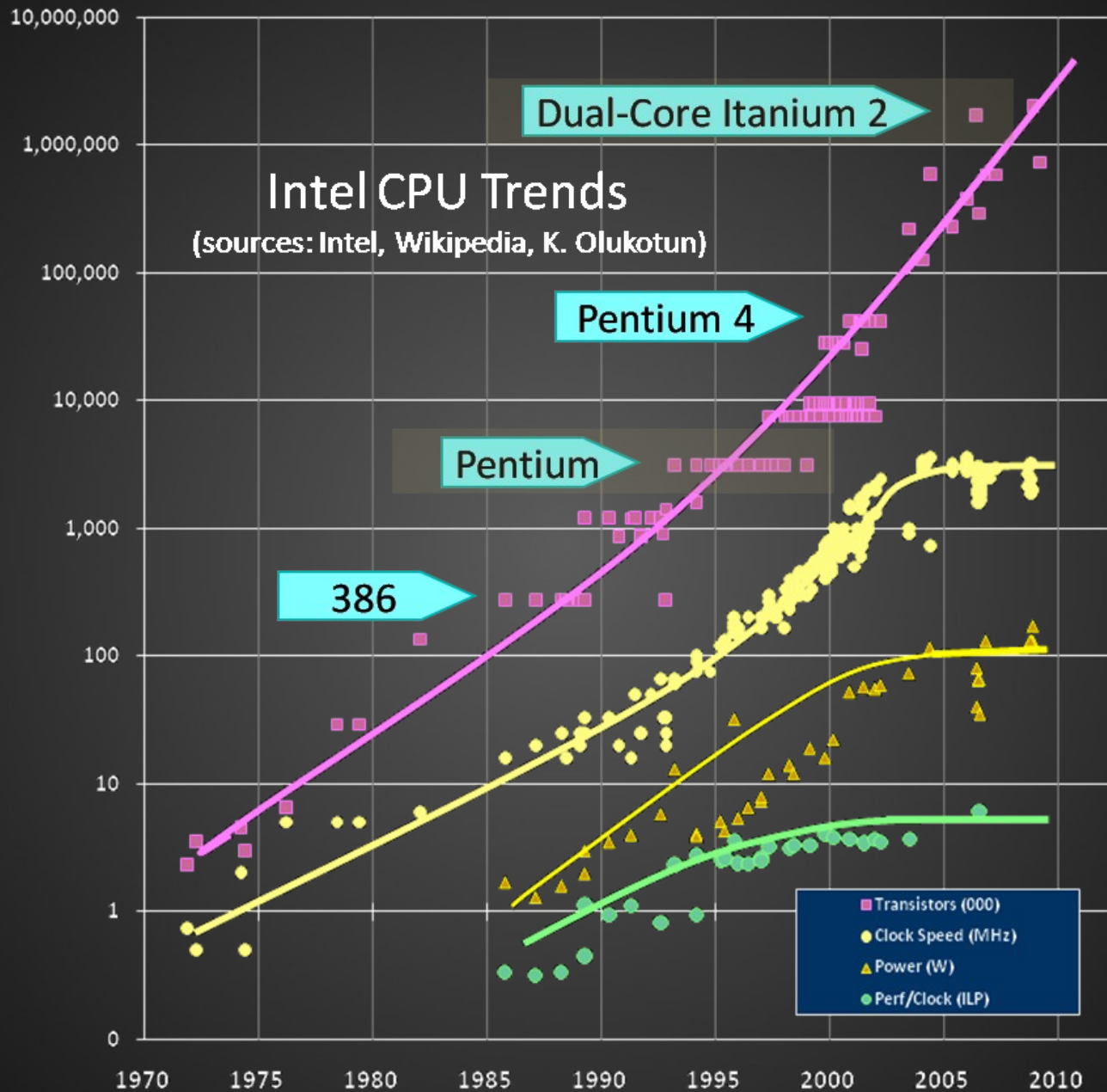




Διπλωματική Εργασία

Τροποποίηση της υλοποίησης διαδικτυακών πρωτοκόλλων λειτουργικού συστήματος για την εκμετάλλευση πολλαπλών επεξεργαστικών μονάδων





Το πρώτο βήμα: η BGL

• Ένα μόνο νήμα στον πυρήνα κάθε φορά



Το πρώτο βήμα: η BGL

Ένα μόνο νήμα στον πυρήνα κάθε φορά
Κλιμάκωση μόνο για αμιγώς υπολογιστικές
εφαρμογές



Το επόμενο βήμα: περισσότερα κλειδώματα

- 1. Εντοπισμός συχνά διεκδικούμενου κλειδώματος
- 2. Αντικατάσταση με περισσότερα κλειδώματα



Προβλήματα

LOR

Κλιμάκωση προς τα κάτω

- Bus contention



Οργάνωση της στοίβας πρωτοκόλλων στο DragonFlyBED

Η πρώτη ύλη



Η ιστορία πρωτοκόλλων στο BDD

- Αρκετά γρήγορο
- Αξιόπιστο
- Τεκμηριωμένο
- 20+ χρόνια ανάπτυξης (το τελευταίο δεν είναι απαραίτητα καλό)



tcp_input()

