

# • C a C++

false = 0  
true ≠ 0

char ≤ short ≤ int ≤ long ≤ long long < float < double

→ implicitní konverze vždy do silnějšího typu ⇒ int · long = long

• array int a[4] ← random hodnoty → adresa prvku = báze + index  
int p[] = {1, 2, 3}

• string = array charů končí na NUL = \0 sizeof("ahoj") = 5

• pointer → je to číslo (adresa), ale musím říct, že to je pointer na něco

int v = 8;

int\* pv = &v; & mi dá adresu

\*pv = 4; \*pointer = přístup na tu adresu, nyní v = 4

• reference → ukazatel, co ale ukazuje jen na 1 fixní adresu → jen v C++

int v = 8;

int &pv = v;

pv = 4;

→ v = 4

→ vnitřní a vnější zarovnání alignment

→ typ velikosti X musí být na adrese A: A%X = 0

0B	c	↓	↓	↓
8B	d			
16B	i		j	
24B	k			↓

→ vnější zarovnání na velikost největšího typu v té struktuře  
⇒ aby se to dolo dalo do pole

• struct data {

char c; 1B

double d; 8B

int i, j, k; 4B

}

void funkce (data in, data \*out) {

out->c = in.c;

out->d = in.d;

}

→ v C je vše předáváno hodnotami

→ předání referencí pomocí pointeru

→ přístup k membrám reference

(\*ptr). field ≡ ptr->field

C++ reference  
data & out  
out.c = in.c

• konstanty

const int jeden = 1; ← neměnná proměnná

constexpr int druhý = 2; ← compile-time konstanta - není v paměti ⇒ lepší

• preprocessing - tohle se řeší na úplném začátku kompilace

#include <stdio.h> → prostě se tam naskopí obsah toho souboru

#define N 100

#if def ...

#endif

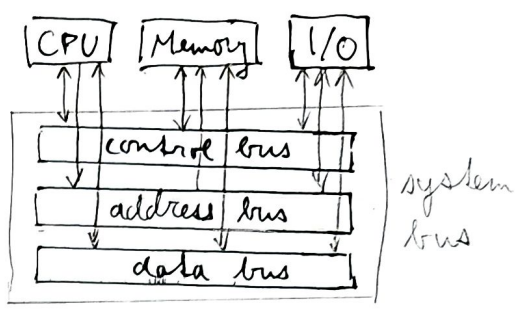
} ngh

# CPU architektura

## Von Neumann

- program a data ve stejné operační paměti
- 1 sdílená sběrnice  $\Rightarrow$  vždy jen 1 transakce
- $\hookrightarrow$  když CPU píše do paměti, tak nemůžeme z paměti psát na monitor

formalý  $\uparrow$



## Harvard

$\Rightarrow$  trochu rychlejší

- oddělená data a code memory  $\rightarrow$  sběrnice z CPU i do code mem
- používá v nejvyšších  $\mu$ C  $\rightarrow$  je potřeba více adresových portů procesoru

## Reálná

- CPU má paměťové řadiče a paměti jsou připojeny vlastní sběrnici přímo k CPU
- $\hookrightarrow$  dokonce má víc kanálů do 1 paměti  $\Rightarrow$  paralelizace
- grafika je buď uvnitř procesoru nebo externí grafická karta
- South Bridge
  - připojený sběrnici k CPU
  - jsou na něj připojeny všechny periférie  $\rightarrow$  připojená přes PCI express
  - $\Rightarrow$  HDD, DVD drive, myš, síťová karta, zvuková karta
- dřív byl i North Bridge pro paměť  $\rightarrow$  pomalý  $\rightarrow$  teď to dělá procesor
- všechny sběrnice jsou peer-to-peer = mezi 2 zařízeními (drive sdílená)

## Co je to architektura?

### Instrukční sada - ISA

$\hookrightarrow$  x86, MIPS, ARM

- specifikace jak se ten procesor má chovat, jaký má instrukce, registry, ...

### HW architektura - implementace nějaké ISA

### Instrukce - specifikuje je ISA

1. load instrukce z adresy v IP
2. decode inst.
3. load operands
4. execute inst.
5. store result
6. increment IP

### Tridy instrukcí

- load: mem.  $\rightarrow$  reg.
- store: reg/imm.  $\rightarrow$  mem. } slow
- move: reg  $\rightarrow$  reg
- aritmetika + logika
- skoky - podmíněné x nepodmíněné
  - přímé = přímo adresa kam
  - nepřímé = adresa kam je na adrese
  - relativní = o kolik se posunout

x86 má HW zásobník  $\leftarrow$

$\rightarrow$  nebo se zásobník dělá přes registry

### • call + return

- při volání funkce je potřeba zásobník, kam si ta funkce uloží svou návratovou adresu  $\rightarrow$  odkud byla zavolána

## • Registry

- obecné - umí aritmetiku, adresaci a složky
- integer aritmetika
- float aritmetika
- adresové - pro nepřímou adresaci ve složce
- branch - pro složky
- příkazové = flags - boolovské hodnoty
- predikátové - 1-bit registry pro každou instrukci - určují, jestli se provede
- aplikacní - speciálně určeny pro nějakou sadu instrukcí
- systémové - nastavení vlastnosti toho procesoru
- vektorové - paralelně dělá 1 instrukci na více dat

if-else nepodporující složky

## • Jména registry

- 99% je pojmenovaných přímo (EAX, r01)
- drív občas relativně více vizku zásobníku

## • Aliasing registry

- překrytí registry → x86: EAX 

Ax	
AH	AL
- snaha se tomu vyhnout → složité překladače

## • x86

- není ortogonální ⇒ specifické instrukce pracují se specifickými registry  
EAX = akumulátor, EBX = base, ECX = count, EDX = data
- ortogonální ISA = libovolná instrukce může použít libovolný reg (výjimky)  
↳ třeba očíslování r0-r31 a jsou ekvivalentní
- segmentové registry
- flags + IP, má HW stack

## • MIPS

- r0-r31 = 32-bit obecné registry, ale některé mají speciál význam
  - r29 = stack pointer - ukazuje na vrchol zásobníku
  - r30 = frame pointer
  - r31 = return address } slouží pro volání funkce  
↳ jiné registry pro jal instrukci ~ místo caller
- r0 = zero - je vždy nula

→ nemá HW stack ani flags

- když se zavolá fce, tak instrukce jal složí na tu fci a do r31 zapíše kam se vrátit  
současně někde na zásobník uloží aktuální stav registry a slibuje, že  
je po návratu obnoví - kromě nějakých temp. registry
- ↳ na return value je nějaký speciál registr

↳ preserve registry

## • Instrukce MIPS

$a \leftarrow b \text{ op } c$

formá kódování

→ 1 instrukce = 32 b.

x86

$a \leftarrow a \text{ op } b$

proměnlivá della instrukce

## • ABI = Application Binary Interface

→ specifikuje, jak a na co se mají které registry používat

→ většinou to vydává autor ISA

→ všechny kompilery ho musí dodržovat, jinak by nebyly kompatibilní

→ říká aliasy registrů

→ jaké registry jsou preserve = volání fce je nezmění

## • Příznaky

- nemají je všechny ISA

- běžné jsou zero, sign a carry

- jsou systémové a uživatelské příznaky - x86 je má promíchane - musí se to složitě maskovat

- x86 sada - ne všechny instrukce nastavují všechny příznaky, ...

## • Klasifikace instrukční sady

• CISC - komplexní sada = spousta speciálních instrukcí ⇒ hodně tranzistorů

↳ ty instrukce se pak překládají do mikroódu & mikroinstrukcí - ale na 1 simple instrukci je 1 mikroinstruce

• RISC - redukovaná sada = jenom ty základní instrukce

↳ rychle se to dekoduje + stačí málo tranzistorů

• VLIW - Very Long Instruction Word ⇒ třeba 128 bitová instrukce

↳ ta instrukce se nemusí dekodovat ⇒ fast - třeba v síťových switchích

• EPIC - Explicitly Parallel Instruction Computer

↳ v instrukcích je explicitně co má běžet paralelně - dnes to dělá HW

↳ třeba aritmetika

→ Load-Execute-Store = instrukce pracuje jen s registrami a nesáhá do paměti

## • Příklad instrukcí

- instrukce dekoduje assembler → ale textovým zápisu instrukcí se toly říká assembler

## • Co všechno je v procesoru?

- řadič paměti → R nij vede sběrnice do paměti

- hierarchie cache

- jádro / jádra

- registry

- logické procesory = unit 1 jádra je víc proudů instrukcí najednou → intel → hyperthreading

## • Zjednodušené schéma

- 1 jádro má:

• výpočetní jednotku ~ *execution unit*

• cache - 3 úrovně

~ 102B L1 → slovo stejně rychlá jako registry + dělení na *instrukční* a *data* *doboru*  
~ 1002B L2 → větší + o řád pomalejší + unifikovaný = kód a data  
~ 1MB L3 → sdílená mezi jádry + zase o řád pomalejší

• vlákna - více vláken sdílí 1 výpočetní jednotku :: nebývá vyžitá na 100%  
- ale každé vlákno má svoje registry

- dělá se prefetch = HW odhadování co se má načíst

- L3 je ve skutečnosti rozložená na L3 slices

→ každé jádro má svůj a spojuje je obousměrný ring

→ přístup k L3 cache trvá různé dlouho

← mega rychlá sběrnice

## • Schéma jednoho jádra → Intel Coffee Lake

- Front End čte instrukce a dekoduje je

- Coffee Lake má 5-cestný dekodér = dekoduje 5 instrukcí v 1 kladu

↳ 4 jsou simple dekodéry na simple instrukce → přebírá se na 1 mikroinstrukci

↳ 1 je komplexní na ty komplexní instrukce co se překládají do mikroinstrukcí

- po dekodování instrukce spadnou do "bazénu" a tam čekají na zpracování

↳ Out of Order Execution - celkem chaotické "pool"

↳ instrukce se probírají, ale pořadí musí mít navenek stejný efekt

→ v tom bazénu plavou ty mikroinstrukce, přičemž každá má nějakou "barvu" a čeká, až se jí otevře port příslušné barvy → probere jím a vykoná tu svou operaci

→ každý port umí jímé instrukce

→ nakonec to vše do reorder bufferu, kde se ty výsledky překládají do správného pořadí, aby ten výsledný efekt byl správný

