

# IBM Academic Initiative



## Sistemi Storage ed Architettura di IO in ambiente zSeries

Novembre 2007

**Stefano Ricci**  
**IBM Italia**  
**Systems and Technology Group**



# Agenda

## Prima lezione

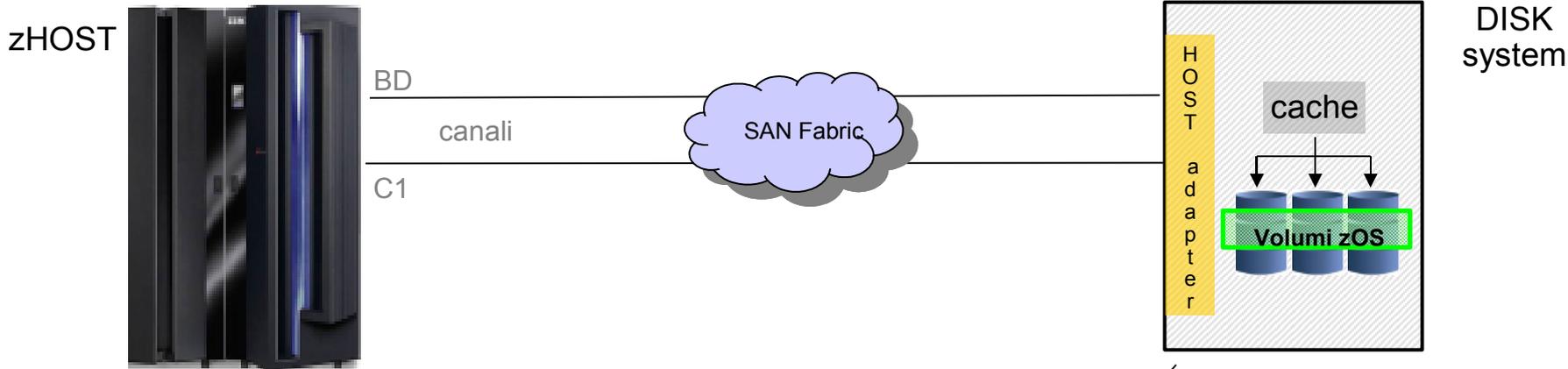
- L'architettura di I/O in ambiente zSeries
- Sistemi Storage a disco

## Seconda lezione

- Sistemi Storage a nastro
- Le funzioni di copia dei dati
- Connettività e trasporto dei dati
- Cenni sulle prestazioni dei sistemi a disco

- L'architettura di I/O in ambiente zSeries
- *Sistemi Storage a disco*
- *Sistemi Storage a nastro*
- *Le funzioni di copia dei dati*
- *Connettività e trasporto dei dati*
- *Cenni sulle prestazioni dei sistemi a disco*

# Indirizzamento di un volume logico in ambiente zSeries



## Generazione dello IO allo zOS (IOCP)

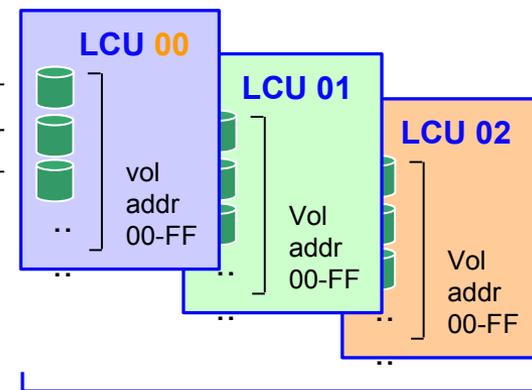
Control Unit =00, Path =( BD,C1), UNIT=DISK

IO device=(0000,00FF),Control Unit= 00,UNIT=3390 B



Indirizzo fisico volumi

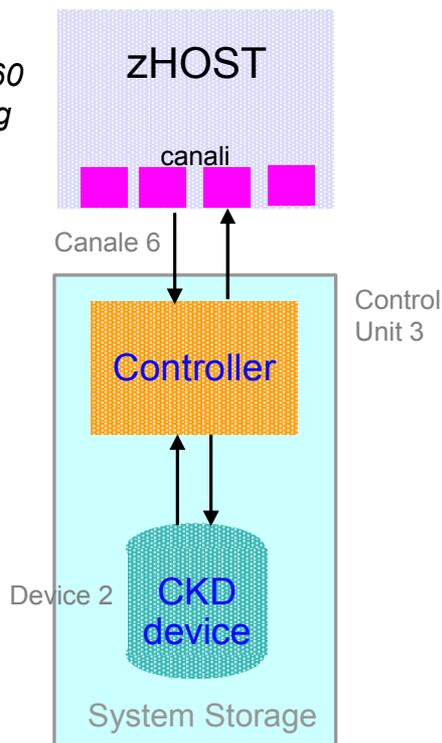
- 00 00 ←
- 00 01 ←
- 00 02 ←
- ..
- ..