4.4BSD (Net/2 Release): 1100 SLOC

DragonFly 1.0: 1900+ SLOC

2η σε κυκλωματική πολυπλοκότητα σε όλο τον πυρήνα

Πολύ εύθραυστη

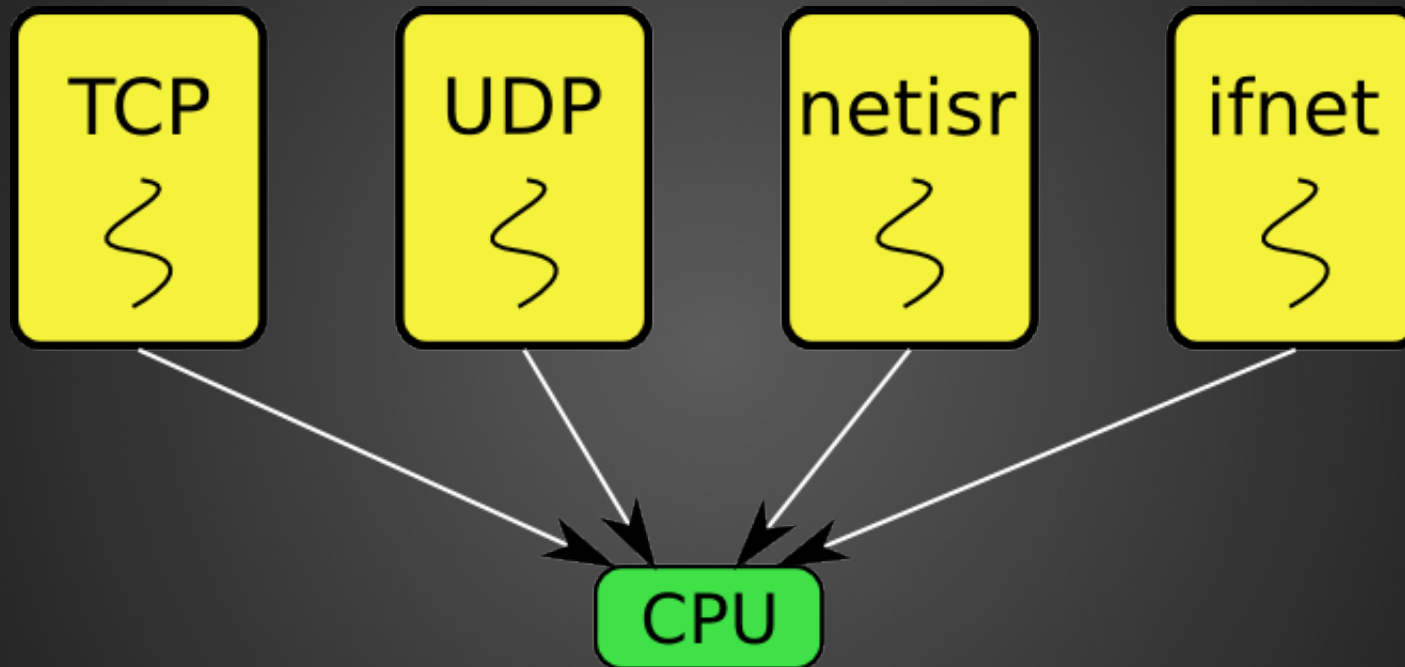


Οργάνωση της στοίβας πρωτοκόλλων στο DragonFlyBSD

Αλλαγές στο DragonFly



Τα νήματα



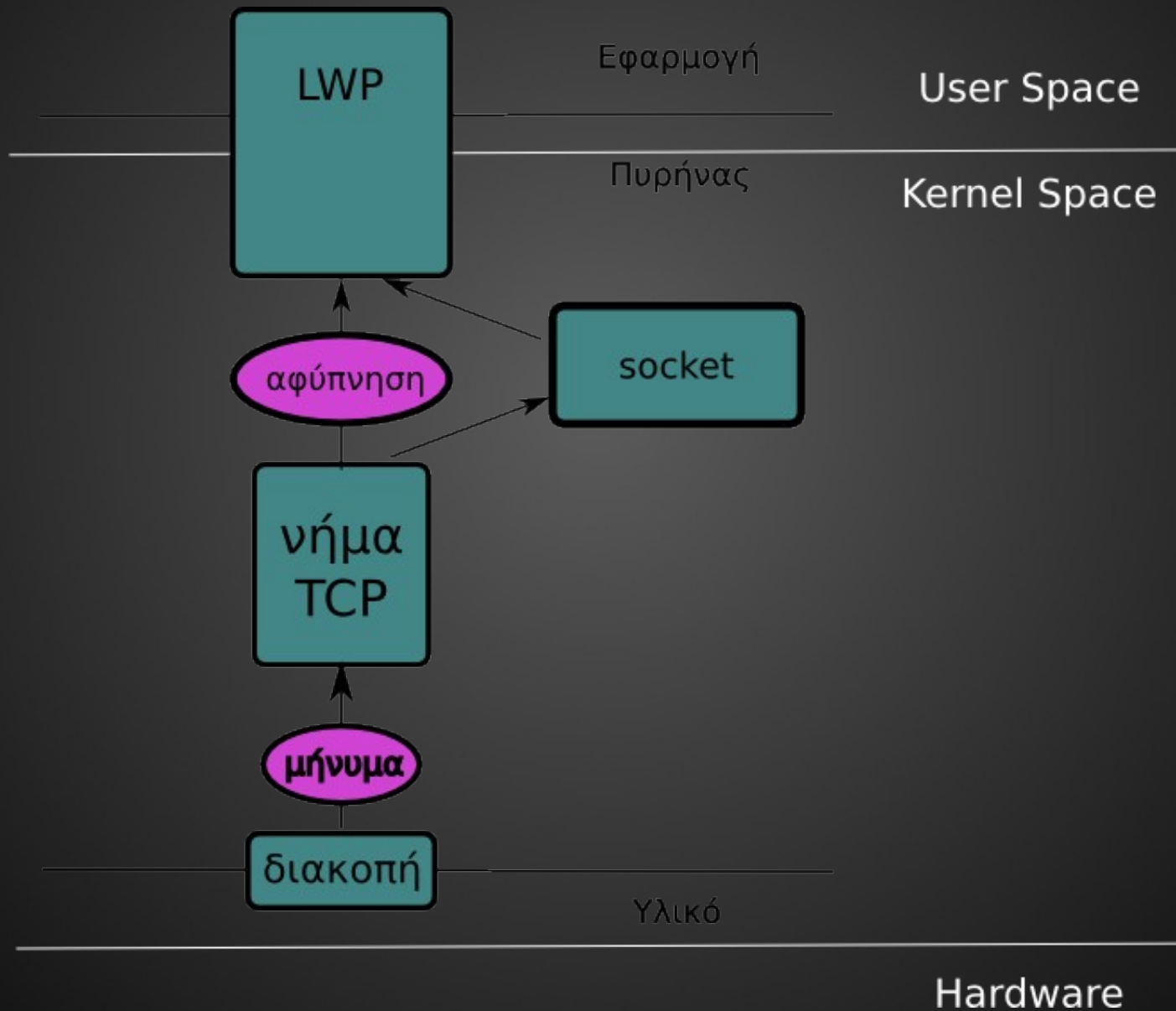


ΕΠΙΚΟΙΝΩΝΙΑ ΜΕ ΤΑ ΝΗΜΑΤΑ ΠΡΩΤΟΚΟΛΛΟΥ

- abort
- accept
- attach
- bind
- connect
- connect2
- control
- detach
- disconnect
- listen
- peeraddr
- recvd
- recvoob
- send
- κλπ