## • CPU pipeline

- vykonávaním 1 inštrukcie sa rozdelí na 14-19 fázi = stages
- v 1 okamihu mám rozdelených viac inštrukcií a každá je v inej fázi
  - ↳ inst. fetch, inst. decode, execution, čítanie z pamäte, zapísanie výsledku, ...
- Každý takto posunú o 1 fázi ďalej
- náběh pipeliney = keďže nabíhajú prvú inštrukciu a ešte neprobujú všetky fáze
- doběh reálne nastane - to penom v SW
- keďže vypadne, tak musí znovu nabíhať - pomalé
- ⇒ podmínene sloky sa snaží CPU predvídať (jejich výsledok) } ≈ 95% to odhadne správne  
↳ napr. cykly
- ↳ pri špatnom odhadu musí celou pipeline zahodiť
- návrh procesoru je jednoduchší a overall to je rýchlejší, ale vzniká latencia

↳ musí prebehnúť rovnako rýchlo

## • Super-scalarita

- v každej fázi sa nachádzajú viac inštrukcií súčasne
- 2-cestný procesor je 2x rýchlejší než 1-cestný
- dnešní procesory majú niekedy 5-cestné
- môžu byť asymetrické → viac simple a komplex inštrukcie
- keďže nastane konflikt (registri, výsledky, ...), tak na sebe by inštrukcie čakajú

za jednotku času

2x viac inštrukcií

## • Cesta inštrukcie

Pipeline decoder → fool → probe portem → reorder buffer

## Paměťová hierarchie

• volatile → reg., cache, RAM

• persistent RAM → non-volatile, rychlejší než SSD, CPU & má přímý přístup

• externí paměť → SSD, flash disky, HDD → je potřeba řadič a sběrnice

• archivní paměť → magnetické pásy

→ je na to potřeba upravit OS a programy

## Adresový prostor

→ očíslovaný slovo ⇒ 32 bit adresový prostor znamená  $2^{32}$  slov

→ dnes 1 slovo = 1 byte

→ fyzicky: byty ve 2D poli

→ řadič paměti nejprve vybere řádek a poté indexuje sloupci v něm

→ pro jinou řádku ji musí odvybrat → to je pomalé

⇒ sekvencím přístup je nejrychlejší

→ doba přístupu

• CAS = # loktů než můžeme adresovat další sloupec v dané řádce

• RAS = čas na aktivaci řádky

## Reprezentace dat

• unsigned int:  $0 - 2^m - 1$

• signed int: dvojkový doplněk:  $-2^{m-1}$  až  $2^{m-1} - 1$

• float: podle IEEE 754 → mantisa a exponent s biasem → CPU to poté převede na nějakou svoji reprezentaci

$$\Rightarrow V = (-1)^{\text{sign}} \cdot \text{mantisa} \cdot 2^{\text{exponent} - \text{bias}}$$

## Endianita

→ big endian = MSB first → ostatní + síťové protokoly

→ little endian = LSB first → intel

→ LLL → LE má LSB na Lowest adrese

## Data Alignment

→ datový typ velikosti X bytů musí ležet na adrese A:  $A \% X = 0$

• vnitřní zarovnání → každý typ ve struktuře je zarovnán na násobek své velikosti

• vnější zarovnání → struktury v poli by musí být aligned → je tam double? 20B ⇒ 24B

⇒ velikost toho strukturu musím zarovnat na násobek největšího typu v ní

→ `sizeof(struct)` vrátí velikost včetně vnitřního i vnějšího paddingu

→ tohle řeším jen pro základní typy ⇒ nebudu to řešit pro struct ve struktuře

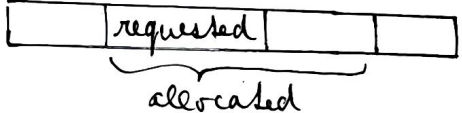
## • Typy proměnných

- global - nachází se při spuštění programu a smáče se při ukončení
- local - alokované na úrovni bloku, vznikají v bloku { a končí v bloku }
- dynamicky alokované - sám si je musím vytvořit → (C++ malloc(), (# new  
↳ instance třídy → sám je musím odstranit nebo to dělá garbage coll.

## • Alokace paměti

- úkol je najít blok nevyužití paměti dostatečné velikosti
- spustím program → před tím, než se pustí main(), tak si runtime toho jázgly od OS vyžádá heap = velký souvislý úsek paměti  
→ z toho heapu pak alokuju  
1, alokace - řeknu si o blok dané velikosti, dostanu adresu = pointer na tu proměnnou  
2, používání bloku  
3, explicitní uvolnění bloku / G.C.  
→ heap si pamatuje zabranou paměť → po větších blocích - min. 64 B

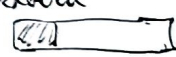
## • Fragmentace paměti

- interní  → chci 48 B ⇒ dostanu 1 blok = 64 B

- externí → mám hodně volného místa rozdrobeného do malých nesouvislých úseků

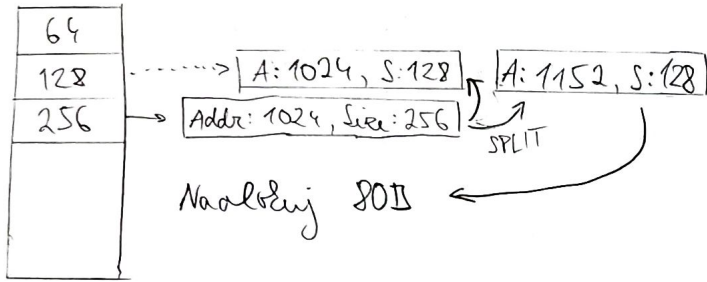
## • Dynamická alokace paměti

- pořad alokuju věci s proměnnou délkou
- jak si pamatuju volné bloky? spojím volných bloků  
bismapa bloků
- alokací algoritmy

- first fit → vezmu první volné místo a rozstřenu ho  
↳ easy & fast, ale může nenechat velké bloky 
- next fit → pamatuju si, kde jsem skončil a z toho místa pokračuju  
↳ trochu méně fragmentuje, ale musí si pamatovat, kde skončil
- best fit → projde to celé a vezme ten nejmenší, kam se ještě fitting  
↳ nechá velké bloky, ale vytvoří spoustu malých děr  
⇒ ta externí fragmentace může být ve výsledku ještě horší
- worst fit → projde to celé a dá to do největšího bloku  
⇒ očividně nejhorší, nicí velké bloky



- Buddy memory allocation - alokační alg. co se actually používá
  - bloky o velikosti  $2^n \Rightarrow$  adresy aligned na  $2^n$
  - mám pole spojáků bloků o velikostech 64, 128, ... 168
  - když mám naalokovat 128B blok a mám nejvýš  $\rightarrow$  prostě ho dám ✓
    - ↳ pokud ho nemám, tak najdu nejmenší  $2^n$  blok kam se to vejde
  - $\Rightarrow$  ten rozštěpím na 2 buddies, jeden si nechám a druhý naalokuju / štepuju dál
  - těch spojáků je konstantně mnoho  $\Rightarrow$  alokace je konstantní



- Dealokace  $\Rightarrow$  počusím se je mergeovat
  - ↳ adresa těch buddies se liší jen o 1 bit - tu velikost  $\rightarrow$  chci ho shrnout
  - $\Rightarrow$  adresa kamovádla je  $Addr \oplus Size$

- dealokace není konstantní  $\because$  musím projít spoják pro daný size
- externí fragmentace je dobrá, ale má to velkou interní fragmentaci

## • Cache

- úroveň cache*
  - = datová struktura, kde mám data co často používám / dlouho se používají
  - HW ale i SW implementace
  - pokud cache děje, tak musím vybrat oběť  $\Rightarrow$  stránovací algoritmy *→ předěje*
  - když začádám o data co tam nejsou, tak se fetchují a přístě to bude fail

## • Cache v procesoru

- spoléhá na lokalitu přístupu = většinou přistupuju k datům co jsou vedle sebe
  - $\Rightarrow$  přednáčítá si ty data  $\Rightarrow$  schovává lokenci přístupu do RAM  $\rightarrow$  95% hit
- při requestu probledává od nejrychlejší L1 po nejpomalejší L3
- koherence kesi = problém když více jader přistupuje ke stejnému datům
  - $\Rightarrow$  ta data se musí přesunout do té společné L3 cache
- když nějaké jádro používá nějaká data hodně, tak jsou hot a jdou do L2  $\rightarrow$  L1

## • Terminologie kesi

- cache line = 1 blok v kesi (není organizovaná po B) - standardně 64B
- cache hit = když requestím data, co mě jsou načesovaný  $\sim 95\%$
- cache miss = když ty data ještě nemám načesovaný  $\rightarrow$  cache line load
- cache line load = načítám data z paměti, pokud je cache plná  $\Rightarrow$  oběť
  - ↳ pokud jsem oběť v kesi zmodifikuju, tak je musím aktualizovat v paměti
- cache line state - používá se MESI protokol  $\rightarrow$  4 stavy kesi kajíny

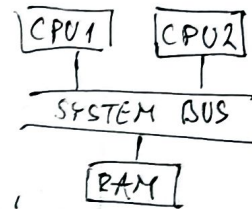
## • Asociační paměť

- RAM ku očištnu cache linama ⇒ dvojce (klíč, hodnota)
- je to vlastně HW slovník ⇒ mega fast, ale stojí to HODNĚ tranzistorů
- snaha je, aby cache byla plně asociační
  - ↳ většinou je třeba 4-asociační = 4 různých cache line mají stejný klíč
  - ↳ větším bitů té adresy, ale vznikají kolize

## • Systemy s více procesory

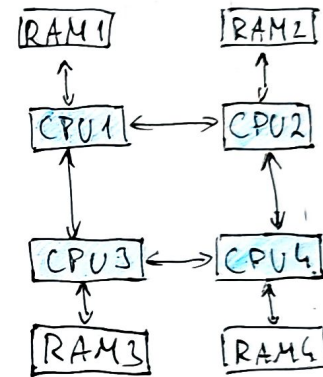
### 1) SMP = Symetric Multi Processing

- 1 sběrnice mezi sdílenou RAM a CPUs
- simple, ale pro více procesorů je ta sběrnice pořád zablokována

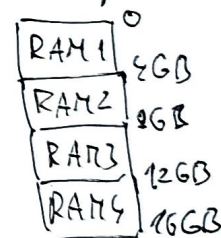


### 2) NUMA = Non-Uniform Memory Access

- Každý procesor má svůj kus RAM
- Aby RAMky musel být dostupné i z ostatních CPUs
- jsou navzájem propojené → je to mega drahý
- CPU mají 1/2/3 pinů pro tu rychlou sběrnici
- doba přístupu do těch RAMek je různá ⇒ NUMA factor
- Aby RAMky tvořily 1 společný adresový prostor
- té dvojici [RAM] - [CPU] se říká NUMA uzal



