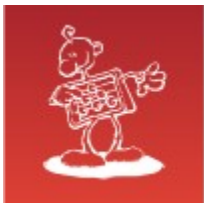


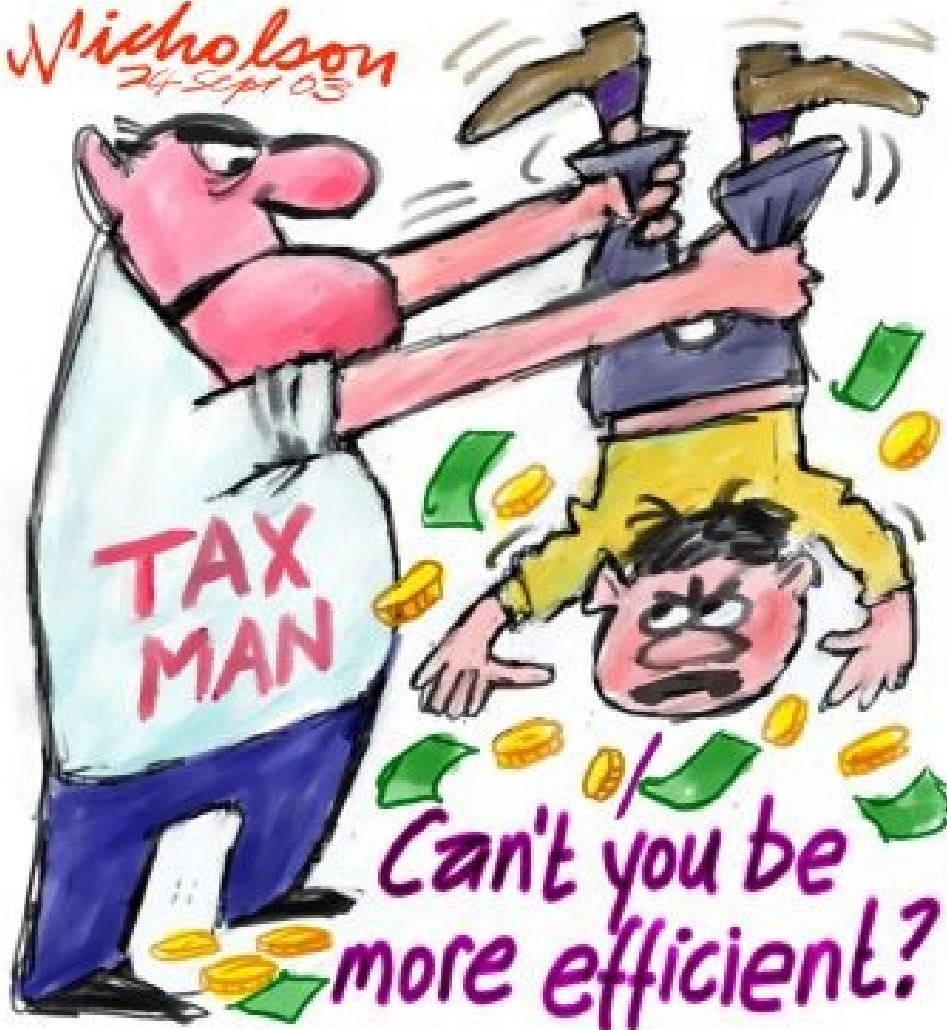
Vrhunsko zmogljiv hardver na odprtokodni osnovi

- **Sergej Rožman**, univ.dipl.inž.; Abakus plus d.o.o.
- Zadnja verzija z morebitnimi spremembami je na naslovu:
<http://www.abakus.si/>





Nicholson
24 Sept 03



Tax Department found inefficient





Vrhunsko zmogljiv hardver na odprtokodni osnovi

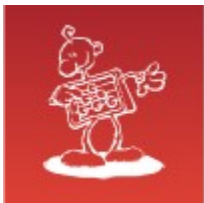
Sergej Rožman, univ.dipl.inž.

sergej.rozman@abakus.si



Poslovna Linux Konferenca
28. in 29. septembra 2009
Portorož | Slovenija





O podjetju

Zgodovina

- od 1992, 20 zaposlenih

Lastne aplikacije:

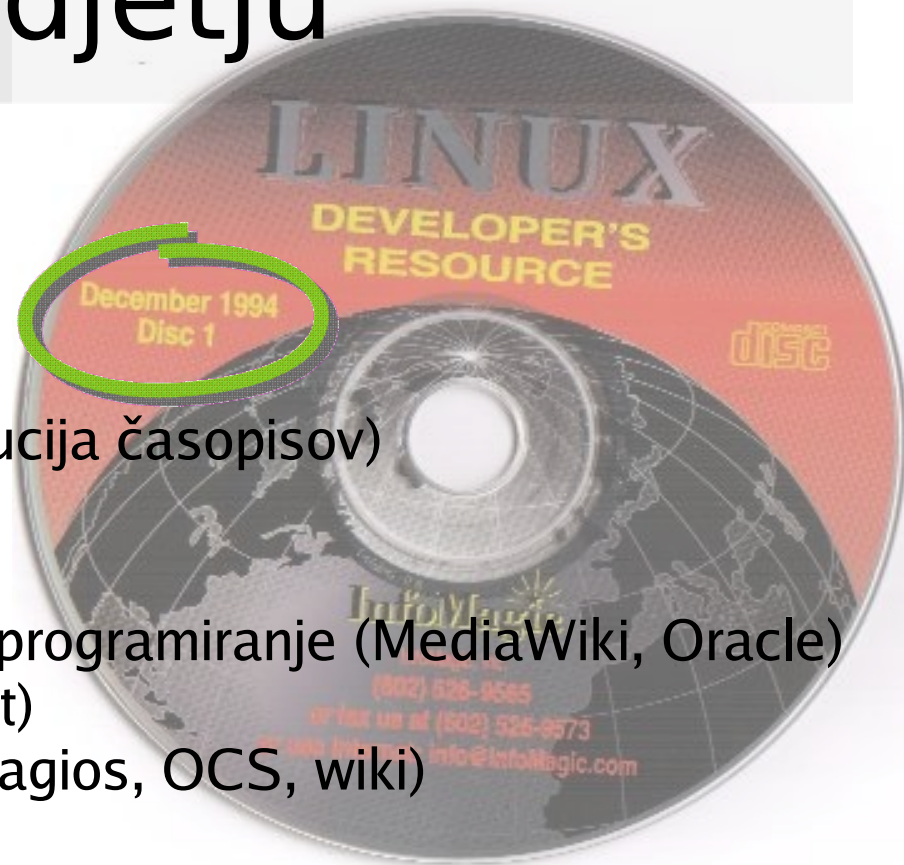
- posebne (letalski prometni sistem, distribucija časopisov)

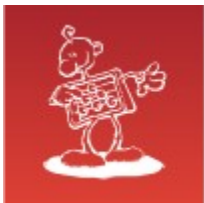
Storitve:

- DBA, vzdrževanje operacijskih sistemov, programiranje (MediaWiki, Oracle)
- omrežja (omr. storitve, VPN, QoS, varnost)
- odportokodne rešitve, nadzorni sistemi (nagios, OCS, wiki)

Okolje:

- od 1995 GNU/linux **(14-let izkušenj !)**
- prenos Oracle na GNU/linux: RDBMS 7.1.5 in forms 3.0 **(pred Oraclom !)**
- **skoraj 20 let izkušenj s sistemi za visoko razpoložljivost !**





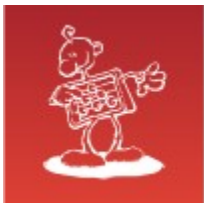
O podjetju



Slika: ekspedit v podjetju FuturaPlus v Beogradu

Aplikativna rešitev za distribucijo časopisnih edicij





O podjetju

Prometni sistem za letališča (slika iz aplikacije)