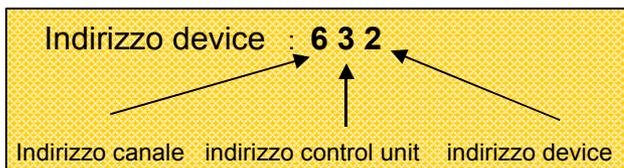
LCU = 00-FF

# Fondamenti dell'architettura di IO in ambiente zSeries

System/360 addressing

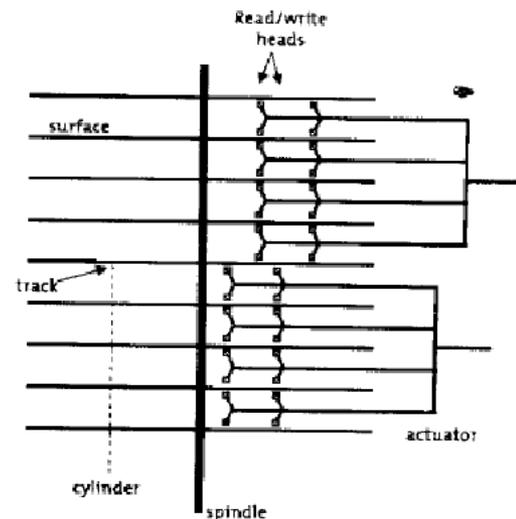


| Architettura                 | Central storage addressing | IO device addressing (per LPAR) | Canali (per LPAR) | Funzioni introdotte  |
|------------------------------|----------------------------|---------------------------------|-------------------|--|
| <i>System 360 (1960)</i>     | 24 bit                     | 16 bit (4k)                     | 256               | <ul style="list-style-type: none"> <li>Concetto di LCU (piu' paths al device)</li> </ul>   |
| <i>System 370</i>            | 24 bit                     | 16 bit (65k)                    | 256               | <ul style="list-style-type: none"> <li>IO multiple sullo stesso canale</li> <li>Selezione del path e riconnessione dinamica (DPS e DPR)</li> </ul> |
| <i>System 370 XA ed ESA</i>  | 31 bit                     | 16 bit (65k)                    | 256               | <ul style="list-style-type: none"> <li>Canali condivisi tra i processori</li> <li>System Managed Storage</li> </ul>                                |
| <i>System 390</i>            | 31 bit                     | 16 bit (65k)                    | 256               | <ul style="list-style-type: none"> <li>Nuova tecnologia di canale (ESCON e FICON)</li> </ul>   |
| <i>z/Architecture (2000)</i> | 64 bit                     | 2*16 bit (65k+65k)              | 256               | <ul style="list-style-type: none"> <li>FICON piu' veloci (Express 2 e 4)</li> <li>Uso del secondo set di indirizzi per gli alias</li> </ul>        |

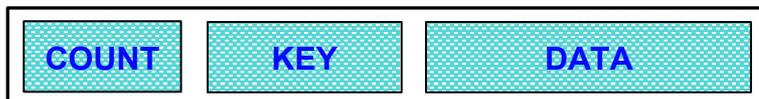


## Architettura CKD - DASD

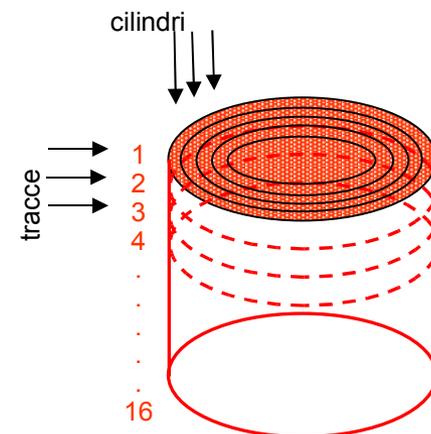
- I dischi CKD sono nati per il mercato mainframe (FBA e' l'architettura nel mondo Open)
- Anche indicati con il nome DASD (Direct Access) consistono in una serie di superfici, disposte in parallelo a formare un **cilindro**, per la memorizzazione dei dati. Dei bracci (actuator) muniti di testine di lettura/scrittura si muovono simultaneamente sulle superfici dove i dati vengono memorizzati in **tracce** concentriche.
- Ogni movimento del braccio munito di testine permette di posizionarsi contemporaneamente su 16 tracce e di trasferire senza altri movimenti meccanici un intero cilindro (costituito dalle tracce che sulle 16 superfici risultano nella stessa posizione).
- 16 tracce formano 1 cilindro; una traccia e' di servizio e le altre 15 sono usate per memorizzare i dati; ogni traccia contiene un numero fisso di Byte (56.664) organizzati nel **tracciato record CKD**. Un record utente e' indirizzato indicando cilindro, traccia e numero record.
- Ogni modello di disco (3390-3, 3390-9, ....) ha un numero diverso di cilindri, quindi fornisce una diversa capacita' di memorizzazione



Il **tracciato record CKD** venne introdotto dal System 360 per avere una corrispondenza univoca tra i comandi eseguiti dai programmi ed il formato fisico dei dati sul disco



COUNT :      Identifica il record number (CCHHR)  
 KEY :        Campo opzionale per la ricerca  
 DATA :      Area dati; lunghezza variabile per ciascun record



## Architettura CKD – formato traccia

Ogni traccia disco e' formattata con una disposizione precisa di record:

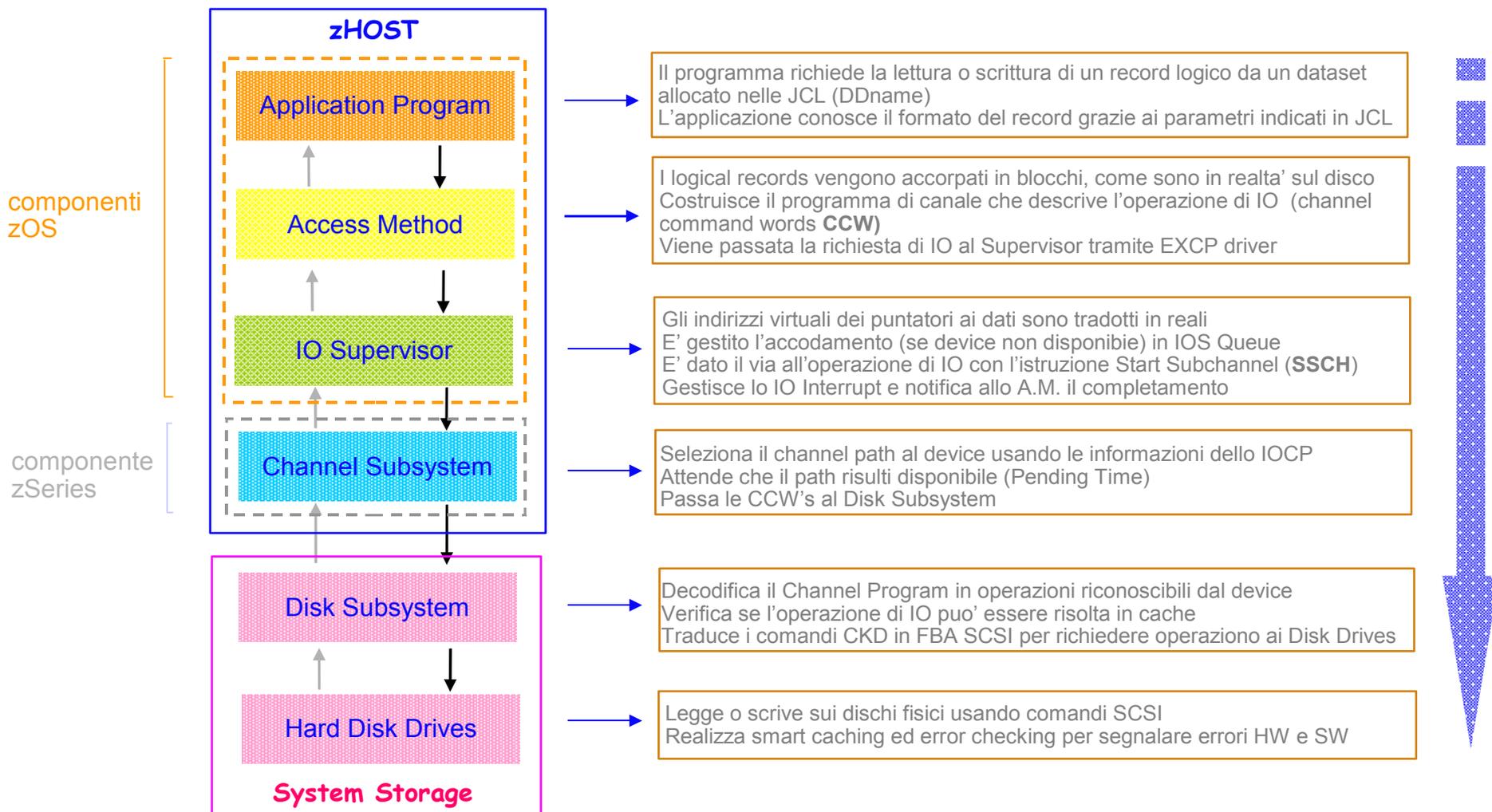
- 'index' e' l'inizio della traccia
- 'home address' e' l'indirizzo fisico della traccia
- 'R0' e' il record zero che descrive la traccia
- Seguono i record di dati in formato *Count, Key, Data* oppure *Count, Data*