Address space

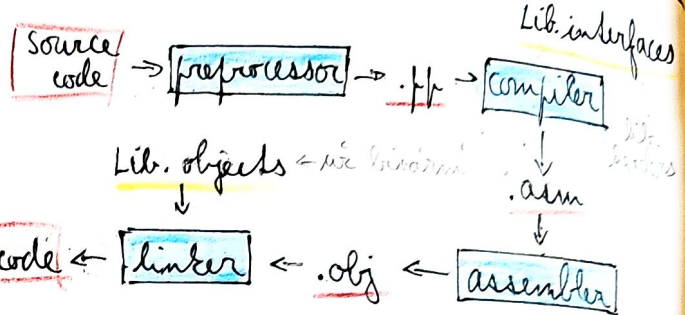


# • Programovací jazyky

## • Překladač

- formálně to je zobrazení slov ze vstupního jazyka generovaného gramatikou do výstupního jazyka generovaného jinou gramatikou / přijímaného automatem
- jazyk má nějaká pravidla a lexikální elementy (while, do, for, ...) <sup>CPU</sup>

1. Preprocessor: zpracuje zdroják a příkazy  
 pro něj → C: #include, #def, make, ...  
 ⇒ vytvoří .pp soubor - stále textový zápis



2. Překladač: dostane .pp a rozbrání knihovnu ⇒ .asm assembler texták

3. Assembler: z .asm vyrobí binární .obj → ta binární forma závisí na OS

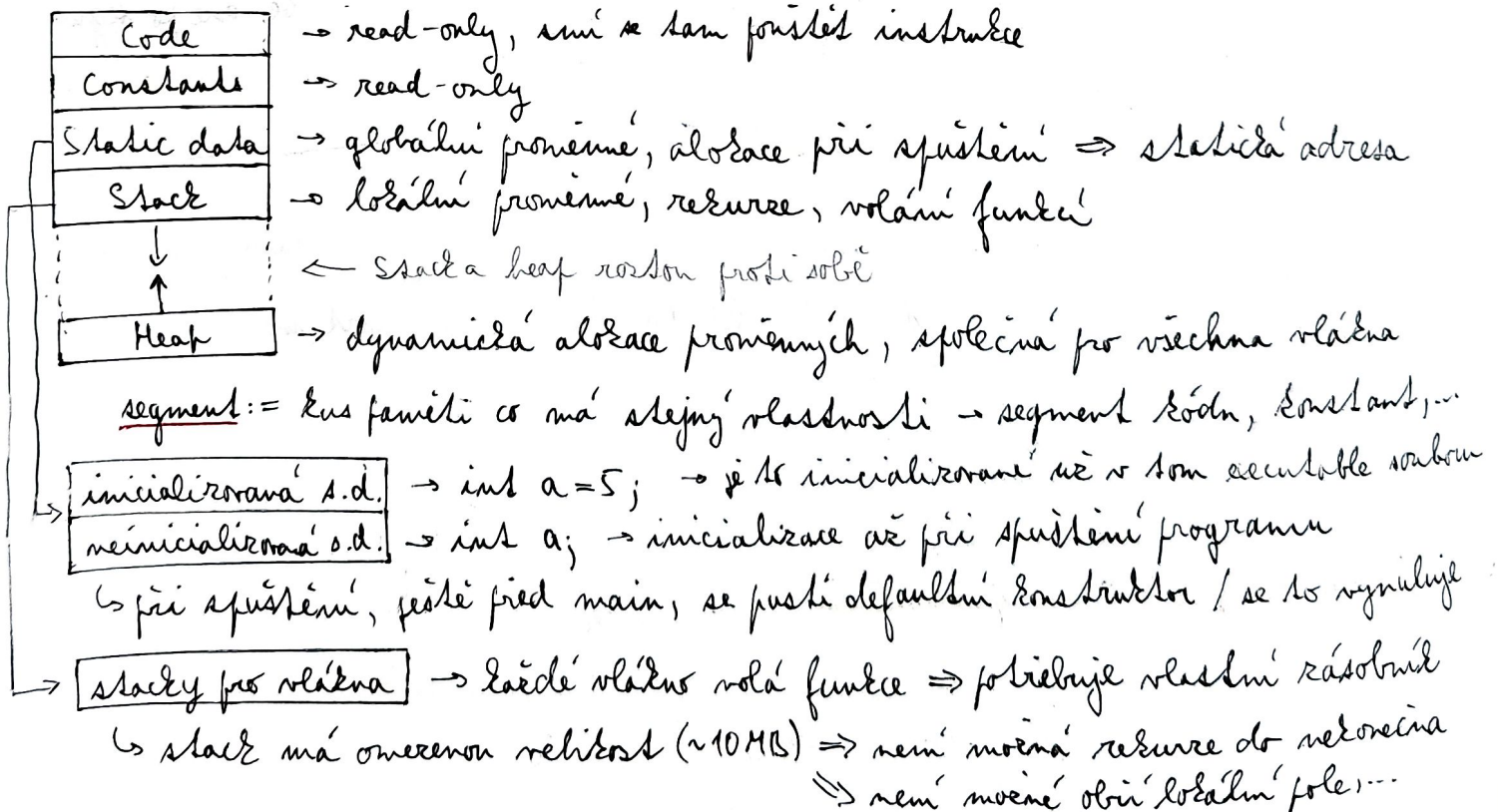
4. Linker: až takhle zpracuju všechny zdrojáky na .obj, tak to spolu s objectama knihovnou dostane linker a vyrobí výsledný executable

→ dneska se 1.2.3. dělá na 1 spuštění překladače

→ OS pak při spuštění udělá z programu proces

⇒ musí načíst do paměti nějakou DS s daty toho procesu

## • Organizace paměti při běhu programu



Knihovny

= kolekcce sкомпilovaných zdrojových modulu v 1 souboru → std. lib. C je 1 soubor  
 → vlastně jen víc objektů

• statické - 1 soubor z více modulu a linker z ní přímo kopíruje ty potřebné části do toho výsledného executable souboru  
 .lib

• dynamické - do toho executable souboru si jenom poznamenám, že potřebuji nějakou fci z této knihovny a až za běhu to sežene loader, nahraje do paměti a upraví příslušné adresy v tom executable  
 .dll

→ statická vyrobí větší .exe, ale ty funkce tam jsou zabudované → rychlejší  
 → před všechny programy využívají nějakou std. lib., takže dává smysl dynamická

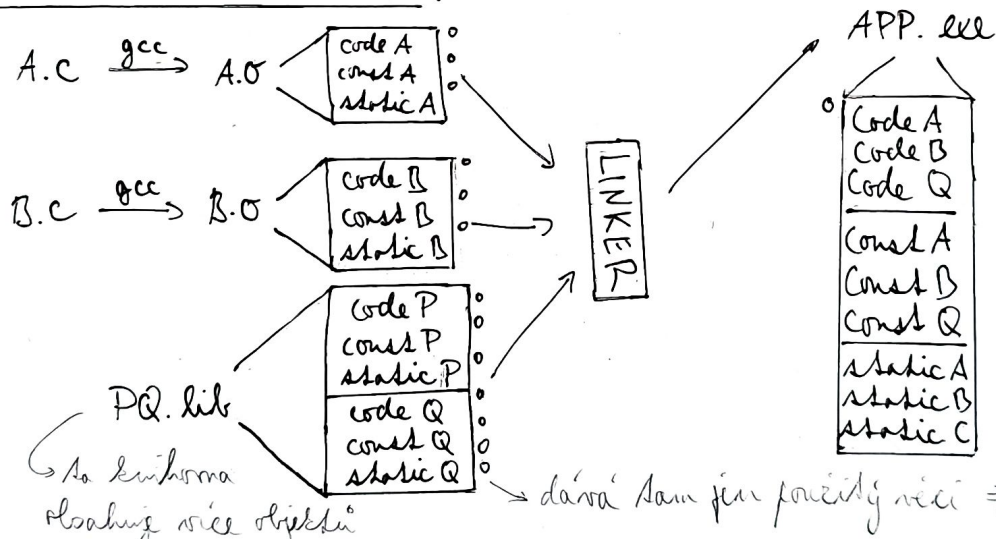
Linkování

→ linker posílá všechny ty přeložené objekty + s. knihovny + rozhraní d. knihoven  
 → vyrobí ten výsledný executable pro daný OS → ten z něj umí vygenerovat proces

Loader

→ část OS, co má za úkol načíst ten executable a dynamické knihovny do paměti a spustit to

Co všechno linker dělá?



→ seřazení není přímo za sebou, ale po segmentech

→ reorganizace daných knihoven

→ musí se toky řešit to, že před linkováním nevíme, jakou budou mít nějaké funkce, konstanty, ... adresu ⇒ překladač do všech těch instrukcí zapíše relativní adresu od začátku segmentu, ve kterém ta věc leží

→ překladač do toho objektu píše relozace - pamatují si instrukci + ten segment

⇒ linker ty instrukce musí po seřazení opravit → pro linker je základní adresa toho souboru 0

⇒ že každé té relativní adrese přičte základní adresu toho jejího segmentu po seřazení

⇒ loader při spuštění dělá další relozace - víc základní adrese toho executable

→ linker hledá entry point = objekt, kde to začne běžet → není to main(), ale std. lib.