Srebrno priznanje
za inovacijo

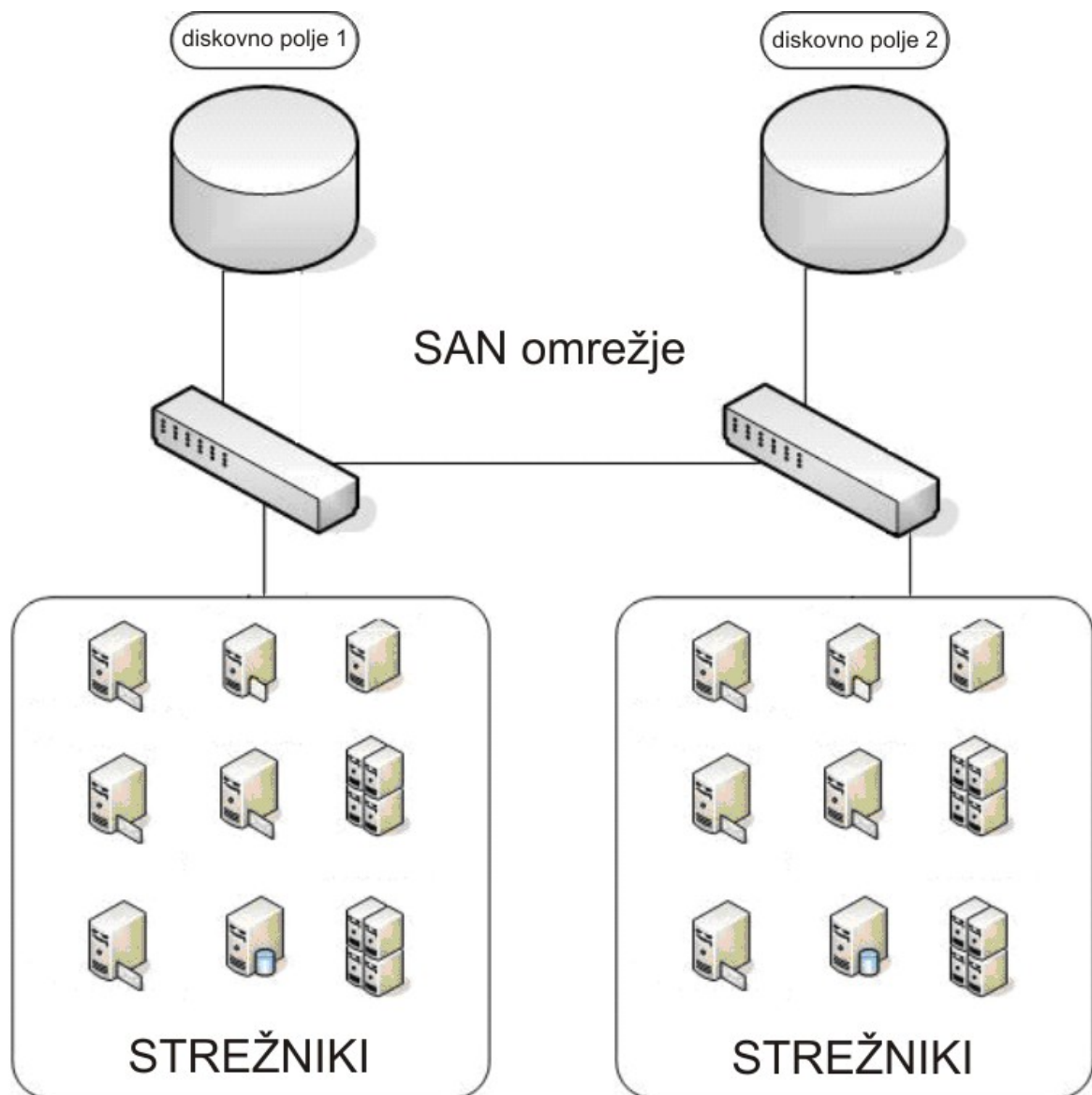


Zanimivost: strežniki na Aerodromu Ljubljana neprekinjeno delujejo že **1221** in **1136** dni (dne: 11. 9. 2009)



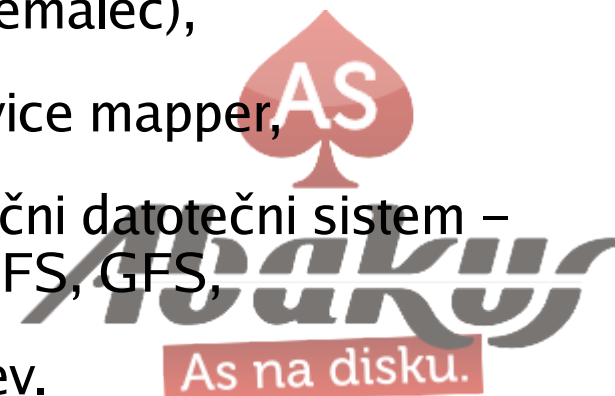


Cilj – zmogljiv hardver zasnovan na odprtih standardih in odprtokodnih tehnologijah



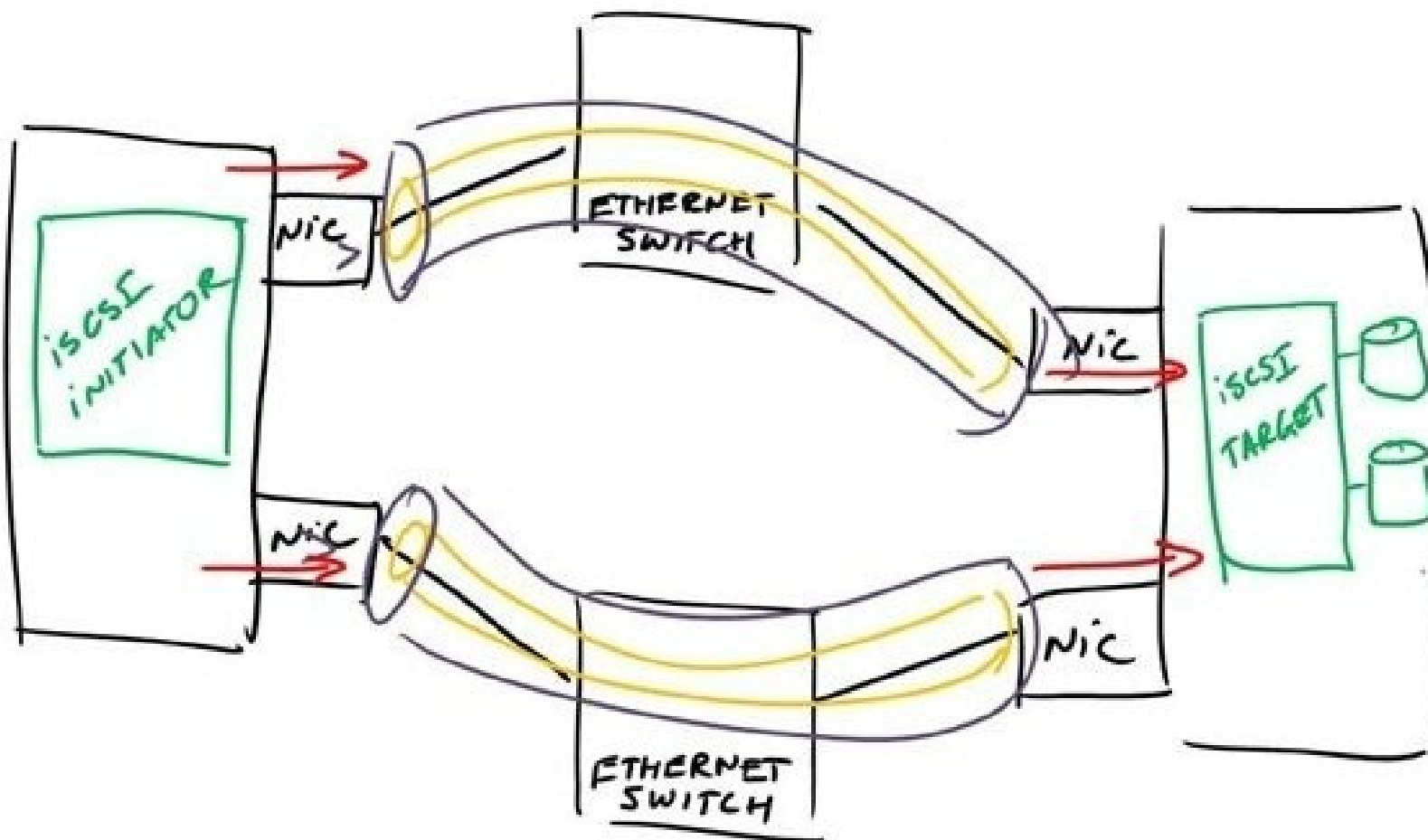
Uporabljene tehnologije:

