

File system Linux di ultima generazione

Linux day 2012, Pisa

File system Linux di ultima generazione

Alla scoperta di nuove funzionalità dei file system, rispolverando alcune vecchie perle nascoste.

Motivazioni:

- curiosità, voglia di imparare, di "smontare"
- gusto di "sporcarsi le mani" con concetti e operazioni a basso livello

I file system sono uno dei fondamenti dei sistemi Linux/Unix

- ogni piccola scoperta/invenzione può innescare una cascata di novità nei livelli superiori

File system Linux di ultima generazione

Un file system è una sequenza di bit su disco.

Non "fa" niente, contiene solo file e cartelle.

File system Linux di ultima generazione

Un file system è una sequenza di bit su disco.

Non "fa" niente, contiene solo file e cartelle.

Niente di più noioso, giusto?

File system Linux di ultima generazione

Un file system è una sequenza di bit su disco.

Non "fa" niente, contiene solo file e cartelle.

Niente di più noioso, giusto?

SBAGLIATO

File system Linux di ultima generazione

Un file system è una sequenza di bit su disco.

Non "fa" niente, contiene solo file e cartelle.

Niente di più noioso, giusto?

SBAGLIATO

Mostrami il tuo codice e nascondi le tue strutture dati, e continuerò ad essere disorientato.

Mostrami le tue strutture dati, e normalmente non avrò bisogno del tuo codice: sarà ovvio.

Eric Raymond, The Cathedral and the Bazaar [1997]

File system Linux di ultima generazione

Un file system è composto da strutture dati su disco, strutture dati in memoria e codice.

Normalmente si considerano solo le strutture dati su disco: il resto è dato per scontato.

File system Linux di ultima generazione

Un file system è composto da strutture dati su disco, strutture dati in memoria e codice.

Normalmente si considerano solo le strutture dati su disco: il resto è dato per scontato.

File e cartelle sono un'*invenzione*

File system Linux di ultima generazione

Un file system è composto da strutture dati su disco, strutture dati in memoria e codice.

Normalmente si considerano solo le strutture dati su disco: il resto è dato per scontato.

File e cartelle sono un'***invenzione***, così come

- permessi
- hard link e soft link
- named pipes, socket unix
- file speciali: block devices, char devices

File system Linux di ultima generazione

Altre *invenzioni*?

File system Linux di ultima generazione

Altre *invenzioni*?

- file sparsi
- loop devices
- journaling
- file preallocati
- snapshots (istantanee) e transazioni
- extents al posto dei blocchi
- device mapper, LVM, crittografia
- ottimizzazioni per SSD e per NAND

File system Linux di ultima generazione

Altre *invenzioni*?

- file sparsi antico, >20 anni
- loop devices vecchio
- journaling vecchio
- file preallocati
- snapshots (istantanee) e transazioni
- extents al posto dei blocchi
- device mapper, LVM, crittografia
- ottimizzazioni per SSD e per NAND

Scoperte

Scoperte - cosa si può fare con questi strumenti?

Scoperte

Scoperte - cosa si può fare con questi strumenti?

- riformattazione non distruttiva:
conversione di un file system da un tipo ad un altro preservando **tutti** i dati
 - anche se è pieno fino al 95%
 - senza bisogno di backup – ma se qualcosa va storto sono dolori, meglio avere un backup

Scoperte

Scoperte - cosa si può fare con questi strumenti?

- riformattazione non distruttiva:
conversione di un file system da un tipo ad un altro preservando **tutti** i dati
 - anche se è pieno fino al 95%
 - senza bisogno di backup – ma se qualcosa va storto sono dolori, meglio avere un backup

fstransform [disponibile su Debian wheezy, Ubuntu 12.10 o direttamente da SourceForge]

Lo ammetto, voglio farmi pubblicità...

File sparsi

Esistono da decenni.

In un file sparso alcuni blocchi NON sono allocati

- per convenzione, i blocchi non allocati contengono zeri
- permettono di risparmiare spazio
- permettono di creare rapidamente file enormi

Uso tipico

- immagini su file di CD, DVD o intere partizioni
 - assieme ai loop device

Si creano con ***truncate*** o ***dd***

Loop devices

Un loop device è un dispositivo a blocchi virtuale

- permette di usare un file (anche sparso) come se fosse un dispositivo a blocchi.

Il comando per gestire i loop device è **losetup**, e una volta che si ha un loop device si può:

- formattarlo con un qualunque file system
- montarlo e usarlo come una partizione
- eccetera...

Comodo per usare le immagini di CD e DVD senza dover masterizzare niente.

Preallocazione file

2001: posix_fallocate() standard POSIX

2010: fallocate() compare nel kernel 2.6.36
ad oggi solo per: btrfs, ext4, gfs, xfs

- chiede al file system di allocare i blocchi per un certo intervallo di un file
- il modo più rapido di creare un file NON sparso: NON scrive nei blocchi
- garantisce che scrivendo successivamente in quell'intervallo non si avranno errori ENOSPC

Conversione tra file system

2005: convertfs, abbandonato nel 2006

2008: btrfs-convert, solo ext3 → btrfs

2011: fstransform (lo ammetto, voglio farmi pubblicità)

- Trasformano un file system in un altro, mantenendo i dati e senza bisogno di backup
- RISCHIOSI: sono beta, se qualcosa va storto fareste bene ad avere un backup
- Possono essere di due tipi:
 - specifici per un file system, es. btrfs-convert
 - generali, es. convertfs e fstransform

Conversione tra file system

Trasformare un file system in un altro, senza backup e preservando tutto il contenuto.