→ knihovna publikuje jaké má fce, konst., ... ⇒ v. obj sežce public → zliční slovo public

→ zdroják říká, že chce externí fci ⇒ v. obj sežce extern → tohle se nadace exportuje main

⇒ linker postupuje rekurzivně od entry pointu a pro externy hledá publiky ⇒ přidává moduly

## • Behová podpora jazyka = Run-time

• statická - compiler + realizace rozhraní knihoven

• dynamická - konstruktory a destruktory globálních objektů  
organizace paměti + obsah paměti před spuštěním  
volací konvence funkcí  
dynamické knihovny toho jazyka

→ to je potřeba pro linker

→ to je potřeba pro běh

## • Volání funkce

stack frame

→ funkce při zvolání dostane aktivací rámec a v něm pracuje  
nikam jinam nesáhá - kromě globálních proměnných

přípraví volající

return val. parameters
---------------------------

→ malé věci / reference se vrací registry

→ pokud je to nějaký velký objekt, tak volající připraví takle místo

tohle volání  
musí udělat

return addr.
control link
machine state

→ kam se poť má vrátit

→ frame pointer aktivacího rámcu volající funkce

frame pointer

local data
temporaries

→ stav registru co se musí zachovat

prostor volaného

local data
temporaries

→ lokální proměnné

→ když dojde registry

→ je frame pointer co obsahuje frame pointer aktivacího rámcu

→ soubor celý aktivacího rámcu je na zásobníku

## • Volací konvence

→ řídí se jí všechny překladače, aby to bylo kompatibilní

• public name mangling → linker jen hloupě pojmenovává řetězce publik a externí  
"mandlování" jmen  
→ znakování / předělání  
⇒ aby fungoval např. přetížení / více konstruktory

• musí se domluvit kdo co ukládá z toho zásobníku + call/return sekvence

• předávání parametrů → na stacku / registry + v jakém pořadí se na zásobník dávají

• návratová hodnota → buď registrama / přes stack  
→ c je standardní způsoby dělení

• registry - které se musí zachovat (preserve) a které ne (scratch)  
→ jaké jsou role registru

## • Call/return sekvence = zodpovědnosti při volání

→ volající předává parametry a po návratu volané funkce je i smaže

→ volaná funkce smaže svoje lokální data, obnoví machine state, nastaví frame pointer reg. na control link, vrátí se a smaže tohle všechno

→ mazání = posunutí pointeru na vrchol zásobníku

→ existují volací konvence, kde funkce volaná ukládá i své parametry  
→ je to jen pro fázi počítání parametrů

## • Předávání parametrů

→ čistě C umí jen hodnotou

- hodnotou - první kopie a pak je to pro tu volanou fci jako lokální proměnná
  - referencí - předává se adresa, v C++ pomocí & → ta fce ví, že to není lokální var  
↳ na ten stack se dá ta adresa → by instancí funkce je nepřijímá adresou
- čistý C řeší předávání referencí tak, že hodnotou předá pointer  
↳ rozdíl: pointer není fixní + musím explicitně psát tu \*, takže je to otrava

## • Proměnné

= pojmenovaný kus paměti, který drží nějakou hodnotu

→ má nějaký typ (ve staticky typovaných jazycích)

→ kde běží?

- C/C# {
- static data - globální proměnné → abžace při spuštění
  - stack - lokální proměnné → abžace na zásobník při běhu programu
  - heap - dynamické proměnné → abžace na heap pomocí new (C#) / malloc (C)

python { slovník - (klíč, hodnota) = (jméno proměnné, info o proměnné)  
↳ takže používá třeba python, javascript, ... ↳ typ, hodnota, adresa, ...

## • Heap

→ dynamická paměť, kterou si při výpočtu můžu brát

→ není to ten stack, takže se to při návratu z funkce nemazá jako lokální proměnné

→ funkci volaná může něco dynamicky naalokovat a vrátit fci volající pointer

• alokace - evidence volných bloků, alokační algoritmy

→ je to explicitní → např. new list<> → takže vrátí pointer

• dealokace - C, C++ jde explicitně udělat free/delete

↳ rozpoznávalo se na to + musely se řešit exceptions

↳ memory leak = nemám GC a ztratím všechny pointery na tu paměť

C#, Java

• Garbage collector = automatická dealokace, docela pomalé algoritmy

⇒ mega fast abžace

⇒ GC se nezavolá dokud nedojde paměť ⇒ nemusím řešit alokační algoritmy

→ GC je pomalý + není možné abžovat když běží ⇒ freeze

→ hodně programů spoléhá na to, že se GC nikdy nezavolá

⇒ odstraňuje fragmentaci

⊕ → heap consolidation = po tom co GC doběhne, tak ty aktivní bloky se sčepu na stranu...

⊖ → ovlivňuje to výkon, může dojít k zamrznutí v nepředvídatelný moment

• GC strategie - rozdílně se používají podrobněji

→ měřou to být nějaké objekty

• Arasování → projde aktivní seznamy a najde živé proměnné → z nich to jde dál rekurzivně

↳ marking

↳ naonec smaže všechno co není označené jako živé

• počítání referencí → když udělám novou lokální referenci ⇒ ++count

↳ když ta lokální proměnná o ton referenci zanikne ⇒ --count

→ když count = 0, tak udělám free

→ pozor na cykly / odlehle komponenty souvislosti - arasování. Semhle problém nemá

## • Přenositelnost kódu

→ aby to fungovalo na více OS, přeladačích, architekturaích, ...

### 1) CPU architektura

- C, C++ → různé velikosti typů, třeba int je podle architektury
  - C#, Java → fixed velikosti typů → int je všude 32 bit
- někdy je třeba řešit i endianitu dat

→ ale je možné si vybrat, že třeba int bude 32 bit

→ třeba univerzální program jenom jednou

### 2) přeladač → gcc, msvc, clang

- C++ má více přeladačů s různými specializacemi
- ⇒ řešení: používat syntaxi a knihovny ze standardního jazyka
- přeladač se snaží dodržovat ISO standard toho jazyka
- třeba C# nebo Java mají jen 1 přeladač ⇒ ten problém tam není

### 3) OS

- různé sys-cally, ale funkce to má stejné → OPEN, ...
- jen různá jména a interfacy (někdy to může vypadat stejně)
- podmíněný přelod = ve zdrojovém označím co se má přeložit podle OS

→ systémové volání

→ Java a C# tohle všechno dělají sami za runtime ⇒ je to součástí sebe

## • Přenositelnost pomocí Virtual Machine

- Java, C# se přeloží do instrukcí procesoru co neexistuje ⇒ mezikód
- ten mezikód dostane nějaký nativní VM, co je za runtime interpretuje
- ↑ bezpečnost a kontrola → zavře se to do sandboxu a ten interpret kontroluje, jestli to nedělá něco co nemá
- řešení problému rychlosti

### • Just-in-Time = JIT

→ jednotlivé funkce, třídy, ...

- postupně za běhu se mezikód přeládá na nativní kód a to se řeší rychle
- ⇒ nemusím to přeladit dvakrát - když volám poprvé nějakou funkci, tak se přeladí podmínky už se zrovna přeládá

### • Ahead-of-Time = AOT

- přeloží se to celé při instalaci a spustí se to už přeložené

Java: bytecode

C#: CLR

## • Operační systémy

→ nemá přesnou definici, ale má 2 hlavní úkoly

1, abstraktní stroj = abstrakce nad HW

→ reprezentovaný kernel API → systémové volání (jsou schované v nějaké knihovně)

→ schováva před programátorem tu HW implementaci a poskytuje to C rozhraní

2, resource manager = přidělování prostředků procesům

→ manageje všechny HW a přiděluje aplikacím HW prostředky

→ ty aplikace to musí nějak sdílet mezi sebou

→ alokace (mem), time-sharing (CPU), abstrakce (disk, síť, grafika ...)

## • Režimy procesoru

1) uživatelský mód → pro všechny aplikace, nemá přístup k nějakým reg. + instrukcím

2) kernel mód → používá ho OS nebo jen nějaká jeho část - zbytek běží v user

↳ kernel = ta část OS co běží v kernel módu = má full access

• přechod kernel → user : OS to dělá triviálně - je na to special registr

• přechod user → kernel → např. některé sys-cally potřebují kernel režim

OPEN  
PRINTF

entry pointy do kernel režimu = způsob jak to přepnout

↳ jsou jasně definované a OS je kontroluje → třeba na to je special instrukce

↳ po přepnutí ta aplikace skočí na nějaké jasně definované místo v paměti, kde pro ni OS připravil ty instrukce

↳ parametry pro ten sys-call (jméno souboru) se předávají registry

→ entry point může být i to, že se ten program pokusí udělat nějakou systémovou instrukci, což CPU pozná že je chyba a předá řízení OS

↳ ta sys instrukce je nějak definovaná pro ten přechod



## • Architektura OS

• monolitická → mohli se dělaty první OS → Linux kernel

- velká globální věc, která celá běží privilegovaně
- nemá to pořádnou strukturu - kromě entry-pointu, ke kterému je nějaká volba servisních procedur co pak vedou dál na utility procedury

⊕ velmi efektivní, hodně rychle sys-cally

⊖ v tom velkém programu prostě budou chyby + když se někdo dostane do OS, tak může dělat co chce a má přístup k celému počítači

⊖ původně se to ani nedalo rozšiřovat, dnes je možné načítat dynamické moduly ⇒ je tam ještě větší prostor pro chyby

• vrstevnatá → evoluce těch starých monolitů

→ organizace OS do vrstev s různými službami a oprávněními

→ vrstva  $n+1$  může využívat pouze služby vrstvy  $n$

⊕ bezpečnější a jednodušší rozšiřitelné

⊖ návrh rozhraní mezi vrstvami je velmi náročný

• mikrokernel

→ snaha je udělat kernel co nejmenší → třeba méně než 1MB

→ nad tím kernelem jsou jednotlivé moduly (filesystem, řadič disku, ...)

→ kernel zajišťuje komunikaci mezi těmi moduly

→ posílání zpráv jako mezi klientem a serverem

⇒ nezabírá celý OS

→ když nějaký modul spadne, tak ho kernel zase spustí

→ kernel i ty moduly mají svoji protected paměť, do které si nevidí

⊕ bezpečné, spolehlivé, snadno rozšiřitelné

⊖ je to hodně pomalé - kvůli posílání těch zpráv

→ windows se snaží být mikrokernel

→ čím mikrokernel už ani nezávisí na architektuře CPU, protože ještě pod ním je Hardware Abstraction Layer

• Zařízení = HW soustava sloužící k nějakému účelu

• device controller = řadič → řídit to zařízení

HW

→ elektricky komunikuje s tím zařízením, řídí signály na dráhu, A/D konverze, ...