- operacijski sistem GNU/linux,
- iSCSI enterprise target – IET (iSCSI strežnik),
- LVM,
- mrežno zrcaljenje – DRBD,
- iSCSI initiator (iSCSI odjemalec),
- device mapper,
- gručni datotečni sistem – OCFS, GFS,
- udev.



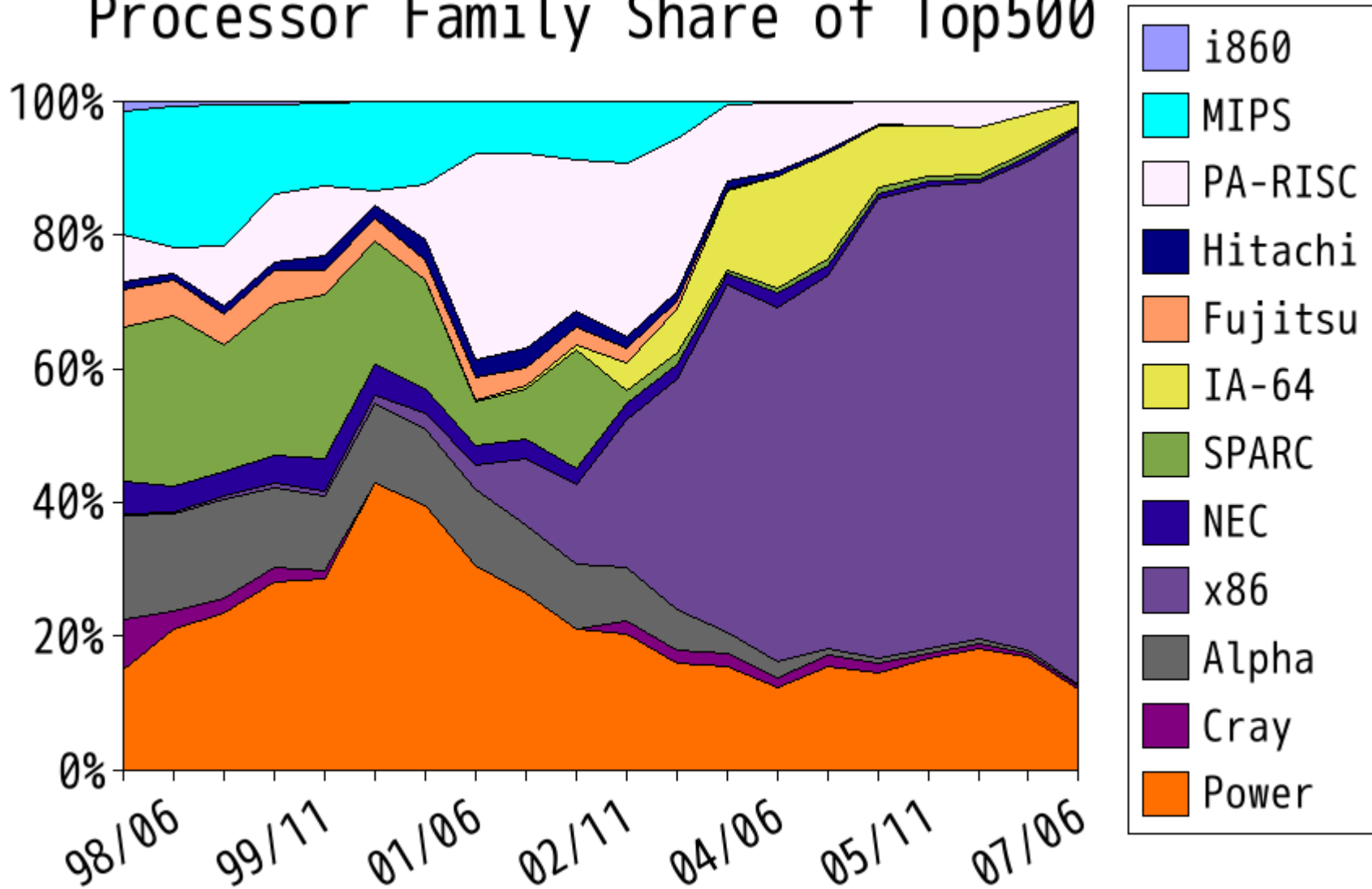


iSCSI





Processor Family Share of Top500

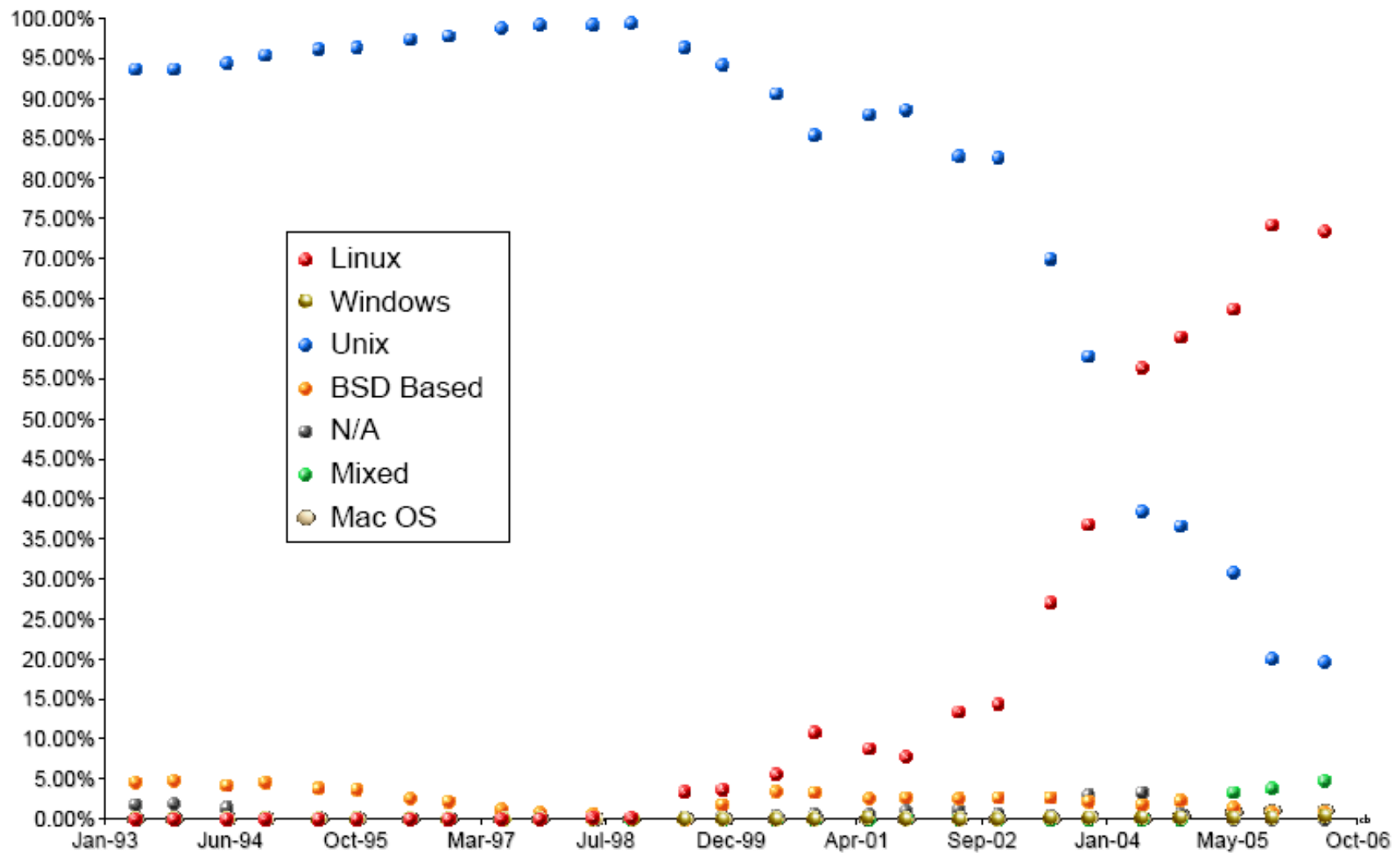


Source: <http://www.top500.org>

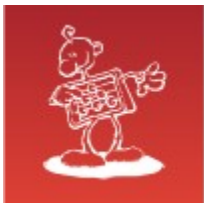




Operating Systems Used On Top500 Supercomputers



Source: www.top500.org



Primerjava lastnosti: fibre channel - iSCSI

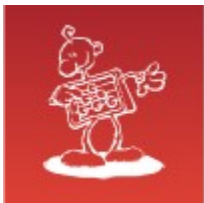
fibre channel

- posebna omrežna oprema
- ponavadi izvedba z namensko strojno opremo
- slabša združljivost naprav (ohlapni standardi)
- običajna izvedba z optičnimi vodi
- **počasnejši prenos**
- srednje razdalje na lastniškem omrežju
- draga rešitev
- v Sloveniji zelo razširjeno (HP EVA)

iSCSI (ethernet)

- standardna omrežna oprema
- v strojni ali programski izvedbi
- z združljivostjo ni težav
- običajna izvedba z bakrenimi vodi
- **hitrejši prenos**
- tudi velike razdalje, tudi po javnem omrežju
- poceni rešitev
- v Sloveniji ni poznano in razširjeno





Primerjava zmogljivosti: fibre channel - iSCSI

fibre channel

- 1 Gb/s (1997)
- 2 Gb/s (2001)
- 4 Gb/s (2005)
- 8 Gb/s (2008)
- 16 Gb/s (2011 plan.)

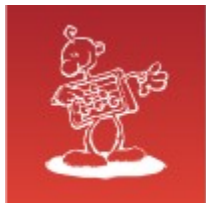
Vir: <http://www.fibrechannel.org/OVERVIEW/Roadmap.html>

iSCSI (ethernet)

- 10 Mb/s (1980)