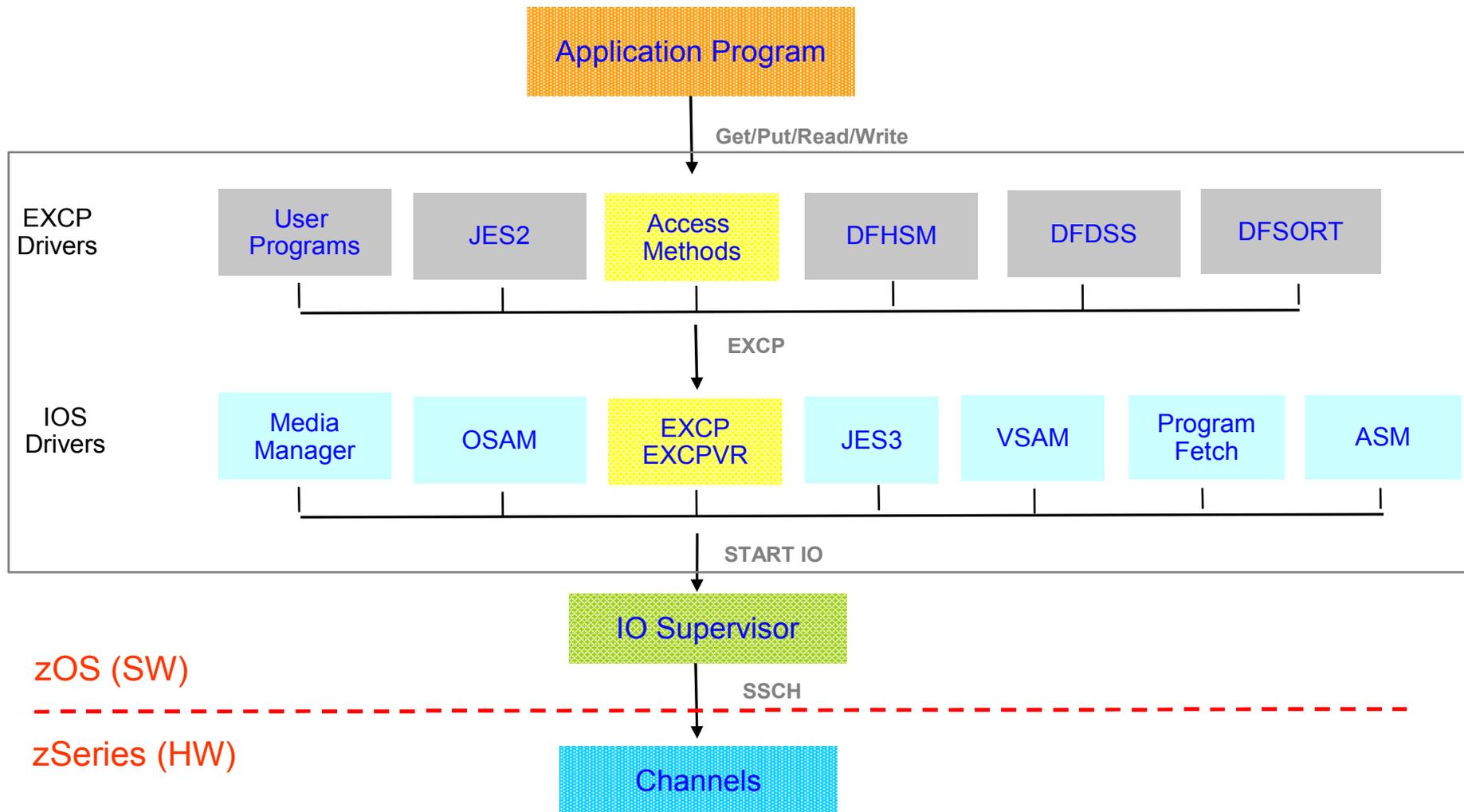
Es: RECFM=FB, LRECL=80, BLKSIZE= 4560  
 ogni Data Area conterra' 4560 bytes, divisi in record logici da 80 bytes  
 ogni traccia (device type 3390) conterra' 11 record fisici per 50.160 bytes

| indirizzamento del disco fisico  |         | Modello | Cylinders | Tracks   | Bytes/volume | Bytes/track   |        |      |       |        |       |
|--|---------|---------|-----------|----------|--------------|---|--------|------|-------|--------|-------|
| <div style="border: 1px solid black; display: inline-block; padding: 2px;"> <table border="1" style="border-collapse: collapse; text-align: center;"> <tr> <td style="width: 20px; height: 20px;">C</td> <td style="width: 20px; height: 20px;">C</td> <td style="width: 20px; height: 20px;">H</td> <td style="width: 20px; height: 20px;">H</td> <td style="width: 20px; height: 20px;">R</td> </tr> </table> </div> | C       | C       | H         | H        | R            | <b>CC</b> – con 16 bit unsigned indirizza max un device type 3390 Mod. 54 | 3390-1 | 1113 | 16695 | 946 MB | 56664 |
|  | C       | C       | H         | H        | R            |   |        |      |       |        |       |
|  |         | 3390-2  | 2226      | 33390    | 1.89 GB      | 56664   |        |      |       |        |       |
|  |         | 3390-3  | 3339      | 50085    | 2.83 GB      | 56664   |        |      |       |        |       |
|  |         | 3390-9  | 10017     | 150255   | 8.51 GB      | 56664   |        |      |       |        |       |
|  | 3390-27 | 32760   | 491400    | 27.84 GB | 56664        |   |        |      |       |        |       |
| <b>HH</b> – 16 testine lettura/scrittura   | 3390-54 | 65520   | 982800    | 55.68 GB | 56664        |   |        |      |       |        |       |
| <b>R</b> – numero del record nella traccia   |         |         |           |          |              |   |        |      |       |        |       |