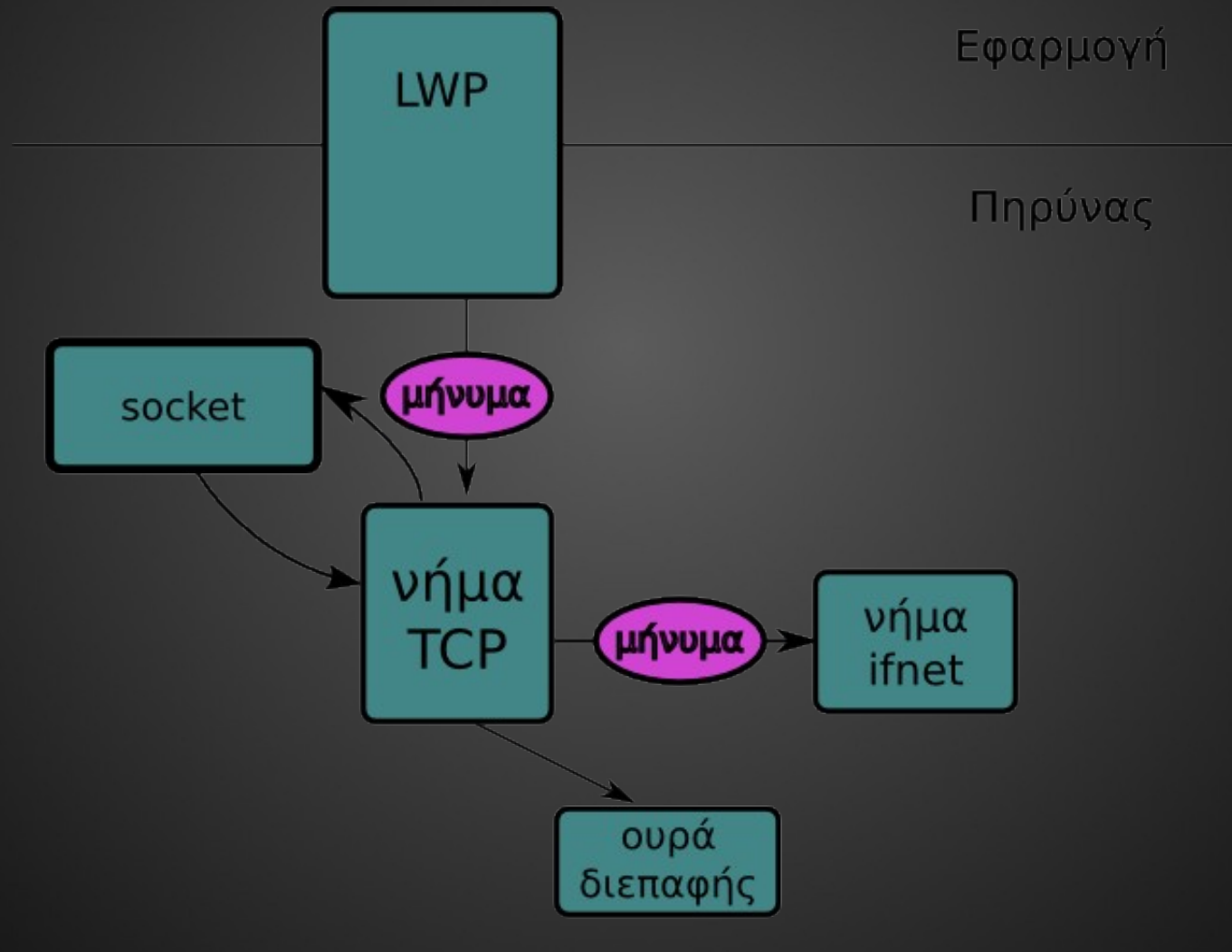


Λήψη δεδομένων





Αποστολή δεδομένων



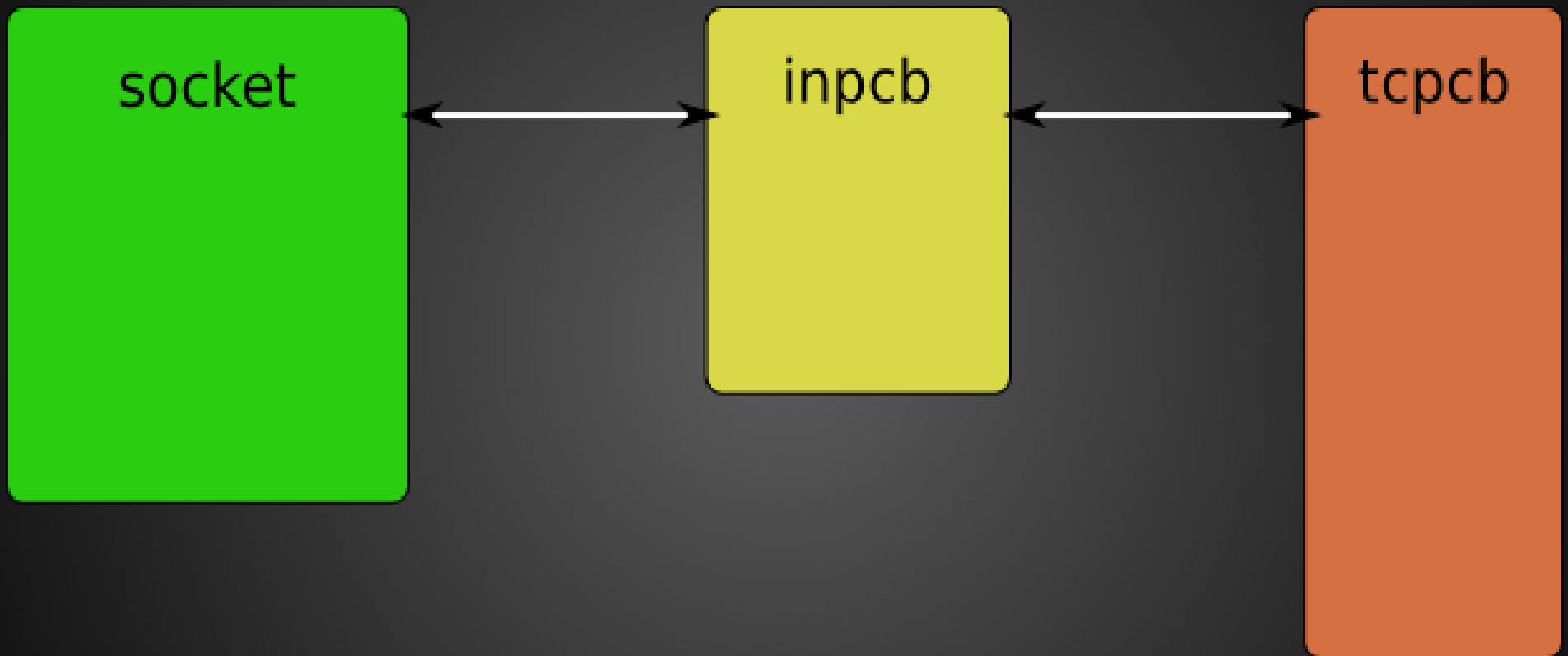


Ζητήματα (1)