- 100 Mb/s (1995)
- 1 Gb/s (f:1998/tp:1999)
- 10 Gb/s (f:2002/cx4:2004/tp:2006)
- 40/100 Gb/s (2010 plan.)

Vir: wikipedia





10 Gb ethernet

Bakreni vodi

- 10GBASE-CX4 (podoben kot InfiniBand) – do 15 m
- 10GBASE-T (cat. 6a) – do 100 m



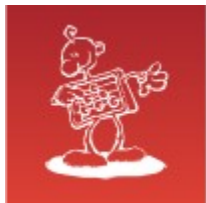
Optični vodi

- od 300 m do 80 km, odvisno od vodov (SM, MM) in optičnih modulov

Proizvajalci

- Myricom
- Chelsio
- Neterion
- NetXen





10 Gb ethernet

Izkušnje:

10 GbE iSCSI pri polni hitrosti polno obremeni dva procesorja ali dve procesorski jedri

Posebne nastavitve (nekaj primerov)

Parametri linux jedra (/proc)

```
net.ipv4.tcp_timestamps = 0      (izklop časovnih oznak)
net.ipv4.tcp_sack = 0           (izklop selektivnega potrjevanja)
net.*.*mem* = 10000000         (dovolj veliki izravnalni pomnilniki)
net.core.netdev_max_backlog = 300000
```

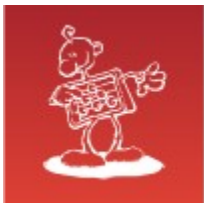
Mrežne nastavitve (ethtool)

```
več procesiranja na prekinitev (ethtool -C $DEVNAME rx-usecs $VALUE)
MTU = 9000                       »jumbo« okvirji
```

Sistemske nastavitve

```
TCP segmentation offload       (nastavitev gonilnika)
izklop irq_balancerja          (prekinitve vedno obravnava isti procesor)
```

Druge nastavitve po priporočilu proizvajalca kartic



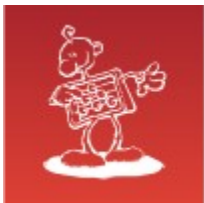
HARDVER





Osnovni model





Modularna zasnova



< 2U – 12 diskov



< 3U – 16 diskov



< 5U – 24 diskov



< 9U – 50 diskov

Moduli:

1U – krmilnik >

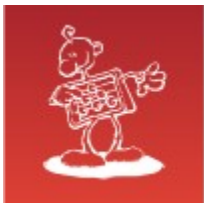


2U – 24 2,5" diskov >



4U – 48 2,5" diskov >





Redundantna zasnova



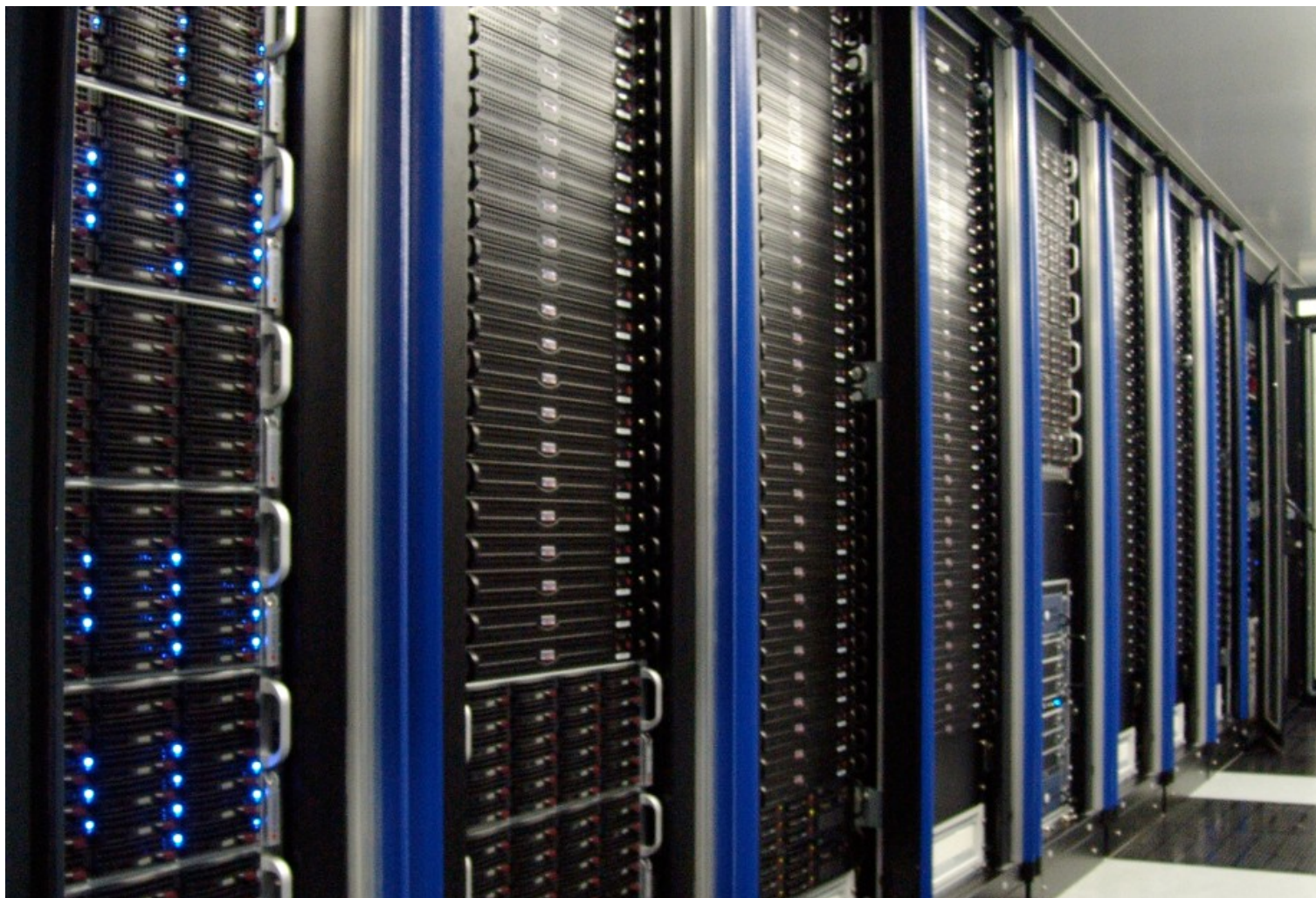
Dvojna povezava do istih diskov

- dvokanalni SAS
- multi-initiator SCSI





Sestavljeno skupaj



As na disku.



SOFTVER

SOFTVERIA
www.softveria.com

AS

Abakus

As na disku.

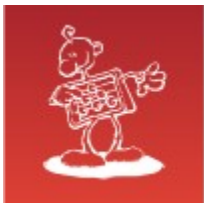


iSCSI enterprise target (IET)

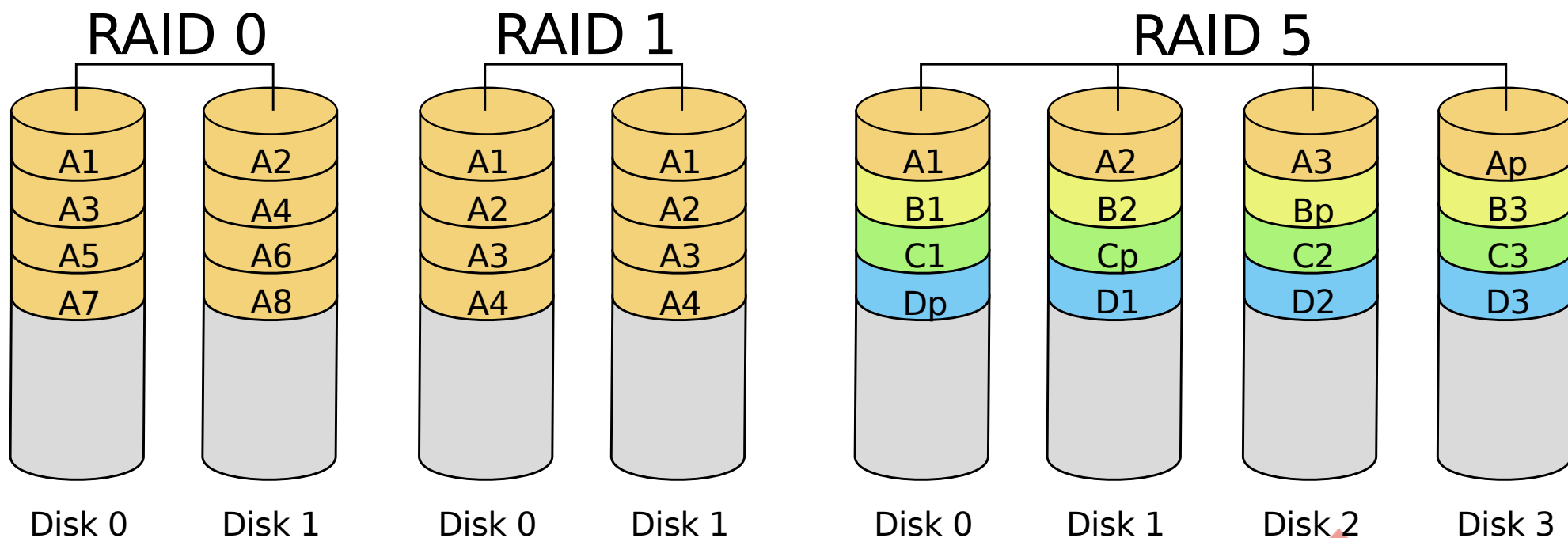
The iSCSI Enterprise Target Project

- RFC 3720
- <http://iscsitarget.sourceforge.net/>
- dinamično konfiguriranje (CLI)
- nima grafičnega vmesnika (GUI) **(še!)**





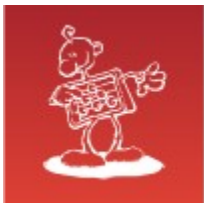
RAID (Redundant Array of Inexpensive Disks)