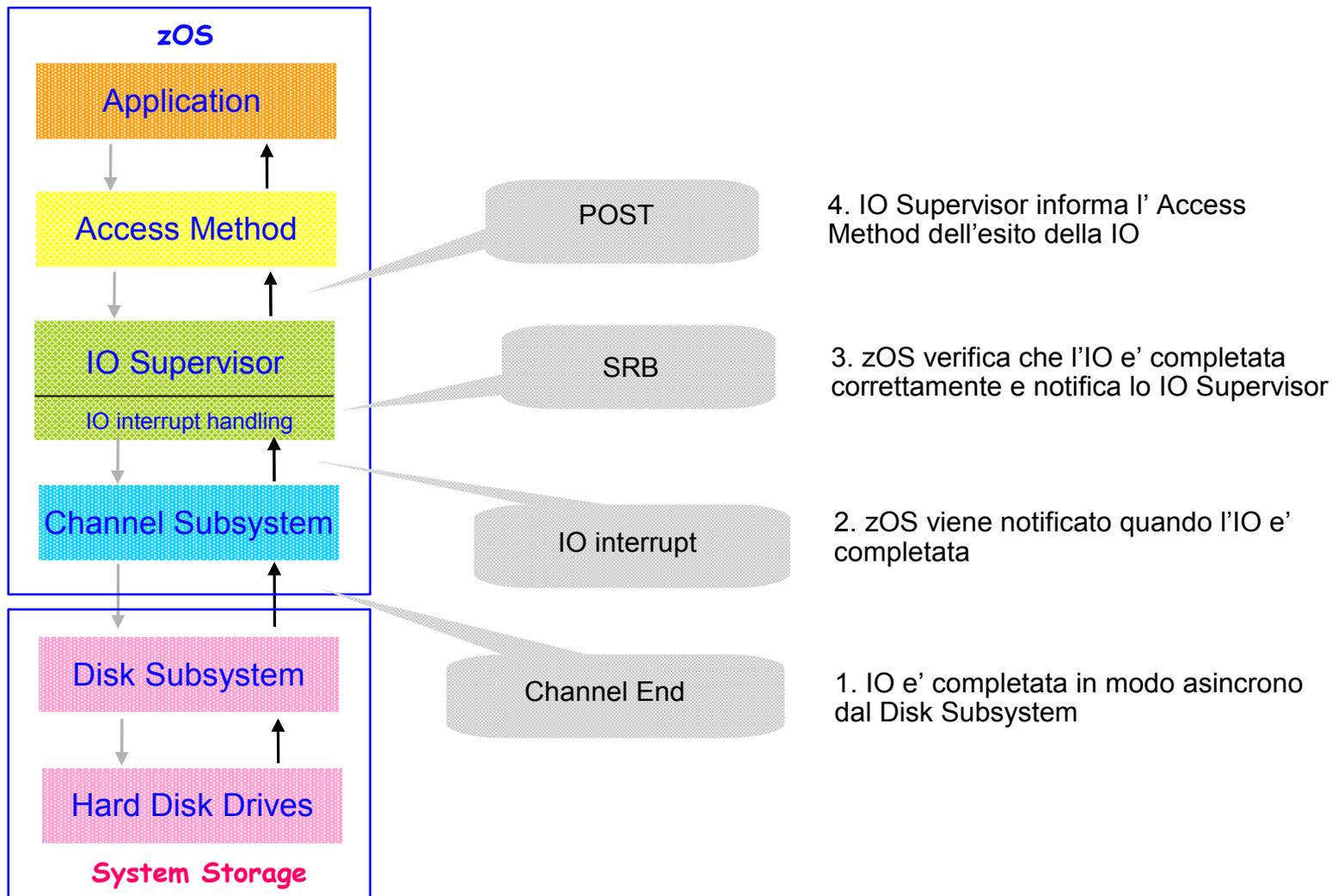
# Flusso logico della richiesta di un dato



## Funzioni zOS per l'accesso ai dati



# IO interruption – l'applicazione ottiene il dato richiesto

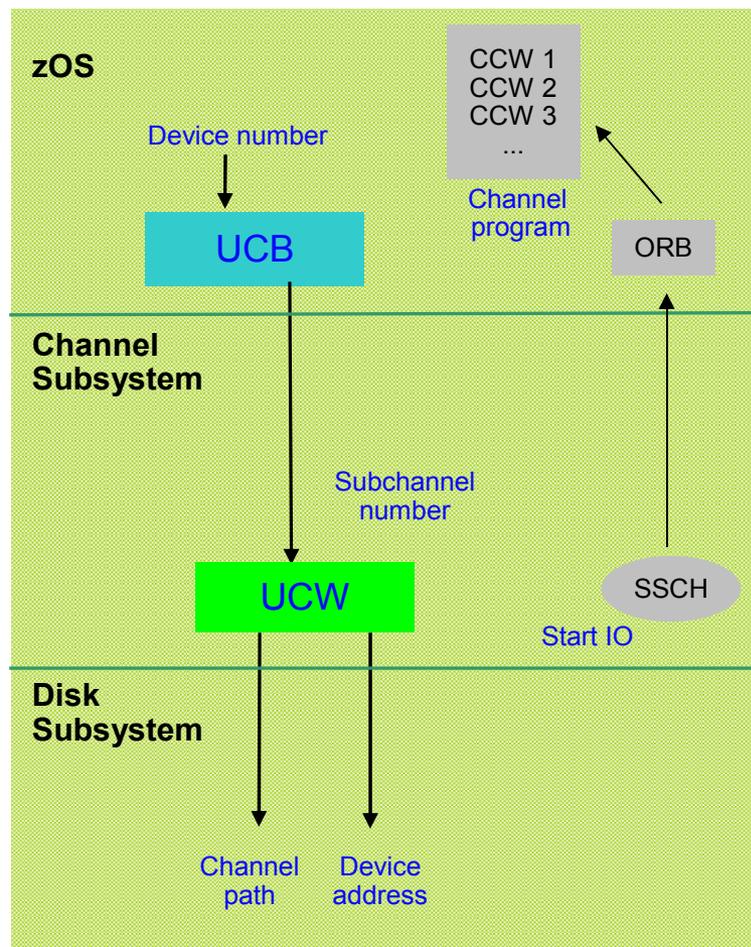


# Start dell'operazione di IO

Lo zOS usa il **device number** per identificare il volume logico definito nella generazione di IO e memorizza le informazioni del device nello **Unit Control Block (UCB)**

Il Channel Subsystem usa i **subchannel number** per identificare i device e le **Unit Control Word UCW's** contengono le informazioni sui device, compresi i path disponibili (equivalenti alle LCU)

Attraverso il canale prescelto viene acceduto il volume logico richiesto identificato dal device address



## Start IO

Quando avviene lo Start di una operazione di IO (IO Supervisor fa SSCH) il **subchannel number** e' passato allo zOS.