→ těch zařízení k němu může být připojených víc ⇒ mají nějakou Apologii

• device driver = ovladač → 1 driver bude pro nějakou konkrétní zařízení, se pro 1 konkrétní

→ část OS, komunikace už probíhá pomocí instrukcí CPU → rapší něco nějakým způsobem

SW

→ poskytuje abstraktní rozhraní pro vyšší vrstvy OS

aby OS věděl jak s ním komunikovat

→ při bootu (BIOS, UEFI) se projdou zařízení a začne dostane nějakou adresu

→ info o těch zařízeních se rapší do tabulky, kterou pak dostane OS

by zařízení dříve bude chvil 1MB

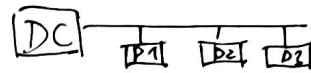
# • Topologie zařízení

DC = Device Controller

• sdílená sběrnice → hlavní dříve

→ musí se řešit arbitrace pro hodně zařízení

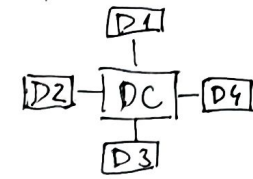
→ protože dot je omezený ⇒ pro hodně zařízení špatně



• star → dnes u disků

→ zařízení je p2p připojené k radičce

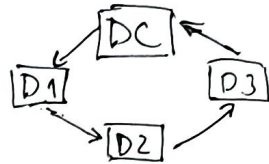
→ nevýhoda je složitější konstrukce radičce



• ring → více se nepoužívá

→ orientovaná sběrnice se zařízení

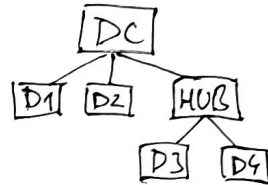
→ levnější řešení, velká latence



• tree

→ DC je v kořeni, může se většit pomocí hubů

→ radič i ovladač musí být schopni ovládat změny topologie toho stromu



## • Komunikace se zařízením = Device Handling

→ tohle řeší OS, většina kódu OS je nezávislá na ovladačích a radičcech

→ programátor může udělat přímo sys-call, ale pro přenositelnost je lepší zavolat std. lib, Samozřejmě

→ pro otevření souboru zavolám funkci fopen(), která udělá syscall OPEN

↳ OS dá tomu souboru handle = ID, podle kterého bude identifikace při READ

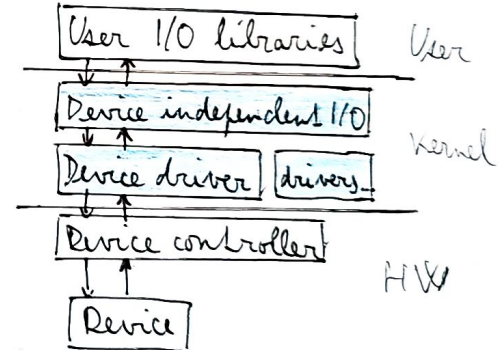
→ ta read žádost se dostane až k ovladači

→ ovladač dá request radičce toho disku

→ ta HW komunikace mezi radičem a diskem může probíhat vícerorát než se požadavek vyřídí

→ až budou data ready, tak dá ovladač té standardní knihovně adresu na ni → ta jí dá mě

→ problematika je ta část, kdy ovladač čte na ty data



⇒ domluva mezi ovladačem a radičem

• polling

→ CPU periodicky kontroluje stav toho zařízení a čte až bude hotovo

→ musím aktivně čekat + to má latenci kvůli periodě

• interrupt

→ zařízení pošle CPU signál, že už má hotovo ⇒ CPU přerušit aktuální proud instrukcí a skočí do obsluhy přerušování - ten to vyřídí a vrátí se

→ radič je drátem připojený k CPU pro ten signál ⇒ CPU na to potřebuje pin

• DMA = Direct Memory Access

→ DMA radič dočká přesunutí ty vypočtený data z radiče toho zařízení do paměti ⇒ CPU nemusí nic řídit ⇒ dříve se čte data nebo opatně do zařízení

→ scatter/gather → DMA dočká z nesouvisejících dat vytvořit 1 blok nebo naopak ten blok rozdělat → třeba když přijde packet, tak ho rozseká na

TCP header, IP header a ty aktuální data

• Typy přerušeni = interruptu

• externí → z venku přes nějaký pin na procesoru  
↳ je možné ho zamaskovat = zakázat / povolit → v kernel módu

↳ když budu dělat nějakou práci

• exception = výjimka

↳ třeba nedefinovaná instrukce, dělení nulou

→ nevoláně vyvolaná nějakou špatnou instrukcí → CPU ji vyvolá sám  
→ u všech těchto výjimek je definováno co se má stát když nastanou  
⇒ třeba se skočí na nějakou pevnou adresu

→ využít pro emulaci nových instrukcí na starých procesorech

↳ když přijde ta instrukce, tak nějak skočím, spočítám to a vrátím se zpět

• fault = ta instrukce není schopna dobehnout  
⇒ udělá se rollback a to přerušeni se udělá před tou instrukcí

↳ může to být kvůli rozložení

• trap = ta instrukce se provede a až pak nastane to přerušeni

• softwarové

→ speciální instrukce co vyvolá přerušeni ⇒ skočí do OS

→ slouží jako entry-point OS pro sys-cally

→ pro ty přerušeni je reálně nějaký řadič, který CPU posílá číslo toho přerušeni pomocí nějakého protokolu ⇒ CPU pak skočí do obsluhy přerušeni

→ adresy těch obsluh jsou buď fixní nebo v interrupt table

• Processing → OS řeší konstantní programy

aktivita procesu



• proces = instance programu, je to nějaký objekt OS

→ vytvoří ho OS pomocí sys-callu

→ OS mu poskytl prostředky a dělá si jejich evidenci, aby si je mohl vzít zpátky až ten proces skončí

prozradit

• vládnos = thread

→ 1 proces může mít více vláden a běžet paralelně

→ je to základní jednotka kernel schedulingu

} Každé vládnos má vlastní zásobník

Kontext procesora = DS pro ukládání stavu vládnos

↳ pamatuje si starý registry → je tam i stack pointer

↳ když CPU přepíná vládnos, tak musí udělat context-switch

↑ přechod na další vládnos

• fiber

→ ještě lehčí jednotka plánování než vládnos

→ každý má svůj kontext a zásobník

} OS o něm není, dělá se to v lib

→ umožňuje kooperativní plánování

↳ fiber 1 řešne "já už nechci běžet" a vybere další fiber

instrukce

## • parent-child processy

- proces může spustit další proces a protože zná jeho PID, tak počkat až dobehne
- dítě si pamatuje PID rodiče → PPID → když skončí, tak mu pošle signál
- zlikvidovat procesy nekillne jeho děti ⇒ adoptuje je inuit proces PID 1  
⇒ zabit strom procesů je celkem těžké

• Scheduler = část OS co přiděluje výpočetní zdroje vláknům ⇒ plánovací algoritmus

• Multitasking - na 1 jádru běží více procesů současně ⇒ je tam nějaké přebíhání

• Multiprocessing - mám více procesorů / jader ⇒ pro ten scheduler je to náročnější  
→ affinita procesu = snaha aby tohle vlákno běželo pořád na stejném jádru  
↳ kvůli lepšího jádra

• Real-time scheduling → jiná verze schedulingu pro nestandardní situace

↳ RT process má začátek (release time) a konec (deadline)

⇒ musí dobehnout někdy v tomto intervalu → tohle je vraga, ale když uplánuje

↳ hard deadline → nemá cenu to dodělat, když jsem to nestihl  
↳ soft deadline → pořád to má cenu dopočítat

⇒ třeba airbagy v autě nebo řízení jaderné elektrárny

## • Stavový diagram plánování

→ vlákno = unit of scheduling = jednotka plánování

→ to vlákno má svoje stavy

• created → čeká to na spuštění (release time)

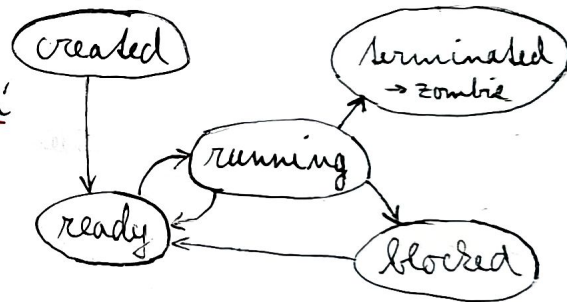
• ready → schopný běžet, ale nemá přístup k CPU

• running → pokud už běží moc dlouho, tak zpět do ready

• blocked → čeká na zdroje - třeba čtení ze souboru / user input zombie stav

• terminated → už skončil, ale rodičovský proces si ještě nepřičítá výsledek

↳ v tomto stavu si OS bere zpět všechny prostředky



• Multitasking → mám 1 jádro a chci na něm stádnat vlákna → jak na to?

→ kooperativní → OS se o to nestará a všechny vlákna spolu spolupracují  
→ možné jenom ve speciálních systémech

## • preemptivní

→ každé vlákno dostane time slice (bílka 20 ms)

→ OS má nějaký externí HW zdroj přerušení → timer

→ pokud to vlákno během toho time slice nestončí nebo se nezabloudí,  
tak na konci time slice dostane interrupt a jde do stavu ready

→ OS udělá context-switch a pustí tam jiné vlákno

## • Scheduling - cíle

- maximalizovat využití CPU a jader → pokud to potřebuju → většinou tohle tohle
- správně rozdělit výpočetní zdroje mezi vláken
- maximalizovat propustnost = # vláken co dokoním za jednotku času
- minimalizovat dobu jednoho vlákna
- minimalizovat čekací dobu v ready stavu
- minimalizovat čas na odpověď v interaktivních aplikacích
  - ↳ když kliknu na okeno, tak musí rychle změnit stav blocked → ready → running a něco udělat

## • Priorita procesu

- = číslo co stanovuje důležitost procesu → a vlákna tu prioritou většinou dělá
- priorita je součet statické a dynamické priority
- statická - přiděluje se při spuštění procesu a může se pouze snižovat
- dynamická - jednou za čas se ready vláknem inkrementuje
  - ↳ když přejde do stavu running, tak se vynuluje

## • Plánovací algoritmy

- nepreemptivní ~ kooperativní → strašně simple
- First Come First Served - fronta, vypočítám vlákno a jde další
- Shortest Job First - musíme předem vědět horní odhad execution time
  - ↳ maximalizuje propustnost = throughput
- Longest Job First
- preemptivní
- Round Robin → jako FCFS
  - je tam prioritní fronta + každé vlákno má time slice
- Multilevel feedback queue → dynamický algoritmus
  - několik front podle toho jak dlouhý timeslice ty vlákna chtějí
  - horní fronta má nejmenší timeslice a spodní největší
  - pro výběr procházím fronty zeshora dolů dokud nenajdu neprádnou frontu
  - pokud se vlákno zablokuje před koncem timeslice tak ho dám výš
  - pokud na konci timeslice pořád počítá, tak níž
  - ⇒ interaktivní procesy skončí v nejvyšší větvi → většinou času jsou zbytečné
- Completely fair scheduler - CFS → Linux scheduler
  - stručně: vybírám proces co ještě nebral moc dlouho a timeslice ~ doba čekání
  - proces si pamatuje svůj virtual runtime = jak celkově dlouho byl ve stavu running
  - jádro má červený - černý strom s ready procesy → indexace podle VRuntime
  - vybírám nejlevější uzel = proces s nejmenším VRuntime
  - timeslice = čas čekání / # procesů ve stromě → to jsou ty ready vlákna

## • Komunikace mezi procesy

→ OS schovává prostředky proceem před ostatními procesy (např. paměť)

→ kooperující procesy si chtějí nějak bezpečně posílat zprávy

• pipes → něco jako síťové sockety ⇒ na jedné straně dávám data a na druhé je beru

