



Puvodni google doc:

<https://docs.google.com/document/d/1Svj0KUPPj1kCWNNPgPTKkiNgOe66P3AwvptjFh8OJfM/edit#>

Pokud nekdo chce ziskat high-level prehled o computer science vctne historie, vyvoje pocitacu a zakladnich logickych obvodu a jak kde zapadaji ruzne predmety co mame za sebou, doporucuji crashcourse computer science (az se budete nudit po statnicich, ted neprokrastinujte!):

<https://www.youtube.com/watch?v=tplctyqH29Q&list=PLH2l6uzC4UEW0s7-KewFLBC1D0l6XRfye>



## Číselné soustavy

**big endian** = číslo uloženo po **bajtech** od nejvyššího bajtu po nejnižší (zleva doprava)

**little endian** = je to naopak -- používá Intel, obecně jsem nabyl dojmu že je více common

**nibble** = 4 bity - odpovídá jedné cifře hexadecimální soustavy (0 - 9, A - F (15))

**desítková** -> **dvojková** soustava:

- delíme dvojkou, píšeme si výsledky modula, pak je prepíšeme od konce a máme binární číslo :)
- nebo jdeme od nejvyšších mocnin -- najdeme nejvyšší mocninu obsazenou v tom čísle, odečteme ji od čísla zapíšeme 1,
  - pak vezmeme nižší mocninu : menší než číslo? napíšeme 1 a odečteme : napíšeme 0 a jdeme na nižší mocniny

binární -> desítková:

- easy

**Zaporná celá čísla** v počítači (nevímejte počet bitů, není důležité teď :D ):

- **dvojkový doplněk** - záporná čísla reprezentována od 1111111... - založeno na tom že když budeme pod sebou odečítat 0 binární od binární jedničky, bude tam přenos a dostaneme právě 11111...
- výhodné pro počítání
- máme pak nejvyšší řád rezervována pro záporná čísla, tedy každé číslo začínající 1 je záporné
- to odpovídá rozdělení rozsahu na dvě poloviny - nejvyšší kladné číslo je pak 01111 -> když přičteme něco dostaneme se do záporných

- tedy když reprezentujeme čísla pomocí  $n$  bitů, dostaneme rozsah čísel od  $\langle -2^{(n-1)}, 2^{(n-1)} - 1 \rangle$  -- kladných čísel o jedno "méně" protože nula

### Scitání čísel ve dvojkové soustavě:

- aktuální bity:  $A_i, B_i$
- **jednoduchá scitacka pro jeden bit:**
  - vezme přenos z nižšího řádu  $C_i$  a aktuální bity scítaných čísel
  - vyčísle součet bitů  $S_i$  a přenos do vyššího řádu  $C_{i+1}$
- **ripple carry adder**
  - zřetězením jednoduchých scítacek můžeme vytvořit scítacku libovolně dlouhých čísel
  - jmenuje se
  - pomale -- jednotlivé scítacky na sebe čekají
- **carry look ahead adder** <https://www.geeksforgeeks.org/carry-look-ahead-adder/>



- založena na analýze přenosu do vyšších řádů
- pro hodně bitů by byla komplikovaná HW struktura, proto se pak taky nějak řeší ty pro méně bitů
- situace s přenosem (carry):
  - generování přenosu
    - $1 + 1$
    - bool variable:  $G_i = A_i \text{ and } B_i$  ( $= A_i * B_i$ )
    - přenos vznikne bez ohledu na přenos nižšího řádu, tj. ovlivňuje výsledek ale nikoliv další přenos
  - propagace přenosu
    - $1 + 0, 0 + 1$
    - bool variable:  $P_i = A_i \text{ xor } B_i$
    - přenos ovlivní výsledek a jede dál
  - výsledek  $S_i = P_i \text{ xor } C_i$ 
    - pokud propagujeme a zároveň dostáváme přenos, máme 0
    - pokud propagujeme ale nedostáváme přenos máme 1
    - pokud nepropagujeme ( $0 + 0$ ) a dostáváme přenos máme 1
    - pokud nepropagujeme a nedostáváme přenos máme 0
  - přenos  $C_{i+1} = G_i + P_i C_i$  (+ je boolean addition tedy or)
    - buď generujeme, nebo dostaneme přenos a navíc přenesíme dál
- pak pro 4 bitovou look-ahead scítacku máme následující vztahy:

The carry output Boolean function of each stage in a 4 stage carry look-ahead adder can be expressed as

$$\begin{aligned} C_1 &= G_0 + P_0 C_{in} \\ C_2 &= G_1 + P_1 C_1 = G_1 + P_1 G_0 + P_1 P_0 C_{in} \\ C_3 &= G_2 + P_2 C_2 = G_2 + P_2 G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 C_{in} \\ C_4 &= G_3 + P_3 C_3 = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0 + P_3 P_2 P_1 P_0 C_{in} \end{aligned}$$

From the above Boolean equations we can observe that  $C_4$  does not have to wait for  $C_3$  and  $C_2$  to propagate but actually  $C_4$  is propagated at the same time as  $C_3$  and  $C_2$ . Since the Boolean expression for each carry output is the sum of products so these can be implemented with one level of AND gates followed by an OR gate.

#### - stromova scitacka

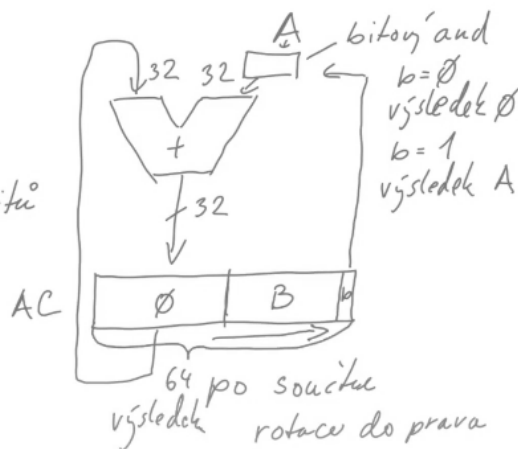
- slozitosť  $2\log n$  -- nebudu zabredavat :D -- ale v podstate se ty zavislosti nejak skladaji do stromu nebo tak neco

#### Nasobeni celych cisel

- podobne jako nasobeni pod sebou
- pokud nasobime A a B -- 32 bitu
- nasobime do akumulatoru který ma 64 bitu
- inicializace: akumulator jsou nuly | B
  - podle posledního bitu (tedy na začátku posledního bitu B) provedeme and s ackem (to nám jej buď vynuluje nebo zkopíruje -- v násobení odpovídá  $B_i \cdot A$ )
  - přičteme výsledek and k levé části akumulátoru a zapiseme
  - provedeme rotaci akumulátoru doprava
  - opakujeme
- po 32 rotacích máme výsledek uložený v akumulátoru
- je to pomalý alg, ve skutečnosti se to přitom dělá jinak :D
- obrázek níže -- je z toho videa převlehl jsem do černobílé aby slo tisknout :)

Násobení čísel

A - 32 bitů  
B - 32 bitů  
 $A \cdot B = AC$  - 64 bitů



## Realna (desetinna) cisla

Realne cislo s konecnym rozvojem se neda vzdy zapsat binarnim s konecnym rozvojem ale naopak se da vzdy!

### Standart IEEE 754 -- 32 bitovy float

- 1 bit **znamenko**, 8 bitu **exponent** a 23 bitu **mantissy** (hodnoty cisla)
- normalizovana cisla:
  - exponent je zvetsen o 127
  - mantissa obsahuje skrytou 1 (protoze vetsina cisel jsou 1.neco)
- nestandardizovana cisla:
  - exponent 00000:
    - cisla mensi nez  $e^{-127}$  a 0
  - exponent 255
    - nekonecno -- mantissa musi byt 0
    - jinak je to NaN
- prevod do IEEE 754:
  - exponent:
    - pokud cislo je mensi nez 1:
      - nasobime dvema dokud nedostaneme cislo vetsi nez 1 -> pak exponent je  $e^{-(\text{pocet nasobeni})}$
    - jinak:
      - to same akorat delime :)
    - nebo udelame  $\text{exp} = \text{floor}(\log_2(\text{cislo}))$  :D
  - mantissa:
    - porovnavame mocniny dvojky od  $-1 \dots -n$  s cislem co nam zbyde po predchozim alg.
    - pokud je cislo mensi -> 0
    - pokud je vetsi nez mocnina -> zapiseme 1 a od cisla odedeme mocninu

### Scitani a odcitani:

- prevedeme na stejny exponent a secteme
- kdyztak musime upravit exponent a tedy i mantissu (aby carka byla na spravnem miste)

## Základní architektura pocitace

CPU je spojen sbernici s RAM (= random access memory = muzeme cist odkudkoliv) sbernice:

- adresni cast - s kterou casti pameti pracujeme
- datovou -- prenasi data
- ridici - co delame - cteme nebo zapisujeme

**Harvard architektura RAM** - rozdelena na instrukce a data

**Von Neumann** - nerozlišuje instrukce a data v pameti

- dnes rozsirena

## **CPU**

registry = rychla pamet v CPU

dulezity registr: PC (program counter) nebo se mu taky rika IP (instruction pointer) = ukazuje kde v programu se nachazime (ukazuje kde v pameti jsme)

## **ALU (Arithmetic logic unit)**

- dela zakladni aritmeticke a logicke operace - kombinace ruznych obvodu a gatu
- jeden ze zakladnich stavebnich bloku CPU

## **Vykonavani instrukce:**

1. fetch - nacte instrukci podle IP do registru, IP++
2. decode - typ instrukce a parametry
3. execute - vykonani -- vysledek se uklada do pameti nebo do registru

=> pro pipelining se pak da jeste rozdelit na vice fazi, zvlaste execute na ALU operaci a zapis do pameti

Instrukce se vykonavaji dokud neprijde interrupt (preruseni):

1. IP se ulozi do nejakeho registru pro navrat z preruseni
2. IP se nastavi podle preruseni

Take muze dojit k vyjimce pri vykonavani (deleni nulou) ci nacisti instrukce (spatna cast pameti)

Procesory podle typu instrukci:

## **CISC**

- intel
- ruzne dlouhe instrukce

## **RISC**

- MIPS, ARM, RISC-V
- 32 bitu = jedna instrukce
- vyhoda - instrukce stejne dlouhe a IP tedy trivialni implementace, povoleny jen skoky po 4 bajtech
- nevychoda: nevejde se tam 64 bitove cislo -- musi se provest dvema instrukcemi

## MIPS instrukce:

Instrukce	Syntax	Operace	Význam
Add	add \$d, \$s, \$t	$\$d = \$s + \$t$ ;	Add: Sečte dva registry $\$s + \$t$ a výsledek uloží do registru $\$d$
Addi	addi \$t, \$s, C	$\$t = \$s + C$ ;	Add immediate: Sečte hodnotu v $\$s$ a znaménkově rozšířenou přímou hodnotu, a výsledek uloží do $\$t$
Sub	sub \$d, \$s, \$t	$\$d = \$s - \$t$	Subtract: Odečte znaménkově obsah registru $\$t$ od $\$s$ a výsledek uloží do $\$d$
Bne	bne \$s, \$t, offset	if $\$s \neq \$t$ go to $PC+4+4*\text{offset}$ ; else go to $PC+4$	Branch on not equal: Skáče pokud si registry $\$s$ a $\$t$ nejsou rovny
Beq	beq \$s, \$t, offset	if $\$s == \$t$ go to $PC+4+4*\text{offset}$ ; else go to $PC+4$	Branch on equal: Skáče pokud si registry $\$s$ a $\$t$ jsou rovny
slt	slt \$d, \$s, \$t	$\$d = (\$s < \$t)$	Set on less than: Nastavi registr $\$d$ , pokud platí podmínka $\$s < \$t$
jump	j C	$PC = (PC \wedge 0xf0000000) \vee 4*C$	Jump: Skáče bezpodmíněčně na návěstí C
lw	lw \$t, C(\$s)	$\$t = \text{Memory}[\$s + C]$	Load word: Načte slovo z paměti a uloží jej do registru $\$t$
sw	sw \$t, C(\$s)	$\text{Memory}[\$s + C] = \$t$	Store word: Uloží obsah registru $\$t$ do paměti
lui	lui \$t, C	$\$t = C \ll 16$	Load upper immediate: Uloží předanou přímou hodnotu C do horní části registru. Registr je 32-bitový, C je 16-bitová.
la	la \$at, LabelAddr	lui \$at, LabelAddr[31:16]; ori \$at, \$at, LabelAddr[15:0]	Load Address: 32-bitové návěstí uloží do registru $\$at$ . Jedná se o pseudoinstrukci - tzn. při překladu se rozloží na dílčí instrukce.