Come funziona?

- creazione di un file sparso grande quanto il device che lo contiene
- formattazione con il nuovo file system
- montato tramite loop device
- spostamento di tutti i file e directory uno alla volta
- spostamento dei blocchi del file sparso fino a coincidere con il device originale (rimappatura)

Detto così sembra quasi facile... magari!

Conversione tra file system

Lo spostamento di tutti i file e directory uno alla volta è lento e rischioso, se il device originale si riempie sono guai: **si PERDONO i dati!**

Si può fare di meglio?

Anziché spostarli dentro il loop device, si possono preallocare! Attenzione ai file system con tail merging

La rimappatura diventa più complessa, ma possibile

Nella prossima versione (0.9.4)

Supporto per SSD

Anche se si usano i normali file system, SSD e memorie USB funzionano meglio – e più a lungo – con alcuni accorgimenti

- disabilitare il journaling
- non usarle per /tmp, /var/tmp, /var/log e simili
- usare l'opzione di mount "noatime" o "relatime" di ext2/3/4
- lasciare un po' di spazio non partizionato
- per gli SSD, usare ***fstrim*** - dice all'SSD di deallocare i blocchi corrispondenti allo spazio libero del file system

Supporto per SSD

Alcuni file system sono ottimizzati per le memorie FLASH, es. SSD e chiavette USB

- NILFS2 – usa un device come unico journal circolare
 - scrive quasi sempre sequenzialmente
 - snapshots e checkpoints
 - garbage collecting per deallocare i blocchi obsoleti
- Btrfs – B-trees + copy-on-write
 - troppo da dire, servirebbe un talk dedicato...

Altri supportano solo le memorie NAND raw

- JFFS2, UBIFS, YAFFS

Altre funzionalità

Extent: raggruppamento di blocchi consecutivi

- risparmia spazio per le bitmap di allocazione dei blocchi
- più efficiente nella gestione dei file contigui

Journaling

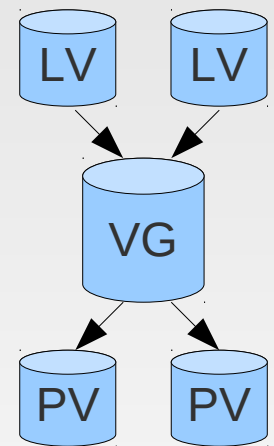
- stranoto, non ha bisogno di parlarne ancora

Device Mapper

- innumerevoli possibilità, es. cifratura

LVM, basato su Device Mapper


- unire più device (PV) per crearne uno più grande (VG)
- suddividerlo in device (LV) anche a caldo
- aggiungere o togliere device (PV) a caldo
- snapshot e transazioni (commit/rollback)



Cosa possiamo farci?

Sono mattoni di base, i soli limiti sono la nostra capacità e fantasia



Burj Khalifa  CC BY-SA 2.0 by Nicolas Lannuzel

File system Linux di ultima generazione

Domande?

File sparsi

Domande per casa

1. Si può creare un file sparso più grande dello spazio libero disponibile?
2. E più grande della partizione che lo contiene?
3. Cosa succede se, scrivendo su un file sparso, si esaurisce lo spazio nel file system che lo contiene?

File sparsi

Domande per casa

1. Si può creare un file sparso più grande dello spazio libero disponibile?

SI

2. E più grande della partizione che lo contiene?

SI

3. Cosa succede se, scrivendo su un file sparso, si esaurisce lo spazio nel file system che lo contiene?

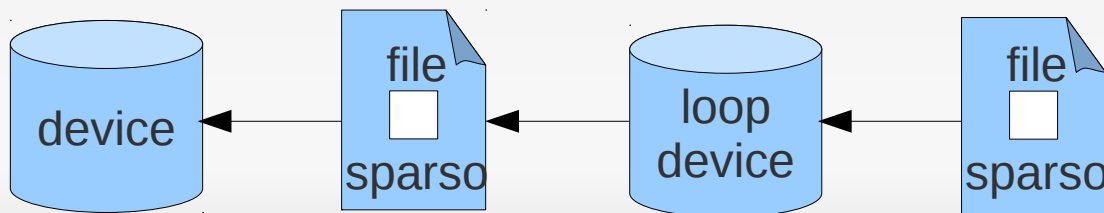
La scrittura dà errore ENOSPC

"No space left on device"

Loop devices

Domande per casa

1. si può creare un file sparso dentro un loop device che usa a sua volta un file sparso?
2. se sì, quante volte si può nidificare questo meccanismo?