• sdílená paměť → dělá se to přes sys-cally

• signály → používá je i kernel

## • Soubory

= data co spolu nějak souvisí / abstraktní proud dat

→ kernel OS nerozumí sémantice těch dat

≈ pohledu OS

## • Identifikace souboru

→ filesystem používá číselné ID

→ soubor má jméno a cestu k němu skrz nějakou stromovitou strukturu

→ takže tam je pro lidi + některé části jména mohou mít speciál význam <sup>např. .git</sup>

(0, 1, 2) = (rodinný, státní, státní)

## • Operace

• open → najde soubor a načte si jeho metadata + vrátí file handle

↳ handle je číslo, které kernel dodá své funkci souborů předloží na ID toho souboru

• close → uloží změny, odstraní ten objekt metadata a vyndá z souborů file handle

• read + write → používáme ten file handle, write se dělá do paměti a až pak na disk

• seek → dělá se abstrakce ukazovátka, ale seek se jenom lineárně posouvá

⇒ každé proces má vlastní tabulku otevřených souborů - je v tom objektu proces

• create, truncate, delete

• flush = zápis změny co jsou načítaný v paměti na disk

→ volá se náhodně a při zavření souboru

• změna atributů

extension napovídá OS/aplikaci co to je

adresa na disku

• Atributy souboru → jméno, velikost, typ, práva, timestamps, ve kterém sektoru běží

• Buffering → předčítání souborů, protože to je slow

↳ OS: načítané sektory se načítají dvakrát

↳ C, C++: žádná podpora si to tedy sama řeší

↳ sekvenční čtení → OS by sektory řeší dopředu = read ahead

• Sync I/O: read() zoblokuje vládku, přečte, vrátí → vládku ready

• Async I/O: read() zahájí operaci ale hned se z toho syscallu vrátí

⇒ vládku běží dál a OS současně čte z toho souboru

→ pokud to nepřičte víc, tak se to vládku dočasně zoblokuje

## • Adresář = sada souborů

- většinou reprezentovaný jako soubor s nějakým special system
- hlavně je to user-friendly a pomáhá to při hledání open()
- někdy si pamatuje nějaké atributy souborů v něm → podle filesystemu
- je nějaká hierarchie a rootem
- operace: create, delete, rename file/subdir, search for name, list members

## • Uložiště souborů

- na disku (externí paměť), RAM, v síti → musí tam být nějaký filesystem
- virtuální soubory → OS pomocí nich poskytuje nějaké funkce navíc /dev/mtd

## • File links

- link (hard link) - více položek v FS ukazují na stejná data → stejné abs. file ID
  - většina operací je transparentních, ale občas je to problematické
- symlink (soft link) - textáček jehož obsah je cesta k jinému souboru
  - je potřeba explicitně říct "follow symlinks"

## • File system = datová struktura v uložisti

- řeší jak a kde jsou data uložena + poskytuje abstrakci souborů pro OS
- musí umět:

1) překládat jména souborů na file ID

→ rozkouskuje tu cestu a rekurzivně se volá na podadresáře ⇒ slow

2) pamatovat si lokaci dat souborů = sekvence bloků

3) management volných bloků

→ spoják, bitmapa, ...

blok = několik souvislých sektorů  
4kB = 8 · 512B

⇒ prázdný sektor soubor zabírá 1 blok = 4kB

lokální FS → na disku → FAT, NTFS, ext2/3/4, XFS, BTRFS

síťový FS → protokol pro přístup k souborům přes síť → NFS, Samba

} někdy je lepší rozdělení

- disk je možné rozdělit do partitions - shodně tam můžou být různé FS
- ⇒ na začátku disku je partition table → velikosti oddílů + FS tam

## • RAID = Redundant Array of Inexpensive Disks

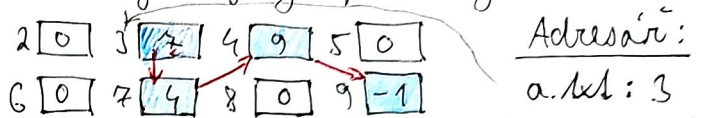
- způsob jak spojit více disků do jednoho → HW nebo to dělá OS
- pro velká data je potřeba (srovnávkou) diskové pole
- je to rychlé a spolehlivé

→ HDD disky

FAT = File Allocation Table

- hodně simple, z dob MS-DOSu
- je tam struktura FAT co se stara o volné bloky a pozici souborů
- na tom disku jsou 2 FATky (kopie) ⇒ záloha když počítač spadne při zápisu
- adresář je speciální soubor
  - v něm je sekvence položek fixní velikosti s atributy
  - pro každou položku si pamatuje číslo prvního bloku ⇒ ve FATce je spoják
- FATka je vlastně pole indexované čísly bloků a  $FAT[\text{blok } n] = \text{blok } n+1$
- ⇒ to pole začíná indexem 2 (0 a 1 mají nějaký speciál význam)

0 = prázdný blok  
-1 = poslední blok souboru



→ na disku je:

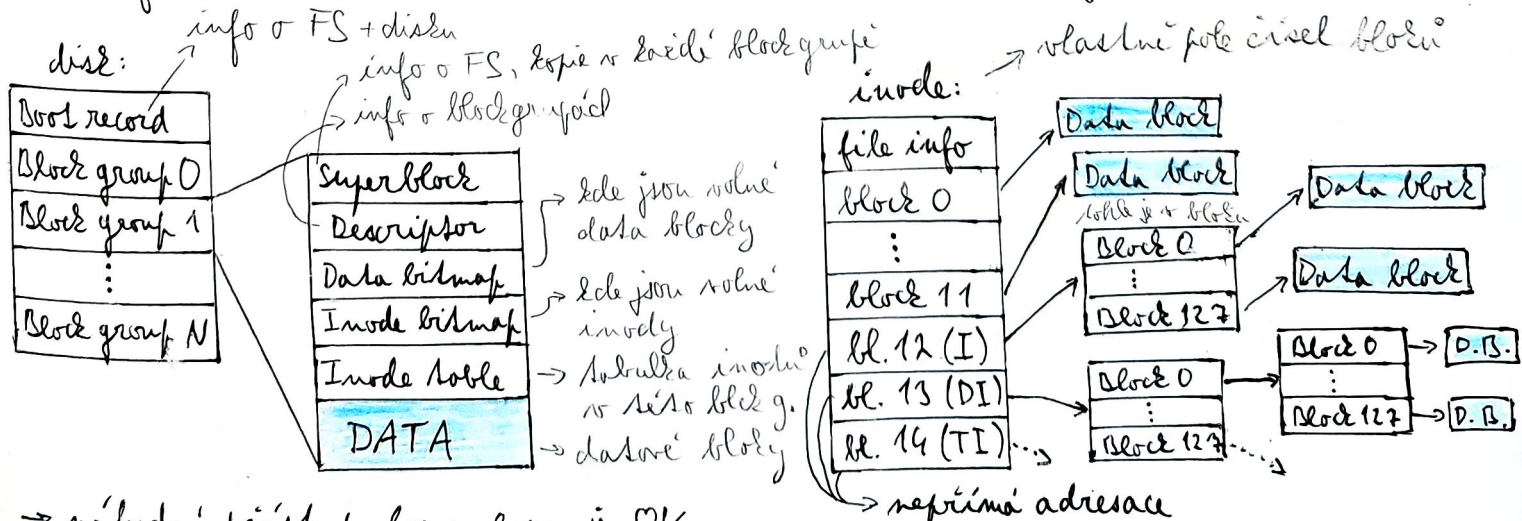


- pro případ, že by se z toho disku bootovalo + info o FATkách a pointer na root a data
- FATka má fixní velikost (je to jen pole) ⇒ umím spočítat jak velká ta fatka bude a kolik bloků zabere
- root má fixní velikost ⇒ omezený # položek
- přístup k datům je simple protože všechno nad tím má fixní velikost

- ⊖ Open je slow protože při probíhávání adresáře ho musím lineárně projít
- ⊖ ten spoják je jednosměrný ⇒ neumím pořádně seek dozadu

ext2 - původně linux a toky celkem simple

- ext3 přidal journals aby se zlepšila disaster recovery
- ext4 umožňuje ukládat větší individuální soubory
- používá index nodes = inode → 1 inode reprezentuje 1 soubor/adresář
- adresář je zase speciální soubor ve kterém je sekvence položek s fixní strukturou → 1 položka = inode number, file name

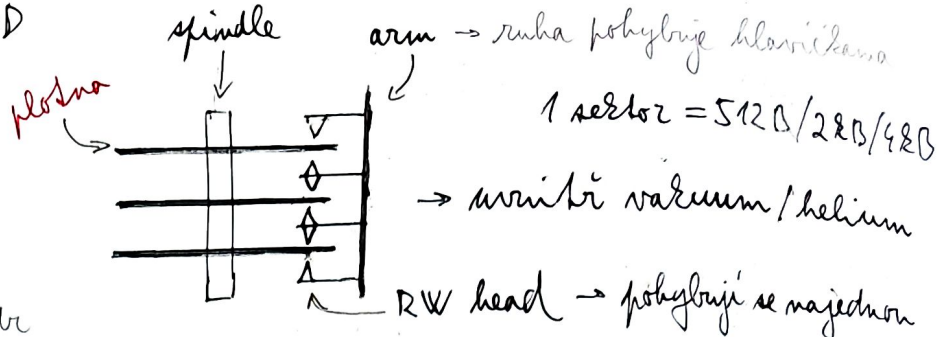
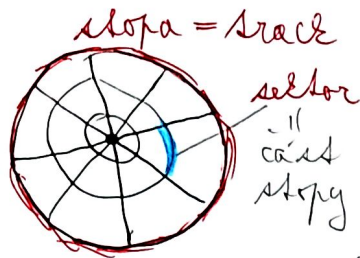


→ náhodný přístup do souborů je OK

→ max. velikost souboru je  $(12 + 128 + 128^2 + 128^3) \cdot 4 \text{ kB} \approx 8 \text{ GB}$



## • Mechanika hard disku - HDD



1 sektor = 512B / 2Kb / 4Kb

→ uvnitř vakuum / helium

cluster = stejná stopa na všech plošnách

block = stejný sektor na všech plošnách

→ rychle bloky pro uživatel FS

• flying height = vzdálenost hlavy od plošny ~ 5 nm

• rotační rychlost = 5400, 7200, 10K, 15K rpm

→ rychle  
→ rychlejší to rychle

→ bottleneck je mechanický pohyb hlavy

→ složitá indexace sektorů ⇒ používá se LBA = Linear Block Addressing

## • Disk scheduling

→ musíme rozhodnout pořadí vykonávání RW requestů

→ dřív to řešil OS, dnes to dělá disk

access time = seek time + čas na rotaci + čas na transfer dat  
→ pohyb hlavy ~ ms  
→ zanedbatelné

→ cíl je minimalizovat access time pro více I/O requestů

• FCFS = First Come First Served → ✓ pro malou záťaž

• SSTF = Shortest Seek Time First → starvation = vzdálené requesty se neplní

• SCAN = algoritmus výstahu → ✓ pro těžkou záťaž

• CSCAN = Circular scan → dojde úplně nahoru a pak se vrátí do přícení a cestou dolů nenabírá pasážery → repeat

• LOOK / CLOOK - jako SCAN / CSCAN ale nenavštíví konce disku ~ pak se vrátí

• FSCAN - 2 fronty, algoritmus pracuje s 1. a nové requesty dává do 2.