Το PCB

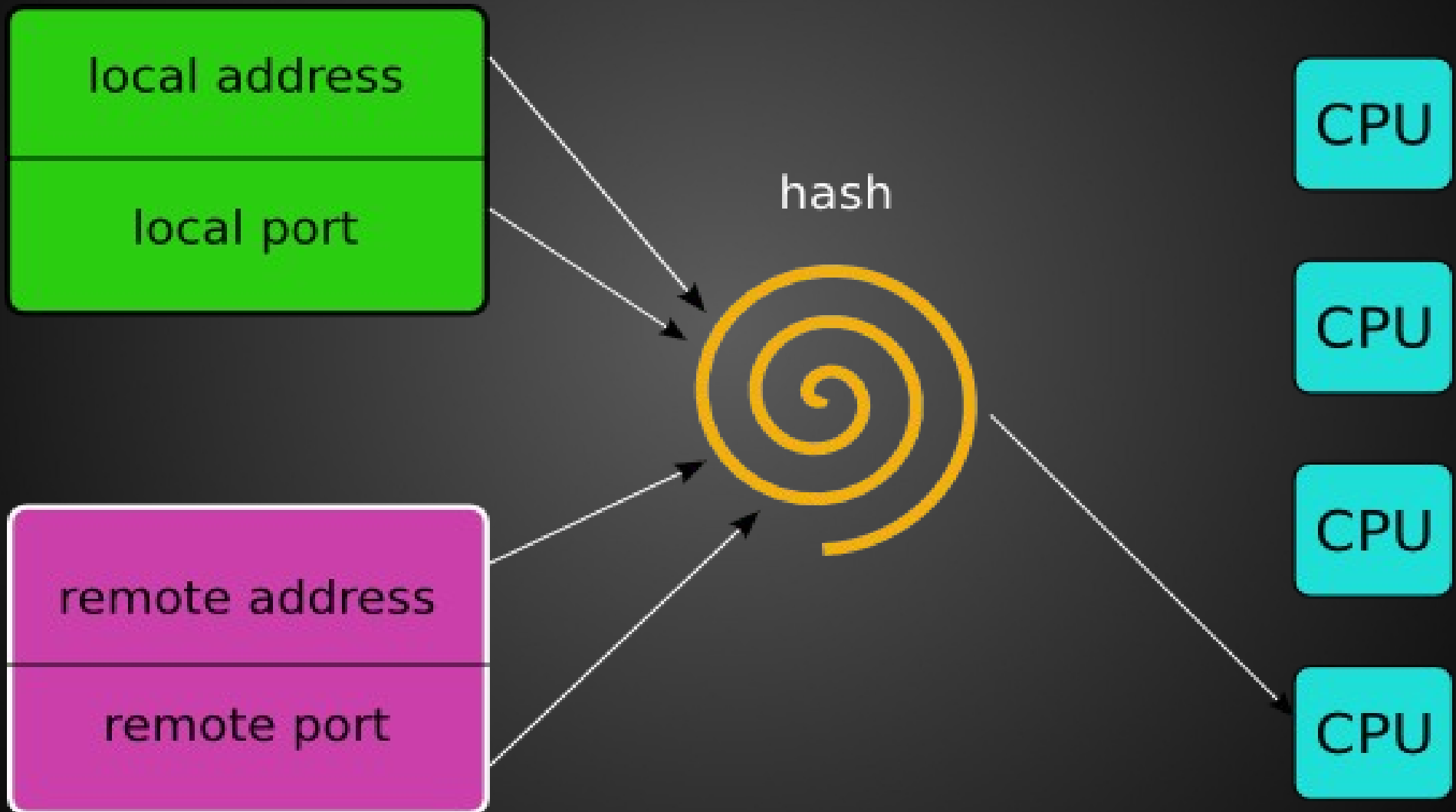


Protocol Control Block



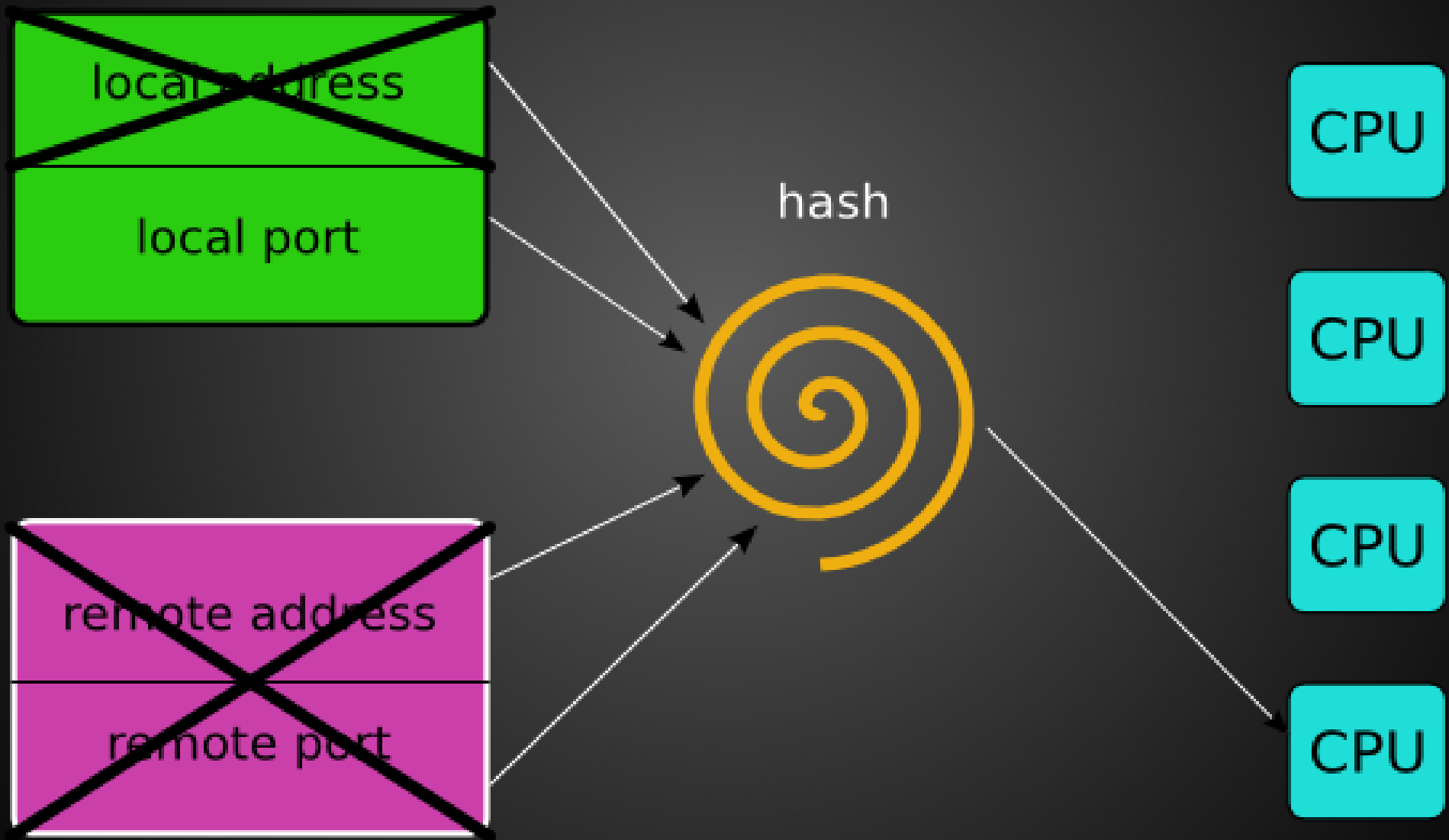


TCP hash





UDP hash





Λύσεις για το UDP

- Χρήση μόνο της τοπικής πόρτας