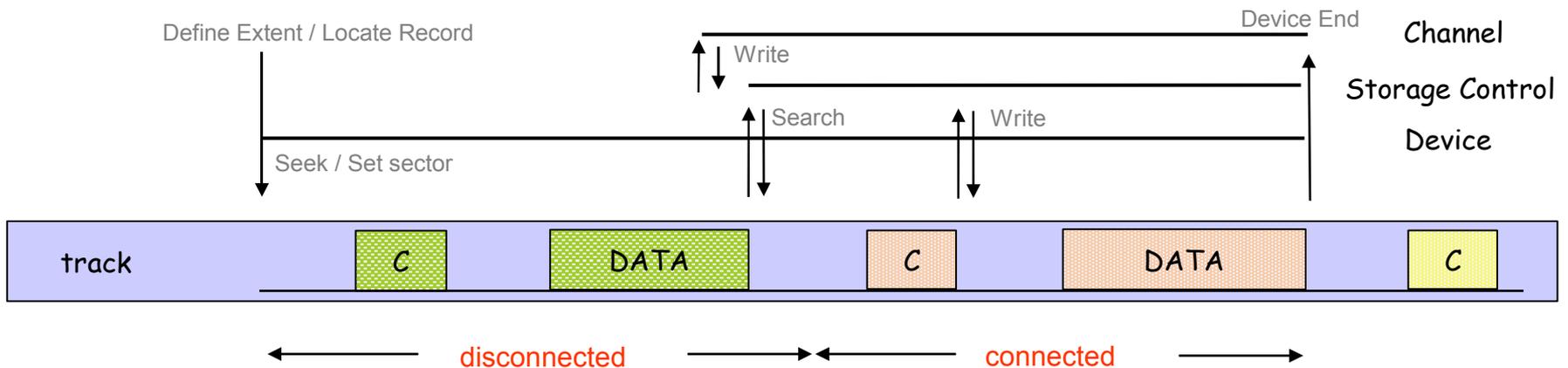
Altre informazioni, come il puntamento al channel program e allo **UCB**, sono contenute nello ORB.

Le **CCW** contengono le istruzioni per gestire l'operazione sul canale e sul device, oltre a quelle per manipolare i dati.

## Esecuzione di un programma di canale non sincrono

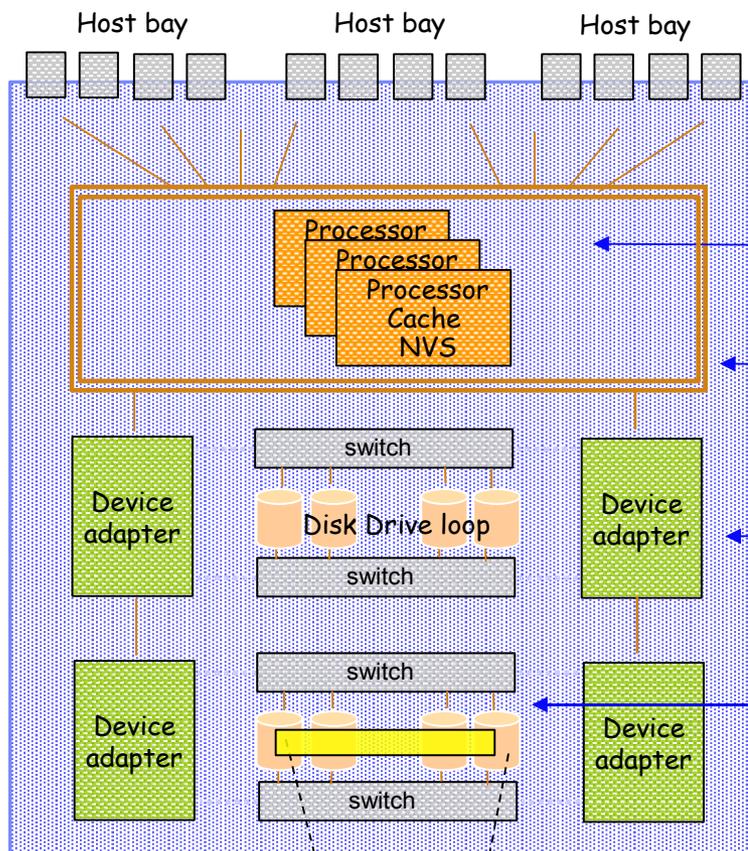
Il tempo necessario al device per posizionarsi sulla traccia richiesta permette la disconnessione e l'acquisizione dei dati dal canale

- I comandi di Seek e Set Sector CCW (ricerca settore) vengono gestiti dal device in modo autonomo
- Prima di ricevere la Search CCW (posizionamento Data Area per scrittura) il device si deve riconnettere alla Storage Control
- Terminata l'operazione la Storage Control richiede la CCW successiva al canale



- *L'architettura di I/O in ambiente zSeries*
- **Sistemi Storage a disco**
- *Sistemi Storage a nastro*
- *Le funzioni di copia dei dati*
- *Connettività' e trasporto dei dati*
- *Cenni sulle prestazioni dei sistemi a disco*

# Disk Subsystem architecture



← **Host Adapter** : processori specializzati (Frontend) per connessione ad host o per funzioni di replica

← **Processori e Cache** : gestione delle operazioni. Per performance e data integrity

