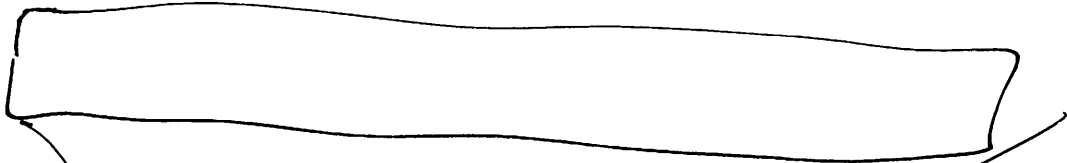# CSE 562

## physical data
### organization

```python
from re import split

with open('data.csv', 'r') as f:
    for line in f:
        fields = split(",", line)
        if(fields[2] != "Ensign" and int(fields[3]) > 25):
            print(fields[1])
```
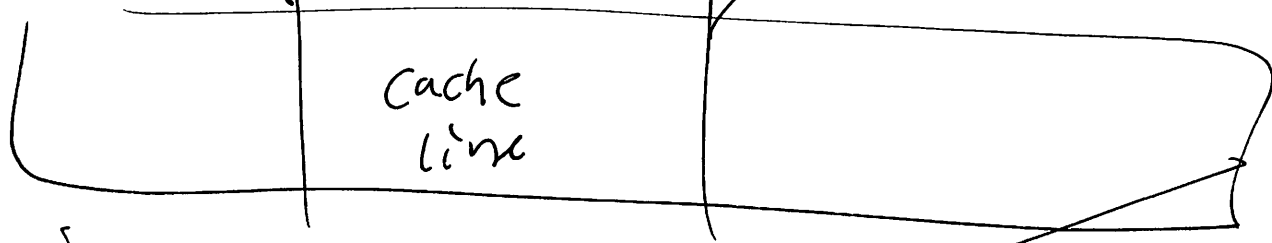
1,Redshirt,Ensign,19,
2,Spock,Lt.,103,
3,Kirk,Capt.,22,
4,Redshirt,Ensign,21,
5,Redshirt,Ensign,18,
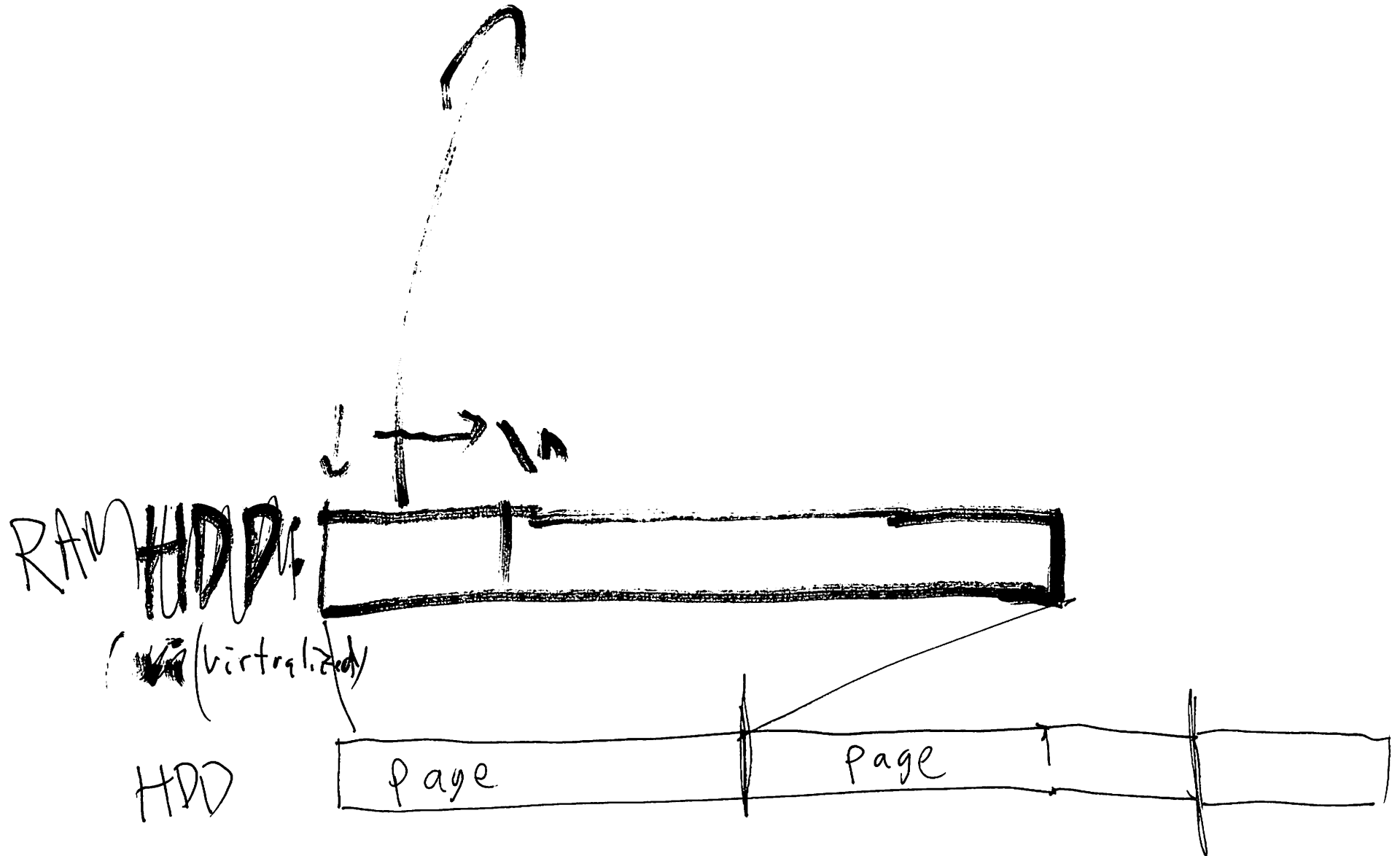6,McCoy,Lt. Cmdr,38,