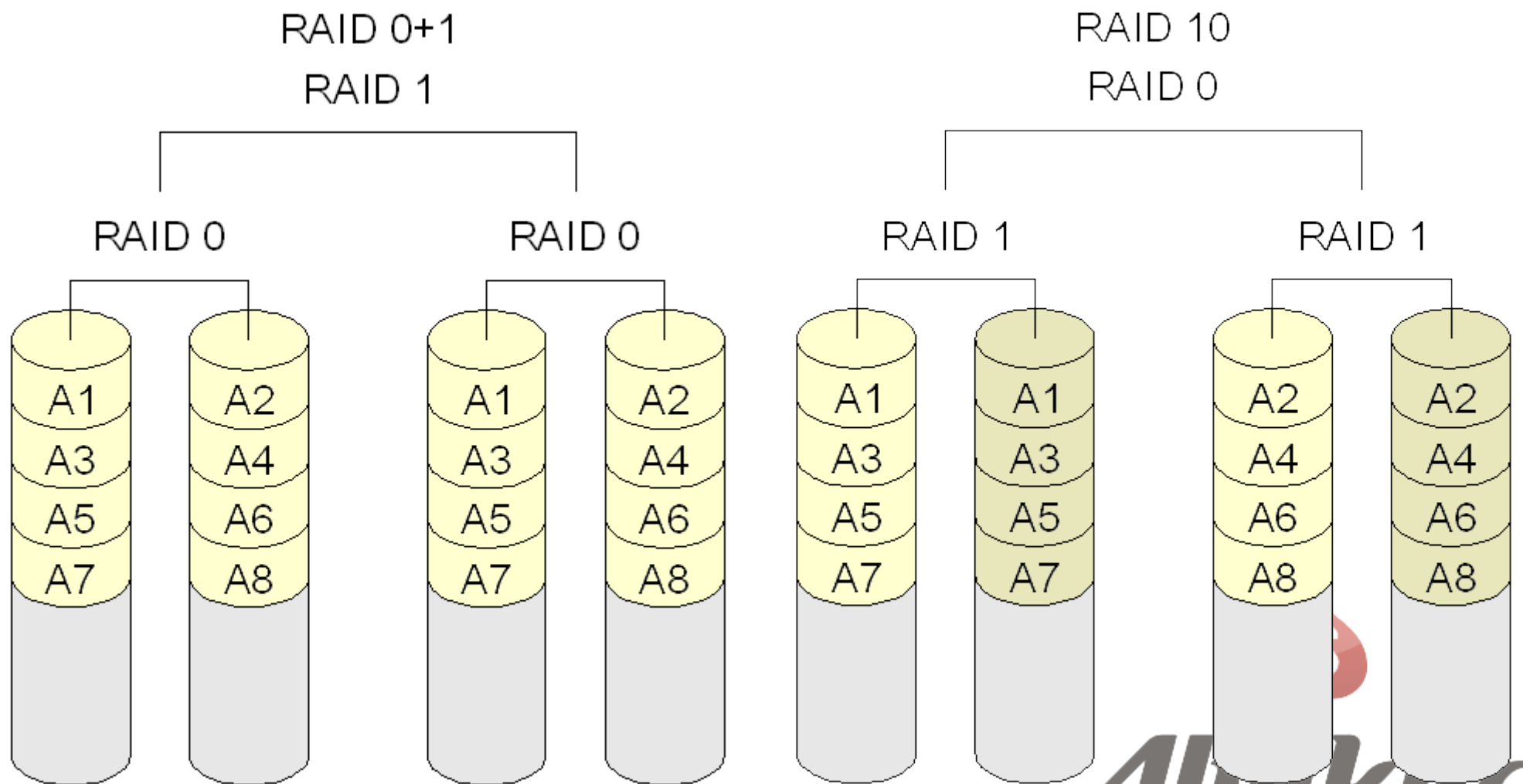
Vir: wikipedia

BAARF – <http://www.baarf.com/>
Battle Against Any Raid Five (Four, Free)





RAID - sestavljene stopnje



Vir: wikipedia



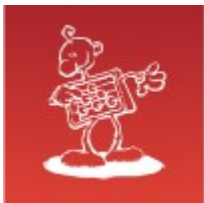
LVM (Logical Volume Management)

Katera funkcionalnost LVM je najbolj uporabna pri SAN?

- »SNAPSHOT«

LVM »snapshot« navidezno v trenutku zamrzne stanje particije, ki jo zato lahko konzistentno prekopiramo/shranimo



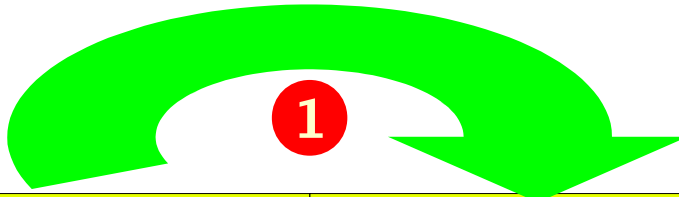


LVM (Logical Volume Management)

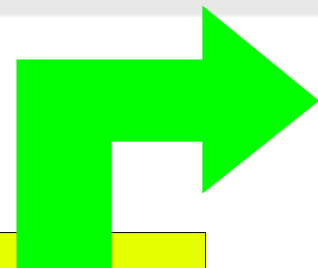
2






1

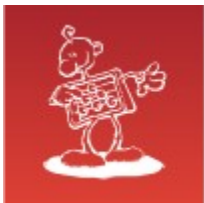


3



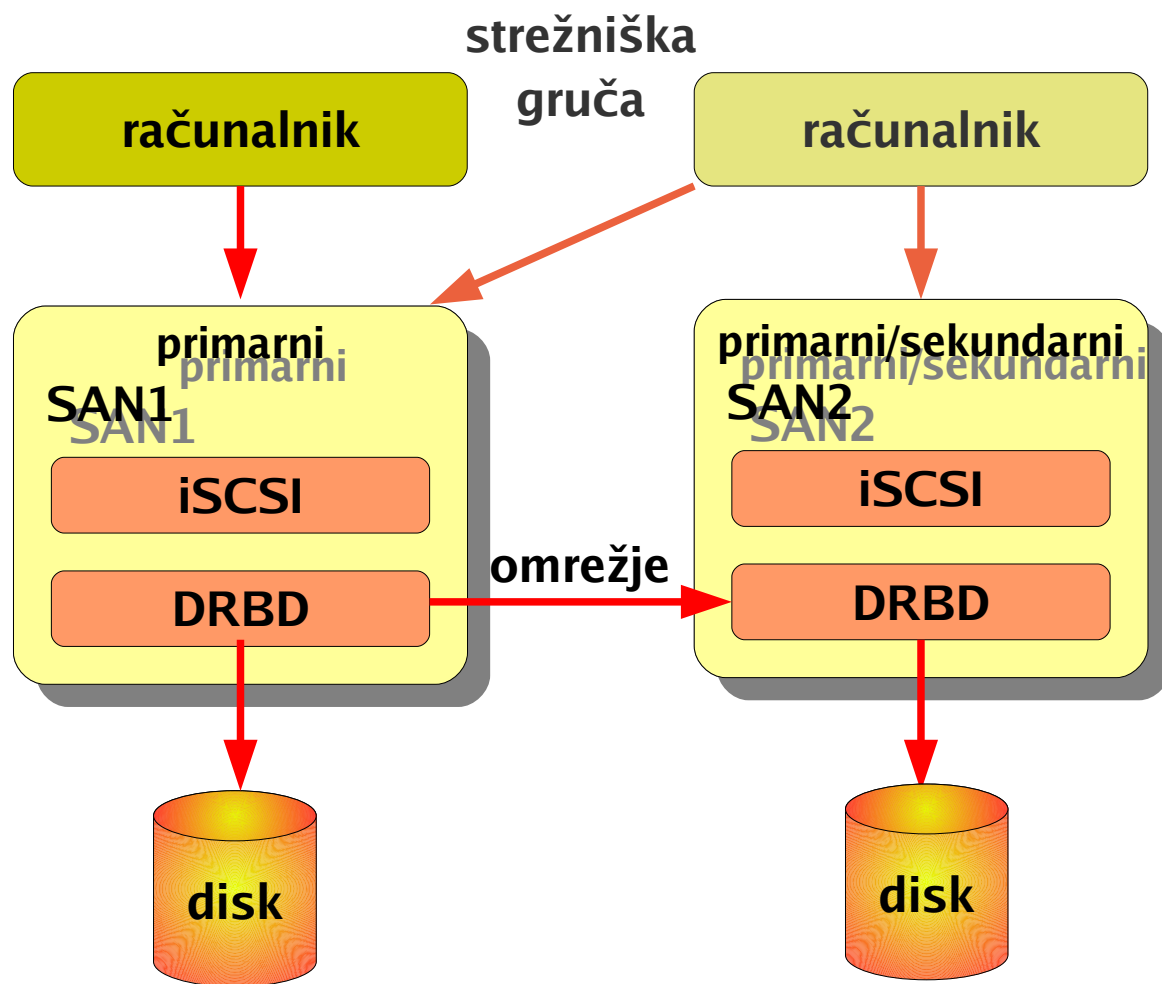
produksijski datotečni sistem /disk (OCFS2)		»zamrznjeni« datotečni sistem /backup (OCFS2)	
logični disk/particija (LV) 500 GB		»snapshot« logični disk/particija (LV) 200 GB	
logična grupa (LG) diskov 3 x 300 GB			
fizični disk (PV) 300 GB	fizični disk (PV) 300 GB	fizični disk (PV) 300 GB	
 disk1	 disk2	 disk3	