## CACHE

Potřebujeme protože procesor jede třeba na 0.3 nsec, ale RAM má průměrné zpoždění 70 - 120 nsec

**L1, L2, L3 cache** -- zleva doprava se zvětšuje jejich paměť (32kb, 256 kb, 8 mb) a zároveň se zmenšuje jejich rychlost (1.2, 4, 12 nsec)

pozn. cena L1 cache je milion korun za GB :D (proto je tak malá)

### Cache:

- skládá se z **bloku dat - 4, 8, 16 slov (slovo = 4 bajty)**
- paměť se načítá z oblasti v paměti o velikosti počet slov \* 4 bajtu, která začíná na adrese dělitelné počtem slov \* 4
- pak zaznam v cache:
  - tag -- první 4 bity, které jsou pro všechna slova stejná
  - blok dat
  - příznakové bity
    - valid 1/0
- cache hit - pokud se v cache nachází tag který se shoduje se začátkem adresy co hledáme, můžeme si z cache vytáhnout dotčený bajt / skupinu bajtů
- miss - není v cache, hledáme v paměti
  - pak něco vyhodíme (podle různých kritérií) a nahradíme

**Plně asociativní cache** = data můžeme dát kam chceme do cache, celý tag se musí porovnat se všemi tagy v cache -- hodně komparátorů a náročná hardware implementace

**jednocestná cache** = data musíme uložit do části cache (množiny) která odpovídá jejich prefixu a tedy i když je jinak prázdná, musíme třeba přepsat užitečnou hodnotu

dvoucestná cache = množiny jsou zdvojeny -- můžeme si vybrat tu prázdnou

vyhoda "cestných" caches oproti plně asociativní je menší počet komparátorů

- když nenajdeme něco v L1 cache hledáme v L2, pak v L3 a pak kdyžtak v RAM
- L3 cache byva společná pro jádra, musí se tedy synchronizovat!

### **Zápis do cache:**

- write through - nezapiseme rovnou do cache ale čekáme na zápis o úroveň výše
- write back - napíšeme data pouze do cache a nastavíme se dirty bit, že jsme data modifikovali => data se pak paralelně s prací procesoru kopírují dál => může nastat problém s konzistencí

## **Stránkování**

k čemu? -> aby každý proces mohl mít svůj vlastní virtuální svět (paměť)

Máme tabulky pro každý proces co značí, kde co je v reálné RAM

-- většinou vícevrstevně.

Více detailů asi v OSY nebo jiných materiálech, nechce se mi to vypisovat :D

## **PIPELINING**

Rozdělíme vykonávání instrukcí na více fází: fetch, decode, execute, memory access (instrukce load/store - čtení a zápis do paměti) a write back (zápis výsledku do registru, load) -> pak můžeme instrukce paralelizovat -> začneme fetch instrukce, pak děláme decode ale mezitím už děláme fetch další instrukce

Může nastat problém se synchronizací - např. stará hodnota v registru - **hazard**

- resimé:
  - pozastavením vykonávání dané instrukce dokud nemá aktuální data
  - forwarding - z mem nebo write back přepošleme rovnou do execute

## **Predikce skoku**

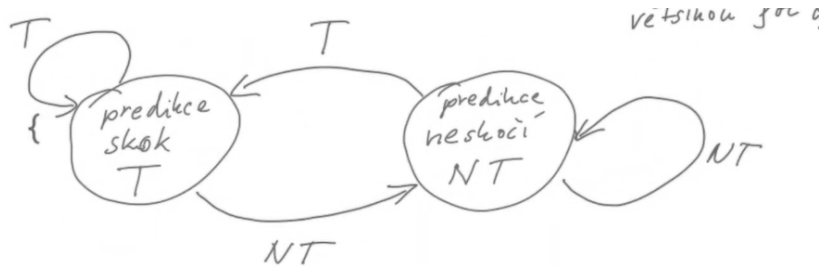
Skok může představovat problém, protože procesory (zvláště ty složitější) si dopředu načítají velké množství instrukcí -> můžeme predikovat, jestli instrukce povede ke skoku nebo ne

T = taken (skocili jsme)

NT = not taken (neskocili jsme)

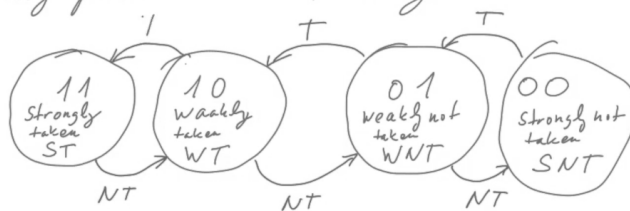
### **Jednabitový prediktor:**

- automat
- když jsme skocili, predikujeme skok, když jsme neskocili predikujeme že neskocíme :D
- nefunguje špatně, protože ve for cyklu mnohokrát za sebou skácame