## • Solid State Disk - SSD

→ bez pohyblivých částí, rychle se opotřebovávají → hlavně zápis a mazání

→ organizace do bloků rozdělených na stránky → 1 page ~ 2-16KB

→ read / write je pro stránky → 1 block ~ 128-256 pages

→ erase je pro bloky → ty stránky se označí jako invalid

⇒ když se rozbitá 1 stránka, tak přijde o celý blok

→ dělá se GC s data consolidation = odstraní se invalid bloky

→ pom special FS co to tak rychle mění & validní stránky se seskupí na straně

## • Virtuální paměť

→ používají ji i instrukce i kernel OS

→ procesy pracují s virtuálními adresami do virtuálního a.p. = VAS

→ operační paměť má fyzický a.p. = PAS, fyzická adresa = 1 číslo

→ překládá se to hardwarově <sup>→ moduluje to OS</sup> v CPU pomocí MMU = Memory Management Unit

↳ je to zobrazení VAS → PAD, ale to mapování nemusí existovat

↳ pokud neexistuje, tak to MMU pozná a vyvolá výjimku typu fault

→ hlavní mechanismy - segmentace (outdated) a stránkování

→ procesy pracují s adresami (nevi, že V) a při každém přístupu do RAM → převod <sup>HW</sup>

→ proč to děláme?

→ dnes už to není moc potřeba

• hodí se větší adresový prostor (VAS může být větší než PAS)

→ přebytečnou virtuální paměť může OS odložit na disk

• bezpečnost - oddělení adresové prostory jednotlivých procesů

⇒ procesy si nemohou sáhnout do paměti

- je možné logicky oddělit segmenty VAS → read only, executable...

## • Segmentace

→ koncepty:

• VAP je rozdělený na logické segmenty - stejně jako paměť běžícího procesu

→ segmenty jsou nějak očíslovány - každý proces má vlastní číslování (čísly se od 1)

• virtuální adresa má 2 části → číslo segmentu, offset v rámci segmentu

→ by instrukce s kterým musí umět pracovat

• operační systém spravuje segmentační tabulku - pro každý proces

→ je to pole indexované segmenty obsahující deskriptory těchto segmentů

→ v deskriptoru je básová fyzická adresa toho segmentu + délka + atributy

→ OS jenom musí zajistit aby když běží nějaký proces byl někde

v registru pointer na tu segmentační tabulku ⇒ CPU si zde běží

⇒  $PA = ST[\text{Segment num.}] + \text{offset}$  → pokud offset  $\geq$  délka segmentu ⇒

→ pokud je číslo segmentu out of range ⇒ vyvolá se výjimka segment fault

→ taky se kontroluje jestli třeba nedělám write do read only segmentu, ...

→ co když dojde prostor ve fyzické paměti?

→ nějaký segment se celý přesune na disk → v deskriptoru je present bit

⇒ pokud přistupuju do segmentu co není present ⇒ segment fault

→ OS ten segment načte z disku a jde se znovu výpadek segmentu

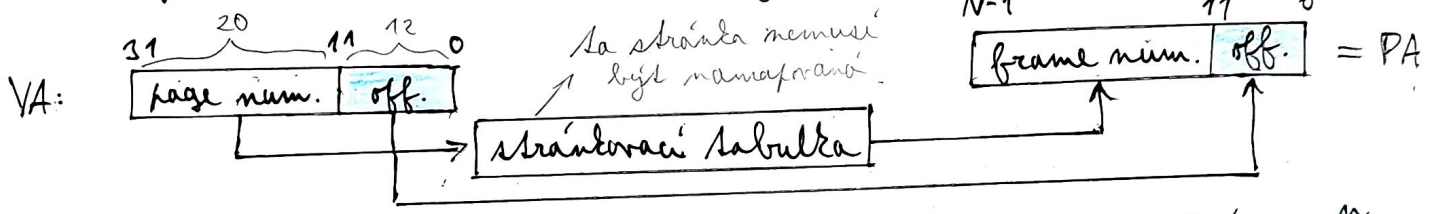
→ spolu se společným managementem fyzické paměti → největší fragmentace  
je tohle hlavní důvod proč už se segmentace nepoužívá

Stránkování

- VAS je rozdělený na stejné velké části (stránky), velikosti  $2^m$ 
  - ↳ velikost VAS je tedy  $2^N$  → pro 32/64-bit CPU  $N=32/64$  - většinou
- PAS je rozdělený na stejné velké rámece = frame
  - ↳ rámece mají stejnou velikost jako stránky → většinou  $4kB = 2^{12} B$
- virtuální adresa je 1 číslo → rozdíl od segmentace

stránkovací tabulka - pole v paměti pro 1 proces

- indexovaná číslem stránky → obsahuje číslo rámce + atributy
- v těch atributech je zase příznak  $P = \text{present}$ , což říká jestli mapování existuje
- pokud mapování neexistuje / jiný problém ⇒ page fault = výpočet stránky
- pro 32 bit VAP a 4kB stránky:

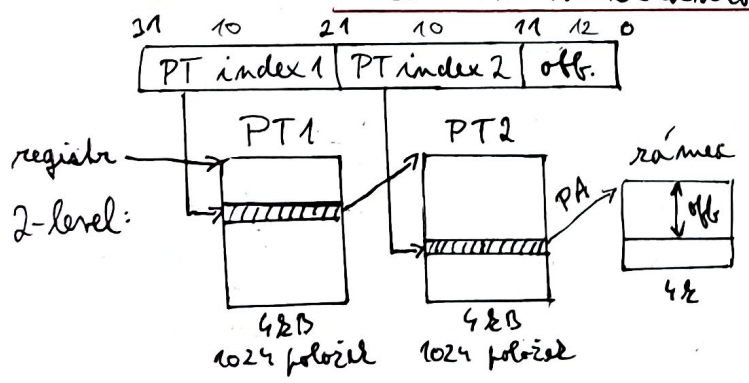


- takhle funguje protože stránky a rámece jsou stejné velikosti a  $2^m$
- VAS a PAS můžou být různé velikosti
- 1 logický segment v tom procesu je rozdělený do více souvislých stránek ale ty rámece nemusí být ve fyzické paměti souvisle za sebou ⇒ když dojde paměť tak nemusím vyhledovat celý segment ale jen 1 stránku

Problémy stránkování

1, velikost - pro 1-úrovňovou tabulku ↑ to je to pole

- ↳ pro 32b VAS je  $2^{20}$  4kB stránek ⇒  $2^{20} = 1M$  položek
- pokud 1 položka = 4B, tak si musím pro 1 proces pamatovat 4MB tabulku
- ⇒ řešení: víceúrovňová stránkovací tabulka → ∴ nepotřebuju celý VAP



- PT1 jsou fyzické adresy PT2
- PT2 jsou už fyzické adresy rámců
- ta tabulka má 4kB = 1 stránka
- PT1 musím mít vždy v paměti
- ty PT2 mapuju podle potřeby

- rámcový PT1 musí mít pořád present bit ⇒ může být page fault na PT2
- 1 tabulka 2. úrovně stačí na adresování  $1024 \cdot 4kB = 4MB$  paměti

## 2, rychlost

- přehled na PA vyžaduje víc přístupů do paměti ⇒ to je mega slow
- ⇒ řešíme to do TLB = Translation Lookaside Buffer
- TLB využívá asociativní paměť ⇒ slovník (stránka, rámec) → víc bylo učen
- využíváme toho, že procesy často přistupují ke stejným stránkám

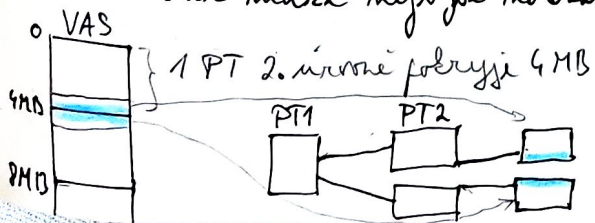
## • Algoritmus pro převod adresy u stránkování → HW

1. rozděl VA na číslo stránky a offset
2. checkni jestli to má v TLB → pokud to tam je tak končíme
3. projdi stránkovací tabulku
  - rozděl číslo stránky na tolik částí kolik je úrovní
  - jdi do PT1 → pokud tam je adresa PT2 jdu dál → pomocí přesné bitů
  - pokud tam adresa další tabulky není ⇒ page fault
  - z poslední úrovně vyvedni číslo rámce a ulož ho do TLB
4. aktualizuj přenosy A = Accessed a D = Dirty ve stránkovací tabulce a TLB
  - A ~ stránku někdo přečetl, D ~ do stránky někdo napsal
  - musím je změnit i když jsem se k té stránce dostal přes TLB
5. získaj fyzickou adresu z čísla rámce a offsetu

## • Jak se řeší výpadek stránky?

- OS interrupt handler musí zjistit kde se stala chyba u toho stránkování
- v který proces běží ⇒ v objektu procesu je pointer na tu stránkovací tabulku
- možné důvody:
  - neoprávněný zápis, čtení ze systémové paměti, ...
  - pokus o přístup do ještě nenaalokované virtuální paměti → proces sám má svoje segmenty a OS musí alokovat tu virtuální paměť
  - vše OK ale neexistuje mapování
- jak vyrobit mapování?
  - OS najde volný rámec nebo vybere oběť algoritmem na výměnu stránek
  - pokud je oběť dirty (změnila se) tak ten rámec uložíme na disk
  - musím zrušit mapování oběti v TLB → TLB nemůžeme měnit, to udělá HW
  - na volný rámec nahrajeme svůj obsah a opravíme stránkovací tabulku
- vrátím se zpátky z handlerem a zkouším tu instrukci znovu

## • Úloha: kolik může nejvíce nastat výpadek stránek když chci načíst 2B?



- by 2B můžou být na rozhraní 2 s.t. 2. úrovně
- můžou nastat 4 výpady!
- je to nejenom na hraně stránek ale i tabulek!