DRBD - mrežno zrcaljenje

(Distributed Replicated Block Device - LINBIT)



- sinhrona replikacija
- asinhrona replikacija
- dvosmerna sinhrona replikacija!
(od DRBD v. 8.0 dalje)

Podobna funkcionalnost: Oracle ASM (Automatic Storage Management) na oddaljen SAN – samo za potrebe Oracle zbirke podatkov



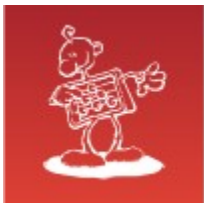


iSCSI initiator

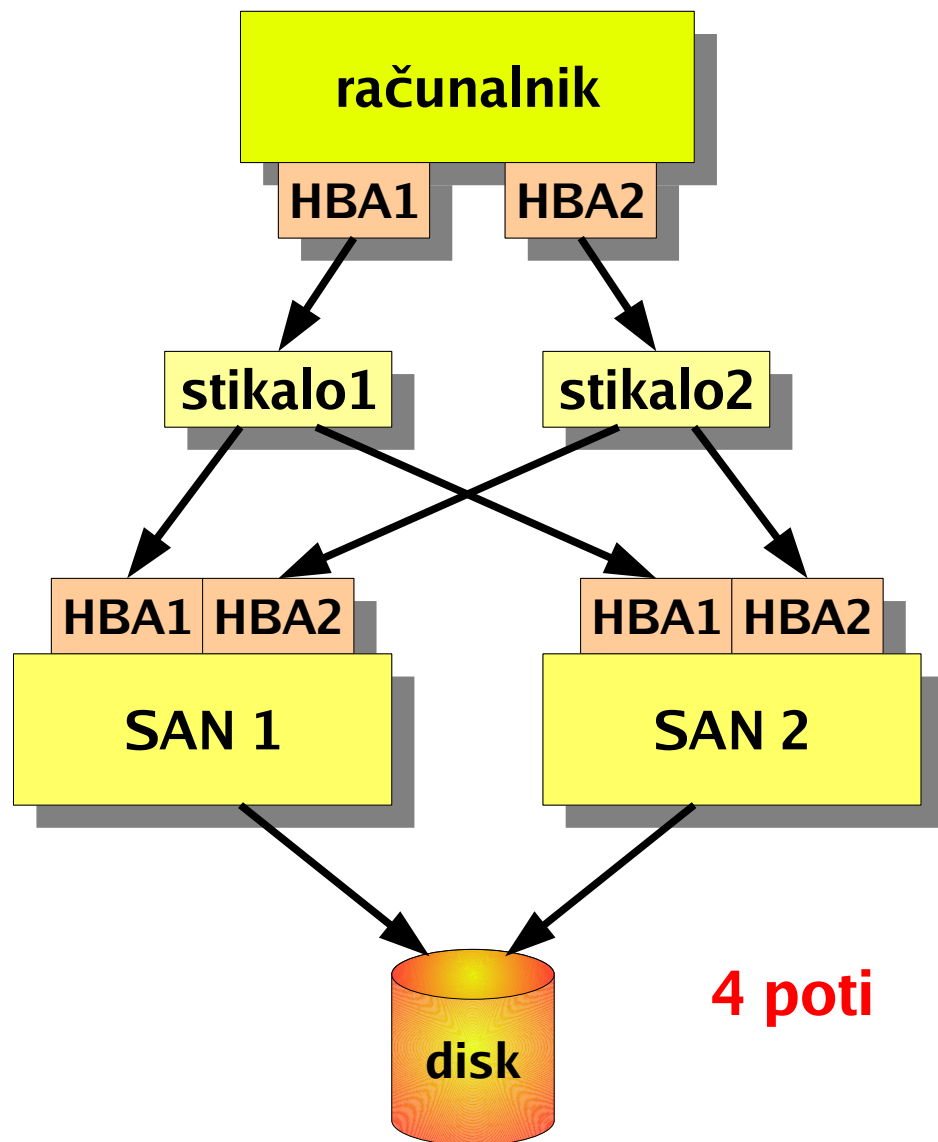
iSCSI klient

- Open-iSCSI (<http://www.open-iscsi.org/>)
- Linux-iSCSI(sfnet) (<http://linux-iscsi.sourceforge.net/>)
(starejši – performančne težave pri branju)
- drugi operacijski sistemi (windows, UNIX, ...)
- VMware ESX Server



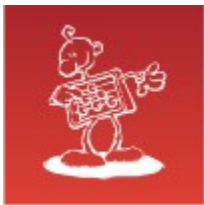


Device mapper (multipath)



Uporabna funkcionalnost:

- »MULTIPATH«
Več vzporednih (redundantnih) poti do istega diska.



UDEV

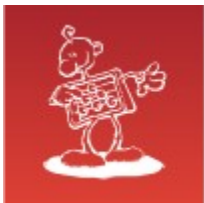
- praktično neomejeno število naprav na /dev (dinamično)
- obstojno poimenovanje naprav

tradicionalno

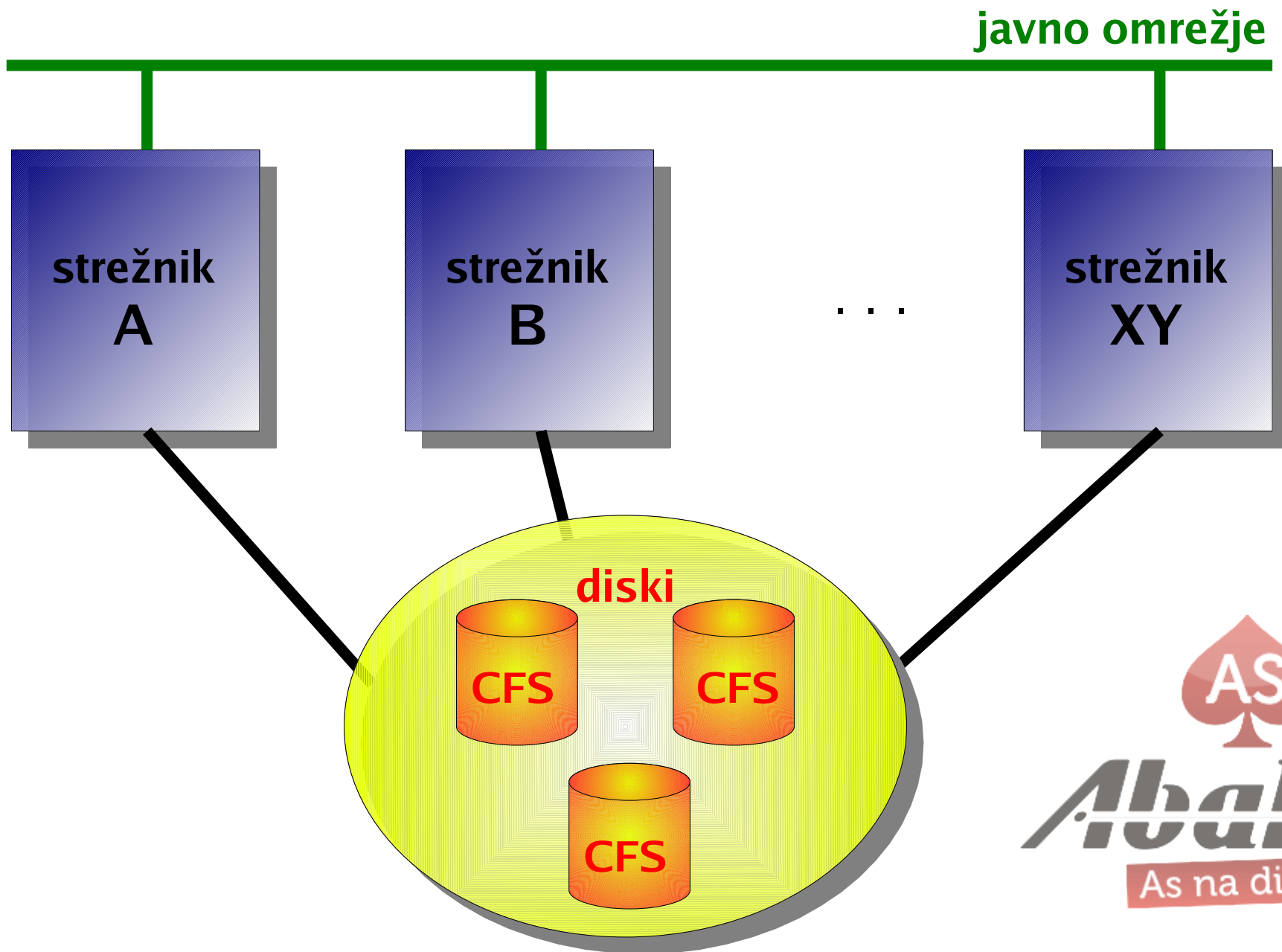
```
Ukazna lupina - Konzola
Seja Uredi Videz Zaznamki Nastavitve Pomoč
Ukazna lupina Ukazna lupina št. 5 Ukazna lupina št. 6 Uk < >
sdb
sdc
sdd
sde
sdf
sdg
sdh
sdi
sdj
sdk
sdl
sdm
sdn
sdo
sdp
sdq
sdr
sds
sdt
sdu
sdv
sdw
sdx
```