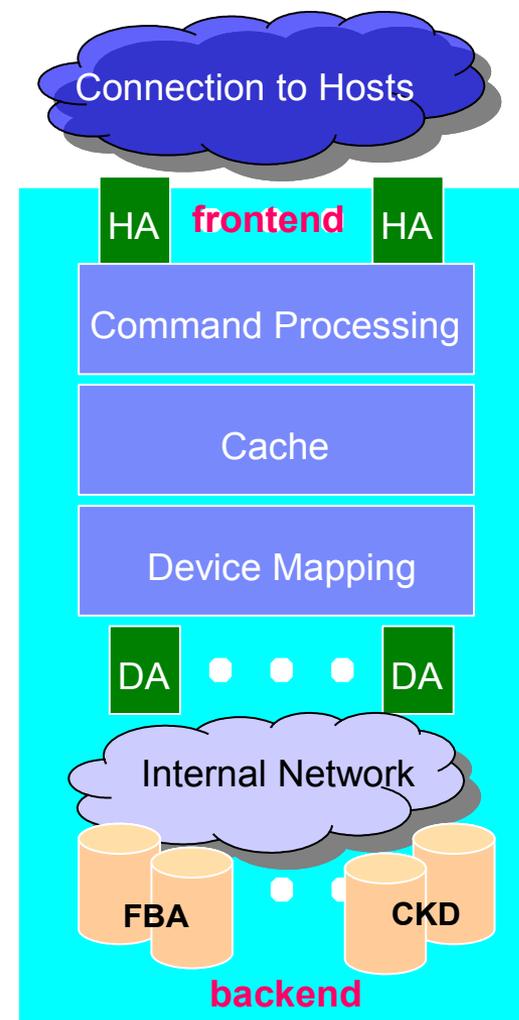
← **Network interna** ad alta velocita' per il colloquio tra Frontend adapters, Processori e Backend drives

← **Device Adapter** : processori specializzati per la connessione ai loop di dischi in backend (DDM) attraverso SSA o switched FC AL.

← **Disk Drives (DDM)** : dischi in architettura SCSI di capacita' e velocita' diversa, organizzati in RAID e connessi ai DA attraverso rete interna ad alta velocita'

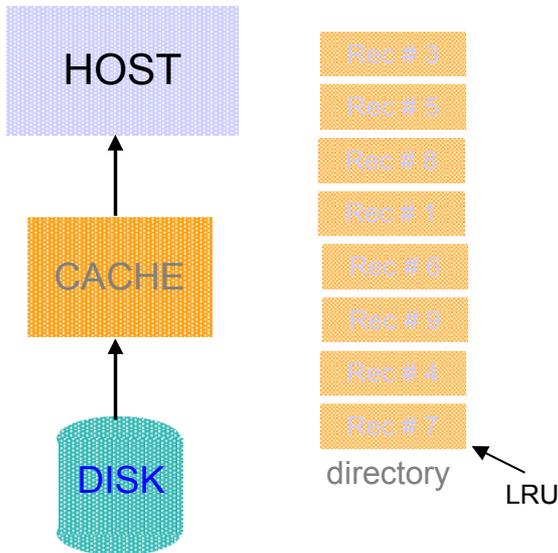
## Disk Subsystem - caratteristiche

- Supporto per ambienti zSeries (CKD) ed Open (FBA)
- Connettività' Fiber Channel/Ficon ed ESCON
- Il formato CKD e' usato per interfacciare i comandi Host ma non e' usato come formato dei dati a disco.
- Virtualizzazione degli SCSI device in volumi logici raggruppati in Logical Subsystem (LCU). Simulazione della geometria 3390 su device SCSI in architettura RAID
- Funzioni microcodificate di copia dei dati (locale e remota) e di ottimizzazione delle prestazioni per l'accesso ai dati (es. PAV)
- Gran parte delle operazioni risolte a velocita' di memoria/network senza coinvolgere i device in backend : tutte le scritture avvengono in cache : 100% write hit
- Operazioni effettuate dai dischi in backend :
  - Stage : lettura dal device per cache miss pre-stage
  - Destage: scrittura a disco di dati modificati dati non rimossi dalla cache
  - Demote : rimozione di dati dalla cache (non modificati)



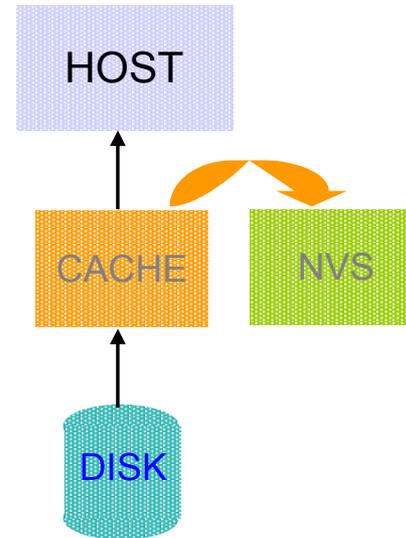
# Caching concepts

## Caching read operation



- La cache e' una memoria elettronica che permette di memorizzare i dati piu' utilizzati ed accederli piu' velocemente
- Mantiene una lista con i blocchi usati piu' di recente cosi' da rimuovere i meno utilizzati
- La porzione di cache volatile contiene solo dati scritti sui dischi in backend