### Dvoubitový prediktor:

Dvou bitový prediktor  $\Rightarrow$  4 stavy



“dele trva jej presvedcit”

Moderni CPU vyuzivaji pro predikci **perceptron**:

vstupem je registr s historií posledních  $n$  skoku

Prediktoru **muze byt vic a muzeme se rozhodovat který pouzit**

## Pripojeni periferii I/O

pro pripojeni periferiim:

- intel -- instrukce in out (spis zpetna kompaibilita) + **sdileni pameti**
- risc -- jen **sdileni pameti (DMA - direct memory access)**

central address decoder - vi, která cast pameti RAM je pro periferie, zna cely system  
X

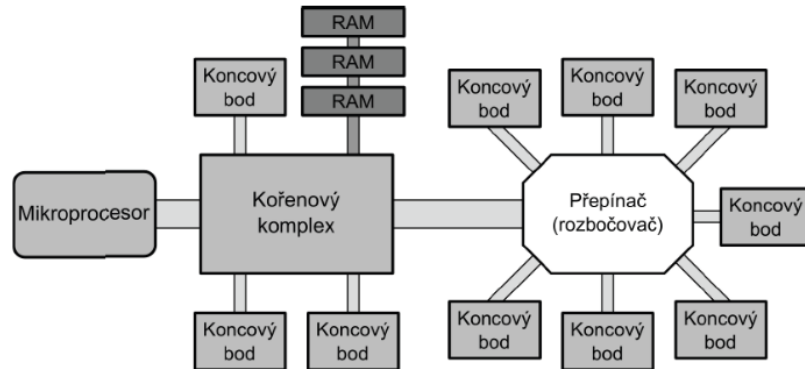
kazda periferie ma svuj decoder - zna svou cast pisecku :D

-> napr LED se rozsviti podle toho co si precte ve svoji casti pameti

- Ruzne sbernice (aka draty k periferiim, menenim napeti vysilaji signal):
  - asynchronni - bez hodin + seriove - prenasime bit po bitu
    - napr SATA -- pripojeni HDD
    - PCI-e (express)
  - paralelni:
    - vice dratu paralelne
    - oproti seriovym umi fungovat v nizsich frekvencich protoze se musi poresit synchronizace dat mezi jednotlivymi draty
    - pr PCI - komunikace pouze jednim smerem (half duplex) - povol - odpoved
      - na jedne sbernici muze byt vice zarizeni s jinou adresou
  - PCI-e (express) - seriova full duplex (vysila i prijima zaroven) - vodici navic

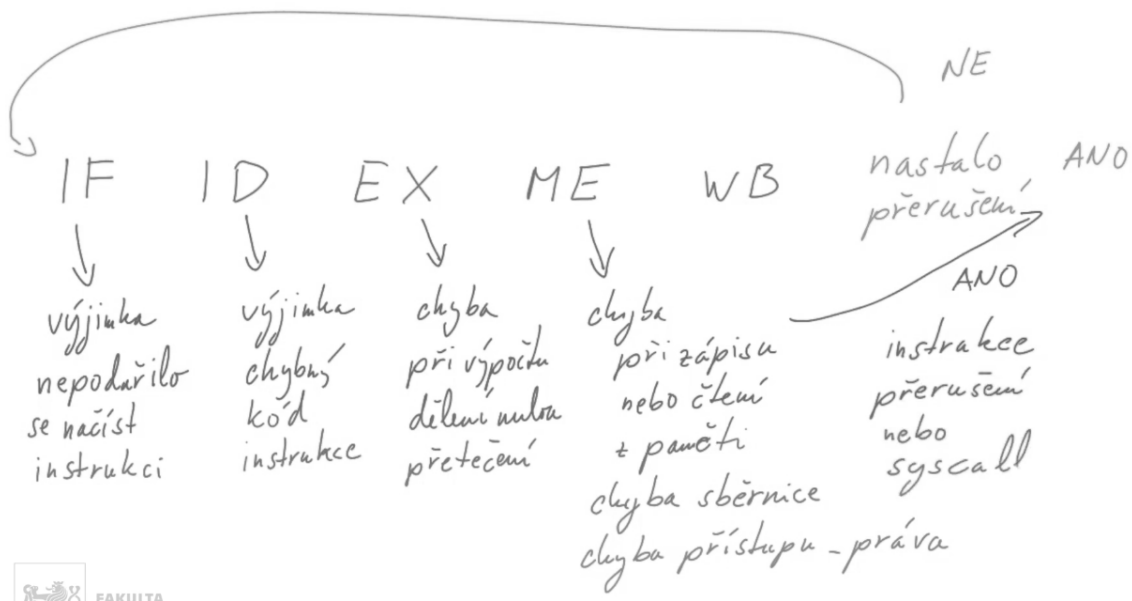


- oproti PCI daleko rychlejší
- point - to - point
  - na CPU je napojen switch který komunikuje s periferiemi
- na obrázku:



## Vyjimky a preruseni

Jake mohou vzniknout:



## Komunikace s periferií:

- CPU udělá preruseni, specifikuje co chce od periferie
- periferie pracuje
- pak da prerusenim vedet CPU ze je done

## Zpracovani preruseni:

- **vnejsi** preruseni
  - při zpracovani nesmi byt naruseny uzivatelske registry, pro program je neviditelne

- asynchronni
- IP (PC) uzivatelskeho programu je zkopirovan do jineho registru - na konci je zpet nastaven
- pri pipelingu je dokoncena pipeline
- **vyjimky**
  - konec programu - nemuzeme dal pokracovat nebylo by korektni
  - chyba operacniho systemu - modra smrt
  - zpracovani je synchronni

## Predavani parametru funkcim v C

Pouziva se **zasobnik** (stack) - jak jinak :D

**implementace zasobniku:**

- push: vezme registr esp, mensi o 4 a pak ulozi pozadovanou hodnotu na adresu na kterou ukazuje esp
- pop: zvetsime esp o 4 a precteme hodnotu byla v pameti kam ukazova esp

**Volani funkce:**

- skocime do funkce pomoci specialni instrukce (zapamatuje se na zasobnik kde jsme byli predtim)
- pak se ve funkci nejprve zvetsi zasobnik o potrebne misto
- parametry pri volani funkce se ukladaji na zasobnik ci u mipsu prvni 4 se ukladaji do registru
- navratova hodnota se da do registru rax (v mipsu aspon)

## Procesory x86

prehled procesoru a jejich historie - asi zbytecne si pamatovat :D

prehled instrukci intelu --taky si nebudu pamatovat

## Preklad C

1. **preprocessing** - textove nahrazeni maker (mohou bt parametrizovane), include se nakopiruji
2. **lexikalni analyza**
  - a. text je rozdelen na lexy -- zakladni stavebni kameny -- cisla, konstanty, symboly
  - b. pravidla jsou vlastne regularni vyrazy -- postupujeme po rade po jednotlivych regularnich výrazech -- od special symbolu po promenne atd
3. **syntakticka analyza**
  - a. zalozena na pravidlech bezkontextove gramatiky
  - b. pomoci prepisovacich pravidel se da dobre zadefinovat napr priorita operaci
  - c. stvori odvozovaci strom

#### 4. semanticka analyza

- a. vytvori mezikod -- neco jako bytecode, multiplatformni "assembler" instrukce
  - i. pr. ocislovane trojice arg1 op arg2
- b. vyuziva stomu vytvoreneho syntaktickou analyzou -- prochazi strom rekurzivne, kdyz dojde nad listy, muze uz udelat prvnι jednoduchε prikazy, z nich hodnoty pak muze skladat do dalsich atd

#### 5. generovani kodu -- prevod mezikodu na assembler

- a. rozhodneme, co se da do jakych registru a jak se dana operace mezikodu mapuje na instrukci konkretnιho assembleru
- b. optimalizace:
  - i. primo na trojicich mezikodu - hledani zavislosti, opakujici se trojice atd
  - ii. pri prevodu na assembler
- c. dostaneme .o soubor

#### 6. linkovani

- a. nezname symboly se musi sestavit
- b. hledaji se zavislosti, symboly ktere musime importovat
- c. **static linking** - vsechno se hodi za sebe, prepocitaji se adresy aby odpovidaly sestavenemu programu
  - i. program muze byt velky, knihovna se kompiluje nekolikrat
  - ii. nepotrebuje zadne vnejsi zavislosti ke svemu spusteni (uzivatel nemusι instalovat zadne knihovny)
- d. **dynamic linking**
  - i. knihovny nejsou pridany do knihoven
  - ii. nevychoda -- musi byt nainstalovana knihovna
  - iii. vytvori se stub u volani funkci z dynamicke knihovny
    - 1. pri prvnim volani se zjistι kde je knihovna, nacte se a zavede se spravna adresa jejich funkci