- Migration + tombstones
- Inpcb replication
 - NxN?



Ζητήματα (2)

Socketbuf



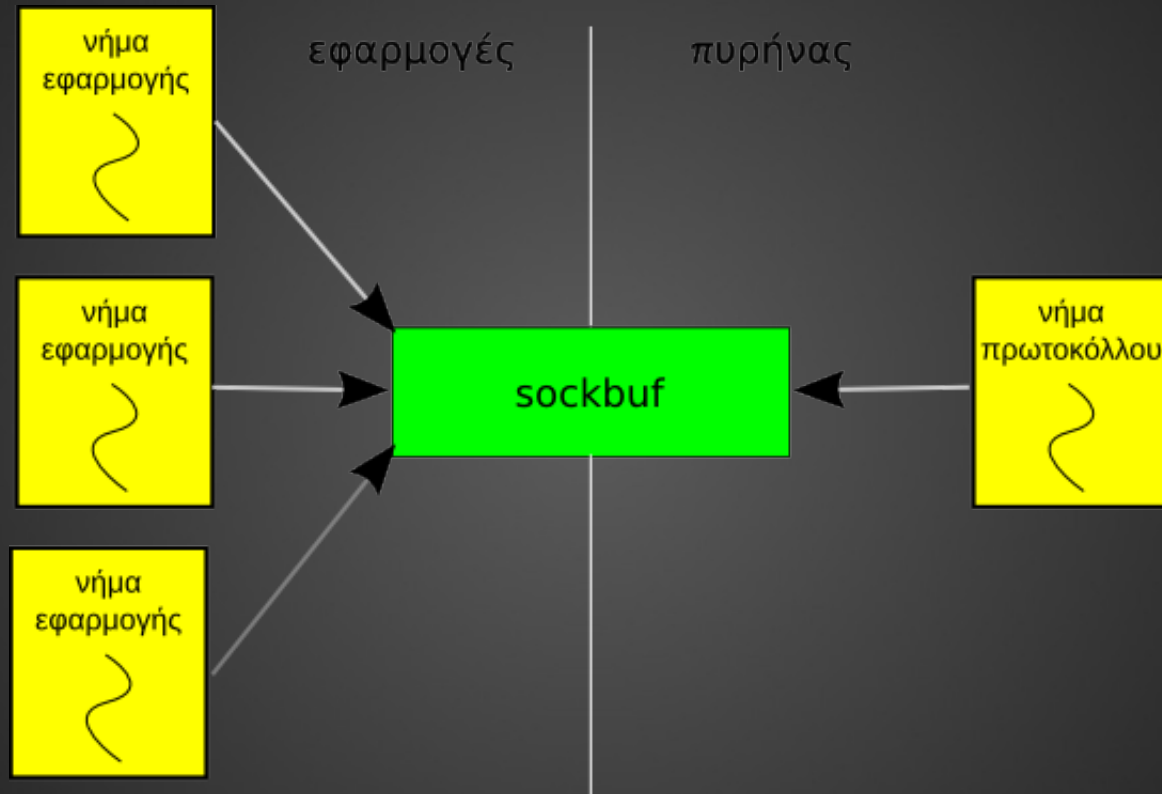
Sockbuf

Ενταμιευτής αποστολής (MPSC)

Ενταμιευτής λήψης (SPMC)