## Caching write : a safe copy

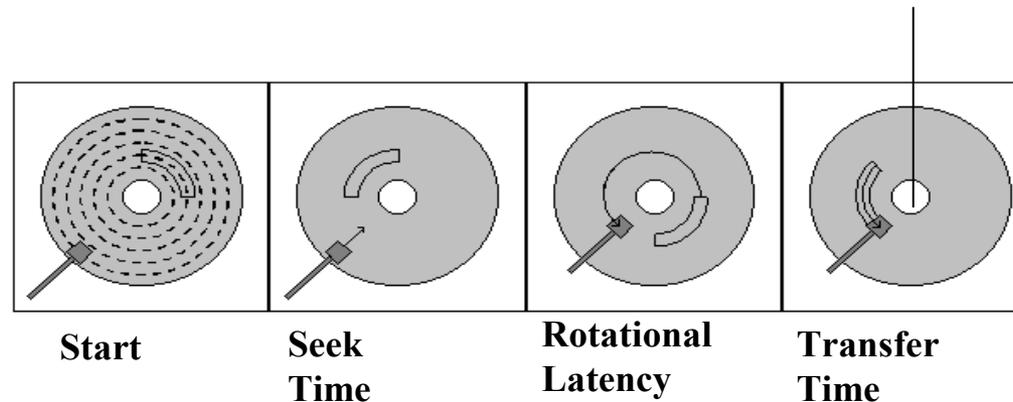


L'integrita' dei dati e' preservata garantendone la loro scrittura permanente:

3. il dato e' scritto in cache
4. Una seconda copia e' scritta nella parte non volatile (NVS) della cache
5. E' segnalata al host l'avvenuta scrittura
6. In maniera asincrona il dato verra' scritto sui dischi in backend

# Disk Drive - fondamentali

- read / write cache hits sono nell'ordine di  $\sim 1\text{ms}$
- physical disk I/O operations sono nell'ordine di  $> 4\text{ms}$  poiche' richiedono l'intervento di componenti meccaniche come il movimento delle testine e la rotazione del disco
- ogni disk drive puo processare un numero ben preciso di I/O al secondo, determinato da:
  - *average seek time* = posizionamento della testina sulla traccia richiesta
  - *rotational latency* = rotazione della superficie del disco finche il primo settore indirizzato passa sotto le testine di read/write heads (avg. time = meta' rotazione)
  - *transfer time* = read/write dei settori di dati (1 sector = 512 Byte)



# Architettura RAID - Concetti

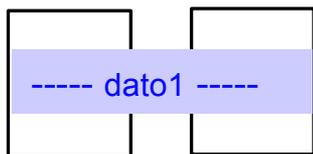
- Utilizzo di dischi economici pur mantenendo una buona affidabilità del sistema complessivo
  - ❖ *mettere insieme più dischi SCSI di size ridotto risulta più conveniente che assemblare un Single Large Expensive mainframe Disk (SLED)*
- Ridondanza delle componenti per permettere di ricoverare i dischi in errore anche a fronte di più eventi contemporanei
  - ❖ *vengono implementate diverse tecniche di RAID per rispondere a requisiti di performance / availability / capacity*
- Tecnica di ridondanza delle informazioni più efficiente del semplice mirroring
  - ❖ *la tecnica del bit di parity risulta più conveniente della duplicazione del dato*
- Fornire un 'data rate' elevato grazie alla distribuzione dei dati su più dischi
  - ❖ *avere più meccaniche che lavorano in parallelo su dati 'striped' garantisce un throughput elevato*
- RAID10, RAID5, RAID6 sono le tecniche maggiormente diffuse

# RAID types

## RAID 0 (striped, no RAID)

Dato distribuito su un set di dischi

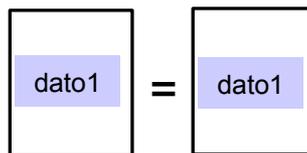
Data rate=elevata  
Availability=bassa



## RAID1 (Mirroring)

Dato scritto su un disco in replica

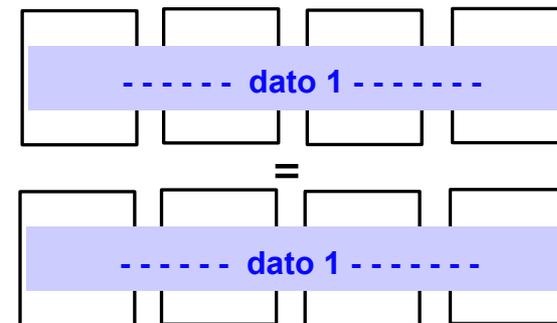
Write overhead=100%  
Read overhead / benefit=0



## RAID10 (striped + mirroring)

Il dato e' distribuito su un set di dischi duplicato

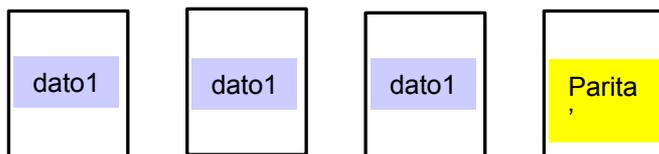
Capacita' utile dimezzata  
Overhead scrittura = 100%  
Availability=eccellente



## RAID 5 (striped + parita')

Data distribuito su un set di dischi

Data rate=elevata in lettura  
Overhead scrittura = da 114% a 400%  
Availability= uso di un bit di parity (alta)  
Capacita' persa = da 14% a 33%



## RAID 6 (striped + doppia parita')

Data distribuito su un set di dischi

Data rate=elevata in lettura  
Overhead scrittura = da 133% a 600%  
Availability= uso di doppio bit di parity (eccellente)  
Capacita' persa = 33%