## • Algoritmy na výměnu stránek

→ využívají se jakékoli situacii kde je potřeba vybrat oběť aby uvolnila místo  
⇒ vybírání rámců, cache, TLB

## • Optimal page algorithm

→ vybere stránku co bude nejdříve nepoužita ⇒ nejmenší page-fault rate  
→ jen teoretický, je snaha se tomu přiblížit

## • Hodiny (clock)

→ rámců se organizují do kruhu a dáme tam ručičku - ukazuje na rámec  
→ využívá Access bit, který MMU nastaví na 1 když na tu stránku sáhne  
→ rámec vybíráme tak, se postupují ručičkou dookola:

1. if  $A \text{ rámec} \neq 0$ :  $A \leftarrow 0$  a jdu dál ← rámec využij
2. if  $A \text{ rámec} = 0$ : vyberu tenhle rámec ← dlouho nebyl použit

## • NRU = Not Recently Used

→ příznaky A se mění pravidelně - třeba 1 za minutu  
→ rámce rozdělí do 4 tříd používání podle příznaků A a D

A	D	Třída	
0	0	0	→ nepoužité + nezminované
0	1	1	→ nepoužité ale musím to reset na čist
1	0	2	→ někdo to používá
1	1	3	

⇒ vyberu náhodný rámec z nejvyšší neprázdné třídy

## • LRU = Least Recently Used → takhle se reálně používá

→ používám minulost pro předvídání budoucnosti

⇒ vybere stránku na kterou nejdelší dobu nikdo nesaahl

→ existují HW implementace

→ SW implementace: → spojít nebo odobit

• zásobník - když na něco sáhnu tak to dám naboru a oběť je na chvilu

↳ takhle je moc pomalé ⇒ apertimace NFU

## • NFU = Not Frequently Used - apertimace LRU

→ každý rámec má počítadlo → někde během, do A.A. by se nevědělo

→ pravidelně měníme A a pokud  $A = 1$  tak počítadlo inkrementujeme

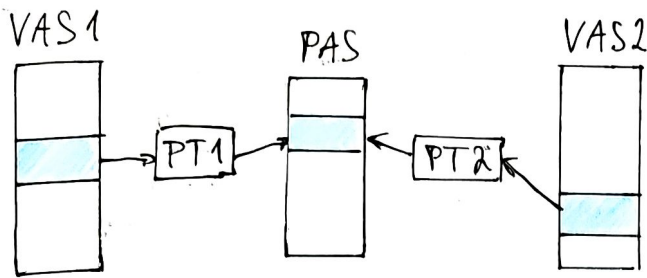
⇒ vybere se rámec s nejmenším počítadlem → aby se hned nevyhodily

→ novým stránkám musím na začátek dát nějaké skóre

→ staré framgy by se nevybraly ⇒ aging: periodicky  $count \leftarrow count/2$

## • sdílená paměť

- více procesů spolu sdílí část virtuální paměti - způsob komunikace
- někdo udělá sys-call "vyrob sdílenou paměť" a další procesy se připojí

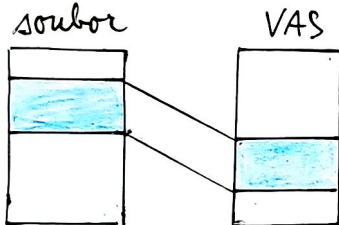


- každý proces má svoji A.A. a ta sdílená paměť je v VAP jinde ale ty rámce jsou sdílené
- ⇒ ve sdílené paměti není možné používat ukazatele → jenom offsety

→ při 2-level stránkování lze mít společnou tabulku 2. úrovně

## • Memory-mapped files = paměťově namapované soubory

- normálně udělám open a pak dělám read a write
- ↳ ten soubor se mi nakopíruje někde do paměti ale pořád musím dělat R/W



- soubor si namapuju do virtuální paměti
- ⇒ do toho souboru můžu přistupovat operacemi procesoru load a store
- ⇒ normálně nastavuju A a D bity

→ když se ta stránka vybere jako oběť, tak pokud  $D \neq 0$ , tak uložíme

### → problémy:

- appendování na konec souboru → jsou tam nenamapované stránky
- soubor může být větší než VAS

## • Virtual machine

- přibližný výraz, teď se budeme bavit o virtualizaci ne Java/C# VM
- v rámci OS umožníme vytvořit chráněné prostředí, které umožňuje spustit další OS
- ↳ je tam iluze toho, že to běží na reálném HW
- v případě problému se o to musí postarat ten host OS - ale většinou se to neděje
- normální instrukce běží na CPU, ale special kernel instrukce (nastavování stránkování) emuluje ten host OS

⊕ izolace VM od zbytku systému

⊕ enkapsulace → VM ~ soubory ⇒ můžeme dělat snapshoty kdyžby to spadlo

⊕ kompatibilita → snadno se přenáší na jiný HW - VM poskytnou virtuální HW

↳ virtuální HW → virtuálně můžeme mít nastavenou různou RAM, počet jader

## • Kontejnery - virtualizace na úrovni OS

- zjednodušená virtualizace → používá se reálný HW ne virtuální
- OS kernel musí umožňovat existenci více oddělených user spaceů
- ↳ na linuxu si v kontejneru napustím windows

## Paralelní programování a synchronizace

- Paralelní počítání = používá se více jader a instrukce se provádějí současně
- Concurrent počítání = multitasking na 1 jádru

### • Problémy paralelního počítání

#### • Race condition

→ když více vláken současně mění stejná data

⇒ výsledek se může lišit podle časování / schedulingu OS

⇒ není to deterministické

→ příklad se spojkem: LL l; ↪ A. next = head;

1. l. add(A)

2. l. add(B)

head = A;

↙ před assign disturbance interrupt

a = 10.

1. a = a + 1

2. a = a + 2

⇒ a ∈ {11, 12, 13}

#### • Kritická sekce

→ identifikují kritickou sekci kódu kde může vzniknout race condition

⇒ řešení = mutual exclusion ⇒ v k.s. - 1 vláknem at a time

#### • Synchronizace - jak udělat mutual exclusion? = zachování integrity dat

→ když se kritická sekce zamýká nebo se vláknem nějak řídí

→ realizace pomocí synchronizačních primitiv

• aktivní - ztrou čas procesoru (aktivní čekání)

• pasivní (blokující) - kernel to vláknem zablokuje dokud není přístup k resursu

↳ pasivní se nemusí vždy vyplatit ∵ sys-call na blokování je drahý

⇒ pokud je race condition vzácná, což je aktivní lepší

→ aktivní potřebují HW podporu

→ používají atomické instrukce → Test-and-set (TSL), Compare-and-swap (CAS)

↳ instrukce co CPU garantuje udělá celou a bez přerušování

int TSL: nastaví rámeček na novou hodnotu a vrátí starou hodnotu

bool CAS: porovná rámeček s danou hodnotou a pokud jsou stejné, což ho nastaví na novou hodnotu, vrátí true/false

## • Spin-lock - aktivní

→ když zámeček = 0 tak je volno

⇒ všichni se ho pomocí CAS snaží nastavit na 1

⇒ ať se vládnou dokončí kritickou sekci tak zámeček nastaví na 0

→ vhodné pro krátké čekání

## • Semafor - blokující

→ počítadlo a fronta čekajících vládek

→ atomické operace UP a DOWN ale je to sys-call → provádí se to třeba pomocí spin-locku

→ počítadlo se na začátku nastaví na # vládek co chceme najednou pustit do kritické sekce

→ když vládnou chce do kritické sekce, tak zavolá DOWN

→ když se jí odchází tak zavolá UP

down: ← je volné místo

```
if count > 0: --count;
```

```
else: ← není volno → jdu do fronty  
queue.push(tohle vládnou)  
vládnou.block()
```

up:

```
++count;
```

```
if !queue.empty():
```

```
  v ← queue.pop()
```

```
  v.unblock()
```

```
--count;
```

→ volným místo

← pokud někdo čeká

← odblokuje ho

← jde kam nikdo místo má

## • Mutex - implementace mutual exclusion

→ atomické operace LOCK a UNLOCK

⇒ binární semafor, kde na začátku count = 1

## • Barriéra

→ když vládnou dorazí k bariéře, tak se zablokuje

⇒ čeká se ať u bariéry bude vícetý počet vládek a pak se všichni odblokuje

→ je to dobré na synchronizaci vládek - tím, že všechny jsou kady

## • Monitor - v programovacích jazycích

→ metody objektu nějaké třídy mají zámeček jako Mutex

⇒ na začátku metody je lock a na konci unlock

⇒ 2 vládnou nesahají současně na stejný objekt

## • Deadlock - vládnou 1 čeká ať vládnou 2 dokončí operaci a to čeká ať 1 dokončí operaci

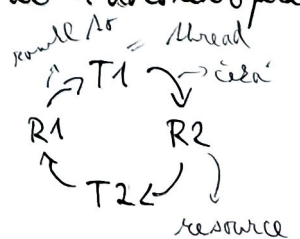
Coffmanovy podmínky:

1. mutual exclusion: prostředky nelze sdílet

2. hold and wait: vládnou drží prostředek & chce další

3. no preemption: prostředek nemůžeme vládnou jin tak sebrat

4. circular wait: v modelovacím grafu je cyklus





# • Klasické synchronizační problémy

## 1, Producer - consumer

→ sled s omezenou kapacitou, výrobou zboží a konzumenty

- co když přijede výrobce a konzument zároveň? → musí se sdílet mutex
- co když se sled naplní? → blokuje výrobce - čeká na konzumenta
- co když je sled prázdný? → blokuje konzumenta - čeká na výrobce

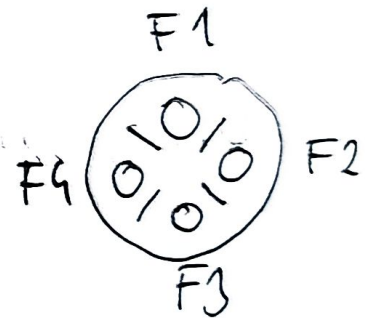
## 2, večerní filozofové

→ N filozofů v kruhu - mají po pravé ruce vidličku a před sebou talíř

→ filozofové buď přemýšlejí nebo jí

⇒ jak se sdílet aby každý někdy jedl

↳ co když dostanou hlad ve stejný moment?



• popodmí pravou a čekaj na levou ⇒ deadlock

• když vlevo není, se vrátím na pravou ⇒ starvation

↳ pořád lůžka, nejsou zabíráni, ale nenají se