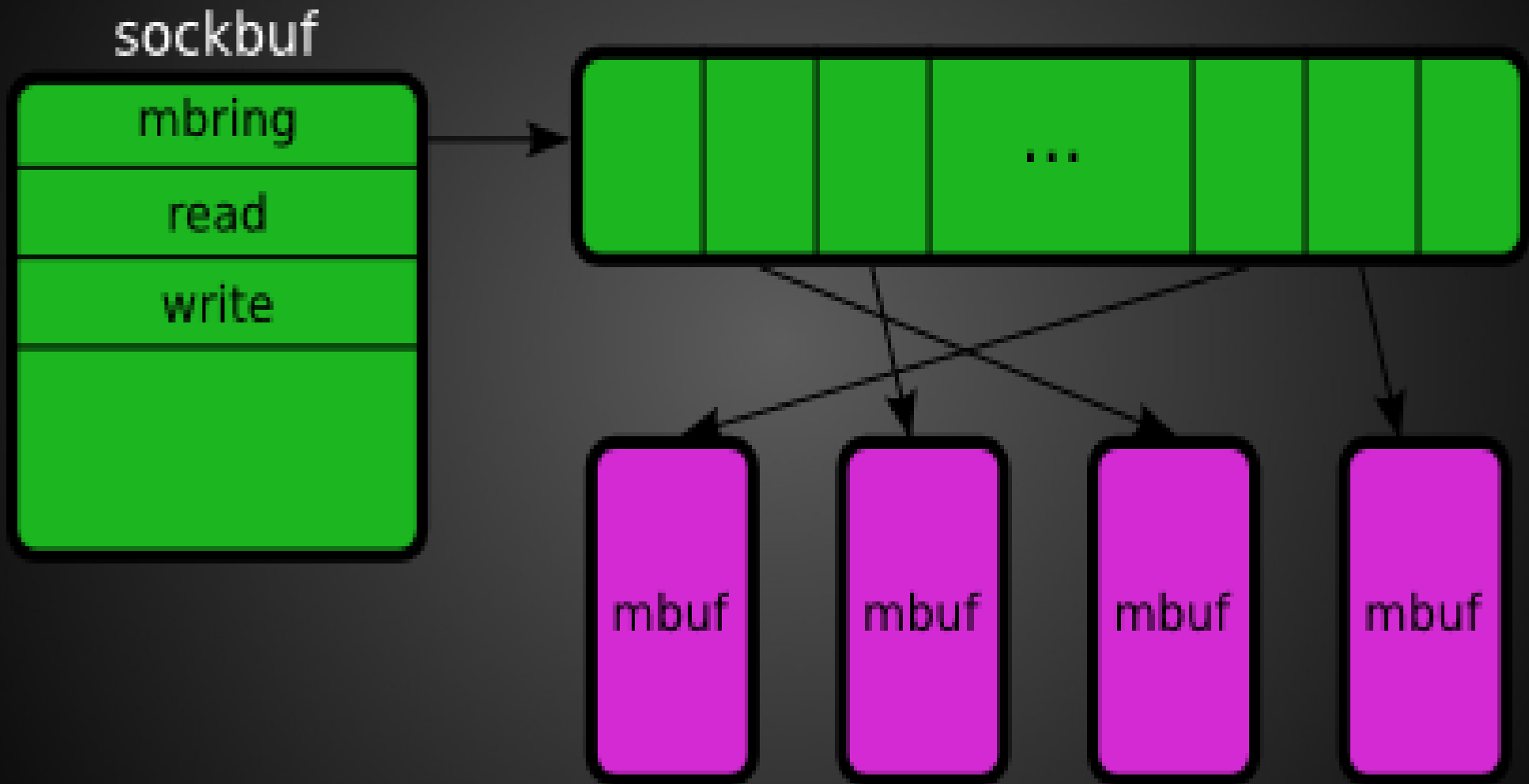


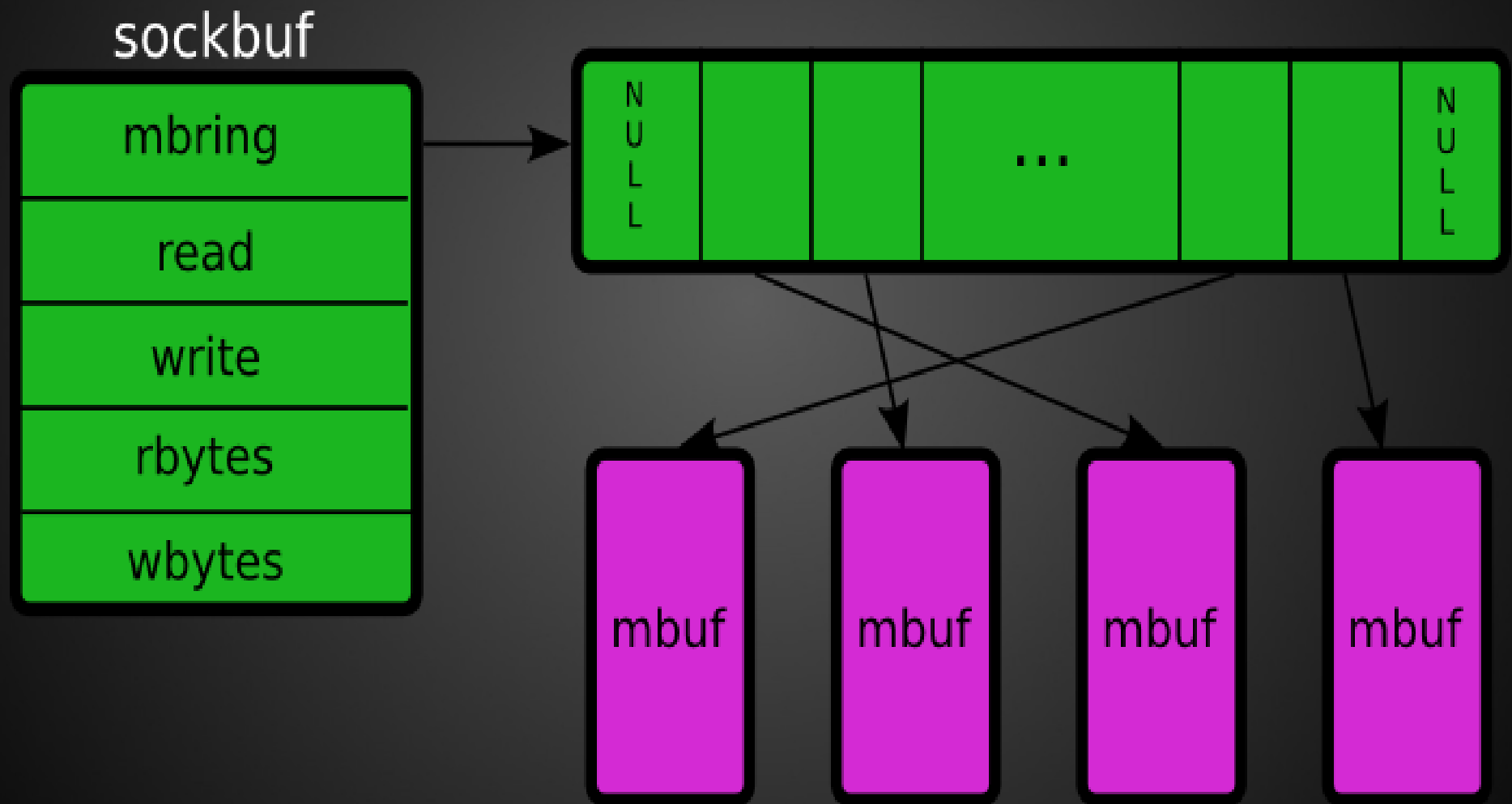
Draw me a picture





Lampport







+/-Κυκλικού Ενταμιευτή

Υπέρ

- Ελάχιστο cacheline thrashing
- Απλή υλοποίηση