Cache

low latency
low capacity

RAM
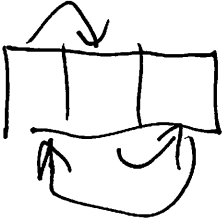
cache
line

Mod latency
Mod capacity

HDD

| Page | Page | Page |

high latency
Large capacity

f : sequence of records
(iterator)

RAM HDD

RAM HDD:

i → n

(virtualized)

HDD | Page | Page | |

Stream
Iterator ☐→☐→☐→◯→☐→☐

Array

Stream
of records [record] [record] [record] · ~~~

Stream
of bytes [byte] [byte] [byte] · — [\n] —

pages

# Idea 1

store fields in a native format

pro: No casting

con: Need to save information about
storage format

```
┌──────┬─────────────┬─────────────┬──────┐
│  4   │      8      │     6       │  4   │
│ int  │ big string  │ big string  │ int  │
└──────┴─────────────┴─────────────┴──────┘
```

↰ Fixed size encoding

Pro: No delimiters (saves 4 bytes)
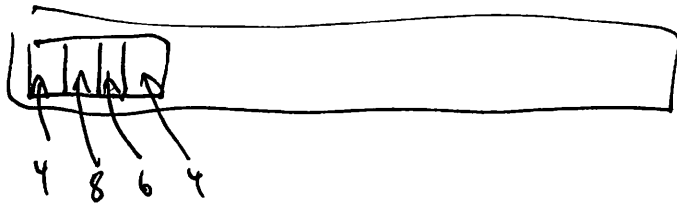
Con: Need to know max length
Wasted space on variable length fields
Not good for adding data

representing records

"CSV" → delimiters

"Fixed Width"

"Per record directory"



Factors

Size of the file

Types of the fields

How much data is changing

output

$\uparrow$

transform

$\uparrow$

filter records

$\uparrow$

read file