Loop devices

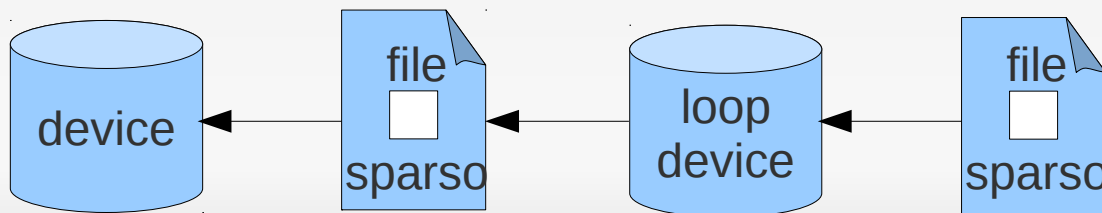
Domande per casa

1. si può creare un file sparso dentro un loop device che usa a sua volta un file sparso?

SI

2. se sì, quante volte si può nidificare questo meccanismo?

**L'unico limite è il numero di loop device
(configurabile al caricamento del modulo)**

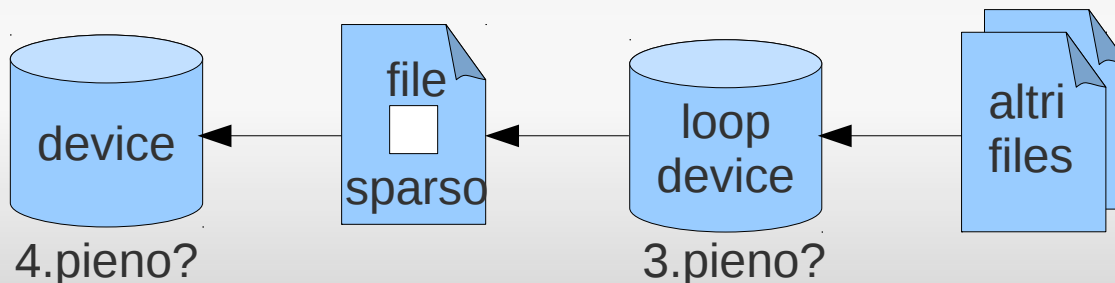


Loop devices

Supponiamo di creare un file sparso, usarlo come loop device, formattarlo e montarlo.

Cosa succede se, scrivendoci dentro:

- 3. si esaurisce lo spazio nel loop device?
- 4. si esaurisce lo spazio nel device che contiene il file sparso?



Loop devices

Supponiamo di creare un file sparso, usarlo come loop device, formattarlo e montarlo.

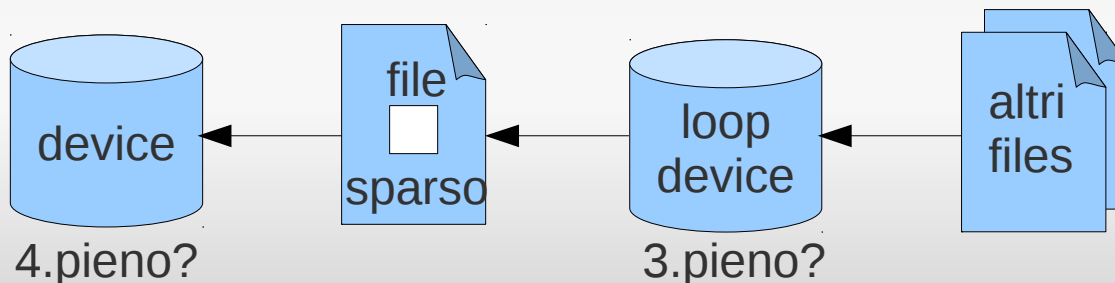
Cosa succede se, scrivendoci dentro:

3. si esaurisce lo spazio nel loop device?

La scrittura dà errore ENOSPC

4. si esaurisce lo spazio nel device che contiene il file sparso?

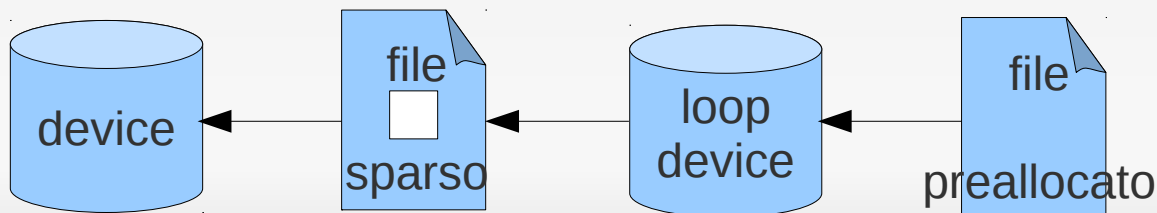
Si PERDONO i dati!! Nessun errore segnalato, solo i log del kernel (dmesg) mostrano il guaio



Preallocazione file

Domande per casa

1. Se si crea un file preallocato dentro un loop device, a sua volta dentro un file sparso, occupa spazio o no?



Preallocazione file

Domande per casa

1. Se si crea un file preallocato dentro un loop device, a sua volta dentro un file sparso, occupa spazio o no?

Occupava spazio dentro il loop device

NON fa aumentare significativamente lo spazio occupato dal file sparso