Κατά

- Σπαταλά μνήμη

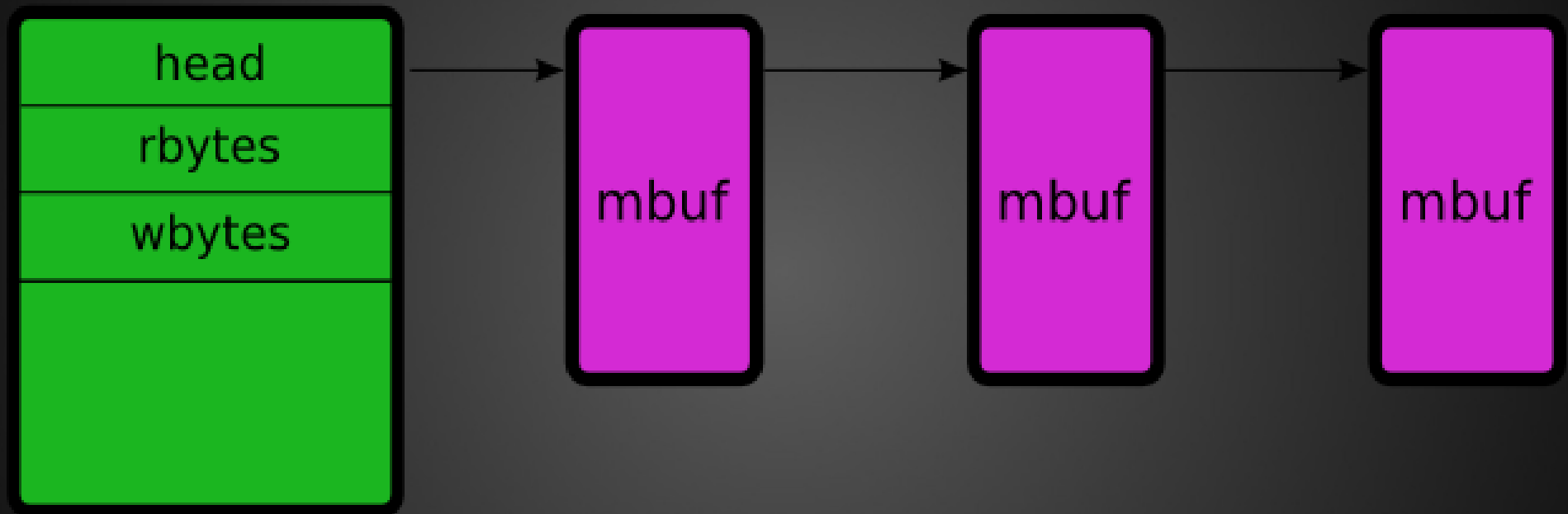
Προοπτική

- Μόνο με δυναμική αλλαγή μεγέθους



M_CORAL (1)

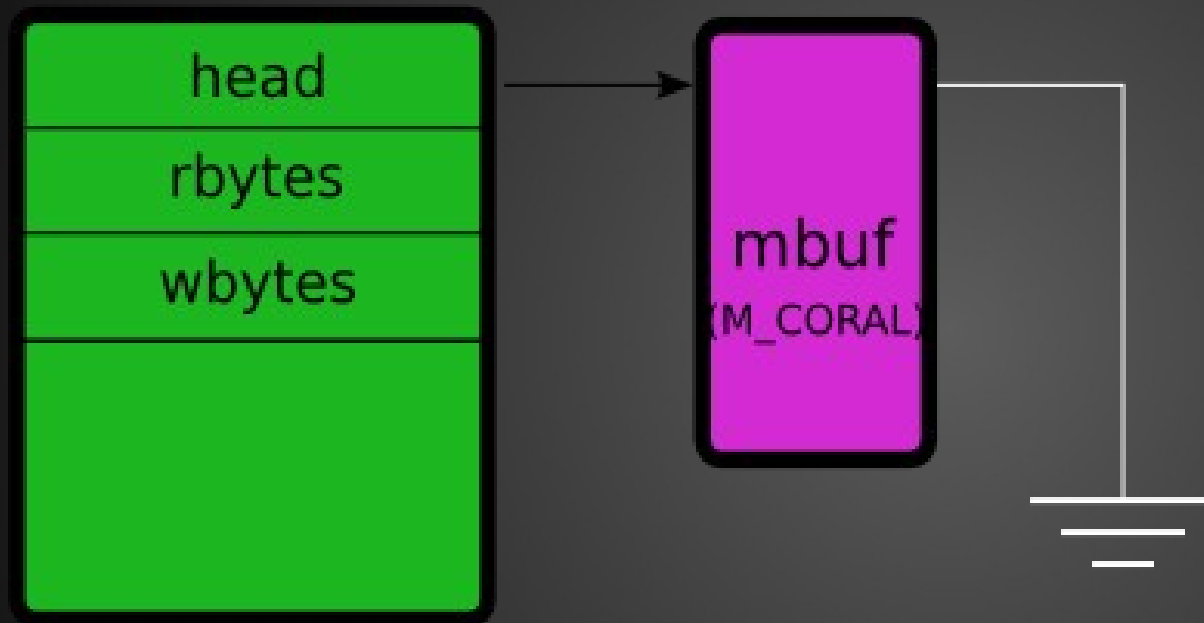
sockbuf





M_CORAL (2)