udev

```
Ukazna lupina - Konzola
Seja Uredi Videz Zaznamki Nastavitve Pomoč
Ukazna lupina Ukazna lupina št. 5 Ukazna lupina št. 6 Uk < >
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:0@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:1@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:10@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:11@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:2@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:3@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:4@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:5@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:6@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:7@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:8@
ip-10.10.1.10:3260-iscsi-iqn.2007-07.si.delo-prodaja:disk.ssan:9@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:0@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:1@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:10@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:11@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:12@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:13@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:2@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:3@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:4@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:5@
ip-10.10.1.11:3260-iscsi-iqn.2007-06.si.delo-prodaja:disk.tsan:6@
```

Gručni datotečni sistem (CFS)





Gručni datotečni sistem (CFS)

OCFS2 (Oracle Cluster File System)

- standardni del GNU/linux jedra od verzije 2.6.16 (20. 3. 2006) – ekperimentalno in od verzije 2.6.19 (29. 11. 2006) – produkcijsko

GFS (Global File System – RedHat)

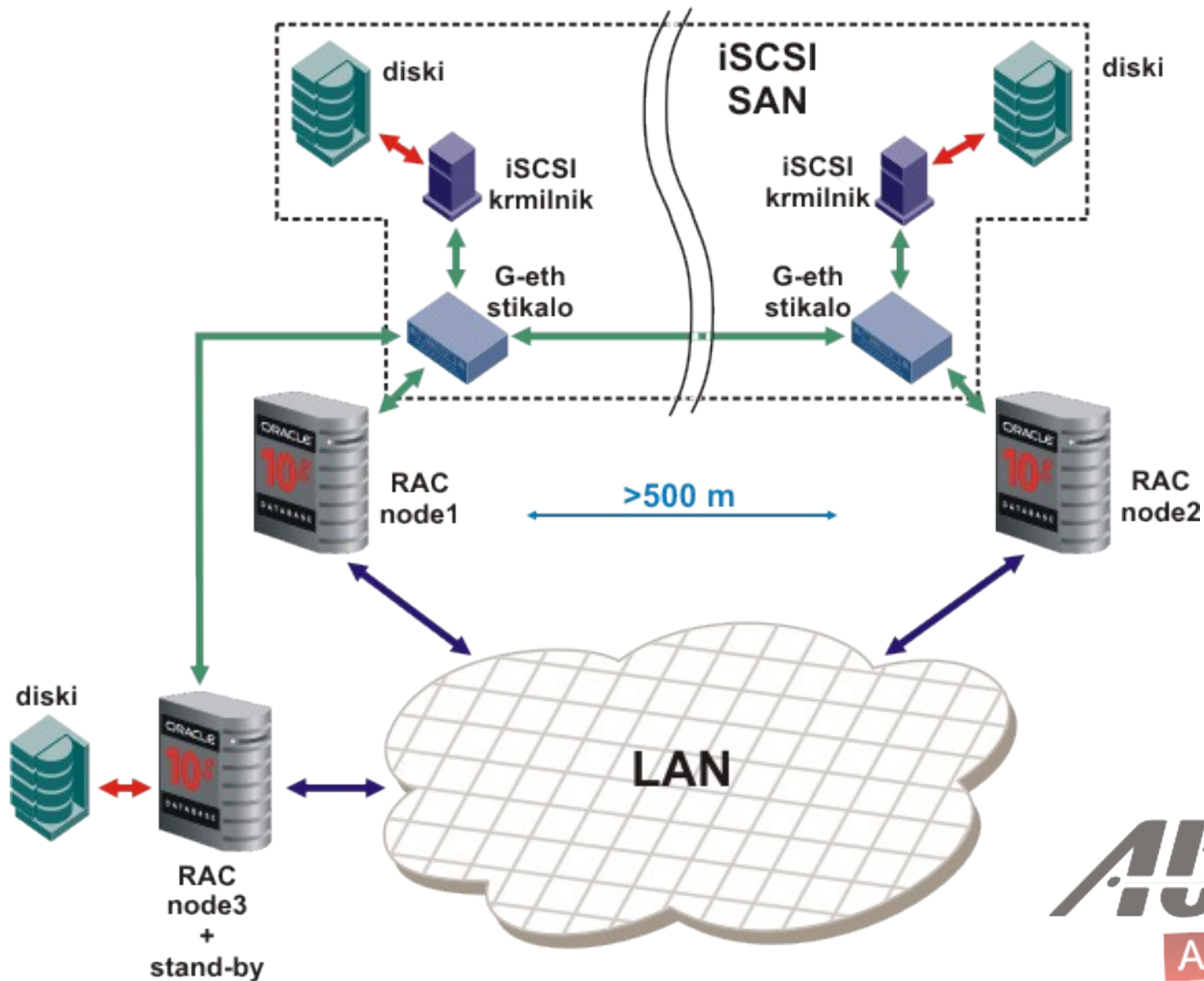
- standardni del GNU/linux jedra od verzije 2.6.19 (29. 11. 2006)

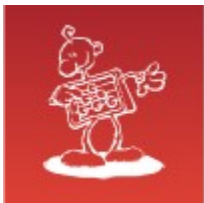
zasnovano na DEC tehnologiji





Uporaba: distribuirana strežniška gruča

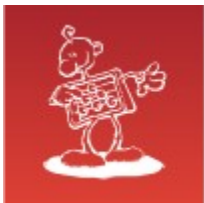




Za konec

- učinkovito hranimo podatke:
 - zanesljivo
 - zmogljivo
 - poceni
- odprtokodne rešitve so preizkušene – na tisoče inštalacij po celem svetu
(DRBD: >60.000 inštalacij, >20 TB diskovnega prostora)





Kako učinkovito hraniti podatke

Vprašanja

Sergej Rožman

ABAKUS plus d.o.o.

Ljubljanska c. 24a

Kranj

e-pošta: sergej.rozman@abakus.si

tel. št.: 04 287 11 14