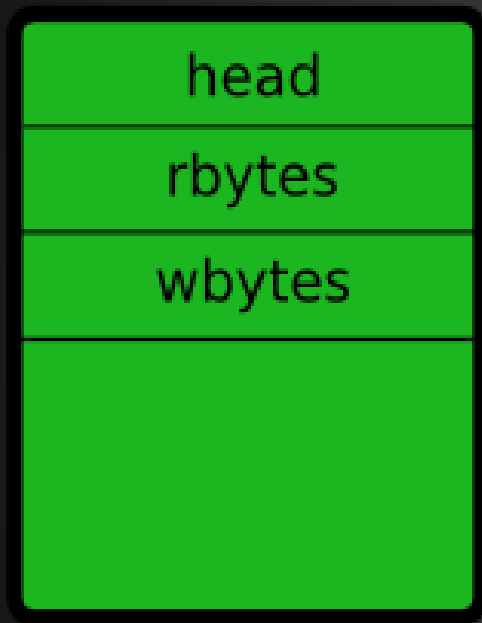
sockbuf





M_CORAL (3)

sockbuf





+/- M_CORAL

Γιέρ

- Δυναμικό μέγεθος
- Δεν σπαταλά μνήμη για ανενεργά sockbuf
- Cache-friendly

Κατά

- Τα mbuf clusters είναι μεγάλα (2K)
- Περίπλοκο CFG

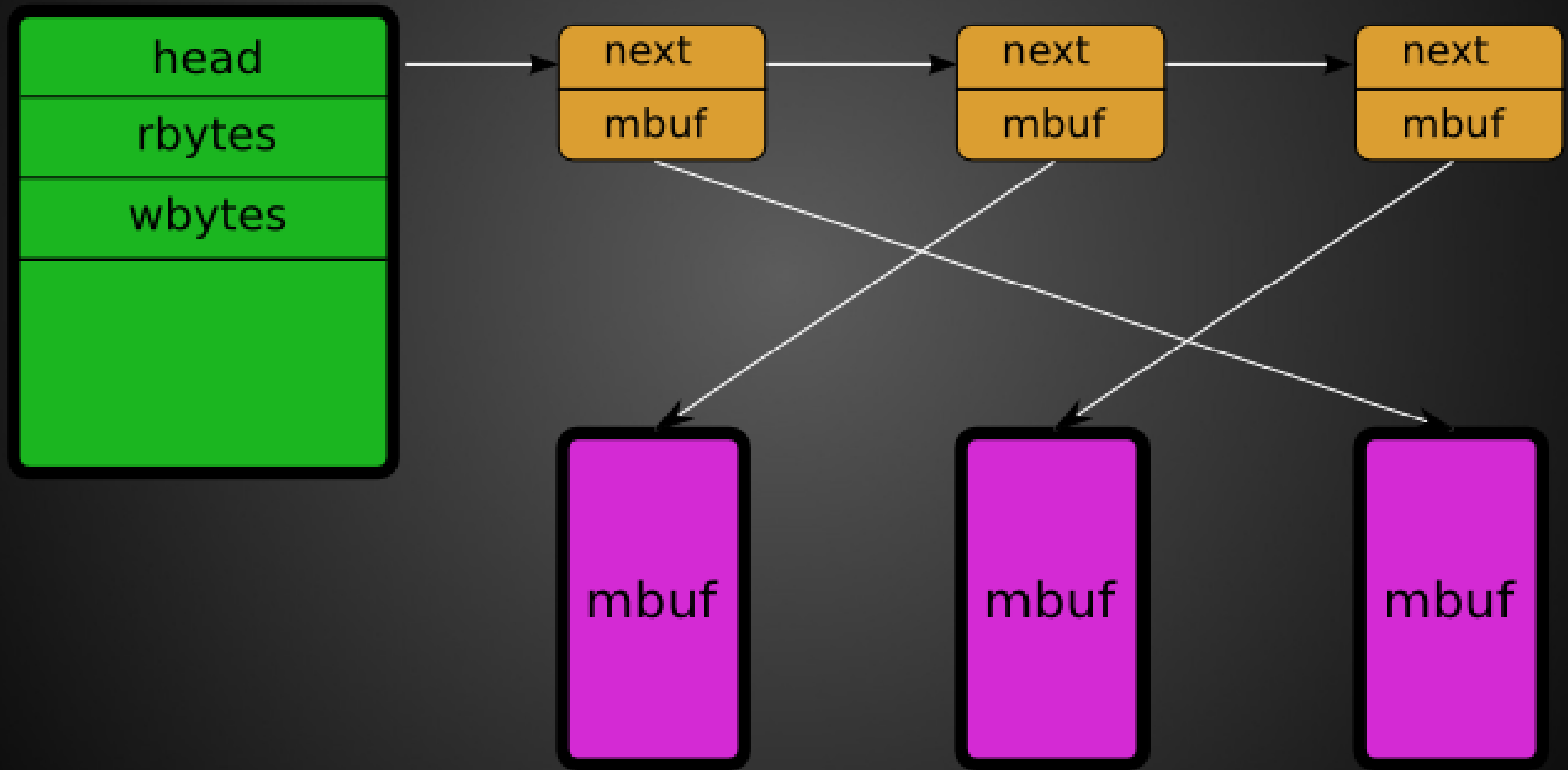
Προοπτική

- Θα χρησιμοποιηθεί μόλις λυθεί το πρόβλημα με τα mbuf clusters



cupholders

sockbuf





+/- cupholders

Υπέρ

- Δυναμικό μέγεθος
- Απλός κώδικας

Κατά

- Larger cache footprint
- Δέσμευση μνήμης

Προοπτική

Θα μείνουν για λίγο ακόμα



Απροσδιοριστία των δεδομένων

«Πόσα bytes περιέχονται στο sockbuf;»

Δεν μπορούμε να δώσουμε ένα νούμερο, αλλά μόνο ένα άνω ή κάτω όριο



Ζητήματα (3)

Races



Races

- Κάποια races τα αφήνουμε να εξελιχθούν
 - π.χ socket options
- Κάποια άλλα είναι επικίνδυνα
 - π.χ SS_CANTRCVMORE



Χαμένη αφύπνιση

```
pthread_mutex_lock(&mutex);
```

```
...]
```

```
again:
```

```
    if (have_data)
```

```
        break;
```

```
    if (exception)
```

```
        break;
```

```
    sleep_on_sockbuf();
```

```
    goto again;
```

```
again:
```

```
    if (have_data)
```

```
        break;
```

```
    if (exception)
```

```
        break;
```

```
    syncmsg(notify_me);
```

```
    goto again;
```



Μετρήσεις

Θεωρητικά εύκολο να δείξουμε την κλιμάκωση αλλά...

Ιδέες:

- Πολλές συνδέσεις με μικρή ροή δεδομένων
- loopback + hacks
- Χρήση 10G



10G

mxge(4)

- + HW csum, vlan tagging
- - RSS, TSO, jumbo frames



Αποτελέσματα

