





# Proměnné

- statické/globalní → jediná instance, umístí je překladač
- lokální - proměnné a parametry + pomocné proměnné
- thread-local → jediná instance per vlákno

Live range = místa v procedurě, kdy je proměnná zapotřebí  
 - zkončí se u skokových tok. prom. bez aliasů

↓  
 do registrů  
 + vlně, je  
 zápis modifikuje  
 celou proměnnou

... = x kdy nepotřebujeme skladovat x → místo něj tam můžete být něco jiného a disj. rozsahem platnosti  
 x = ...  
 registr či zápis.

## Variable splitting/renameing

= rozdělení proměnné s nesouvislým rozsahem platnosti

# Alokace

- statická/globalní → první adresa (překladač rozhoduje, kde to bude)  
 - pokračitel překladače analyzuji, která uložena používají které statické proměnné, aby je umístily dál od sebe pro zabránění false sharingu
- dynamická  
 - knihovni volání - překladač neví

## registra

- Skalární lok. proměnné bez aliasu - **nesmí na to být namířen pointer**

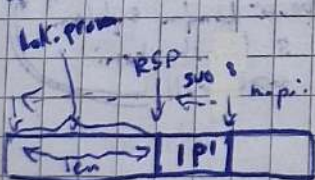
• záložní ková → slotové či aliasované lok. prom. + promy, které se nevedly do reg.

- záložník dle procesoru, nebo emulován pomocí GLOBAL reg. - branch and link instr. { reg=IP, IP=addr  
 ↳ volání funkce + BR reg. zřet

- před reg. alokací je to (strójově) možná ve virt. registrech - optimističtý by instr. pracují s prom., jak by byly v registru

- spill kód = kód, který se tam musel přidat navíc

když se zjistí, že se to nevejde, tak se to musí přelítat  
 - potenciálně se na to pak musí znovu pusht možnost optimalizace



```
CALL addr {
  sub rsp, 8
  mov [rsp], IP
  jmp addr
}
```

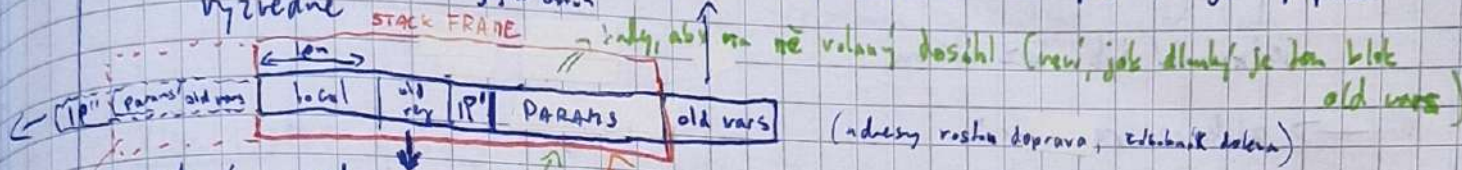
```
RET {
  add rsp, 8
  jmp [rsp-8]
}
```

Sub rsp, len  
 [rsp + Δ] → přístup k prom.

IP = addr  
 není přístup → nejde ho emulovat

# Volání konvence

- povinnost zachovat obsah registrů (všech/úplných/části)
- volající u těch reg., které potřebuje zachovat a nemá záruku, že budou zachovány, pomocí pushů dá před callu na zásobník a pak si je popem vyvede



- volaná procedura musí zachovat hodnoty reg., u kterých má volání zaručeno, že budou zachovány → buď na ně nesáhá, nebo je ona musí pushnout na z.s. a před returnem zase vrátit (nebo si zvolí prostor len a nastane to tam normálně tou)

občas se takhle na z.s. ukládá něco zbytečně - aby se to vyřešilo, tak se to rozdělí reg., které se musí a nemusí zachovat, dělá píť na píť

umístění parametrů - reg. + z.s. starší argument, kterým se řešou typy těch parametrů a on říká, do jakých reg. je dát

- bývají jen vstupní - prostě se po návratu z callu jenom přes ně posune z.s.

odpovědnost za úklid z.s.

umístění návratové hodnoty

- hodnotově - u té části PARAMS - z toho si to volání pak přeposílá do lok. prom.
- pointerem - do PARAMS se dá pointer na data, která jsou v lok. prom.
- reg.

- úprava jména procedury kvůli linkeru - mangling
- např. namespaces, třídy, ...

## Fáze

- instr. selection - typicky 1:n
- instr. scheduling - NP-úplná úloha
- rego alloc - NP-úplná úloha (binární grafy)

# Reprezentace

omezikód střední úrovně

- nesokracení - uvnitř BB jsou grafy
- část. sekv. - uvnitř BB jsou pseudoinstr.
- plně sekv. - "bez BB" - celé je to jeden proud instr.

- Součástí musí být i třeba seznam glob. prom., procedur a další informace nad rámec samotného kódu

omezikód nízké úrovně

- ekvivalentní stroj instr.
- ve výsledku chceme plně sekv. formu, ale optimalizace se lépe dělají na těch část. sekv.
- do plně sekv. formy se přechází až v pozdních fázích

- většinou nemá smysl moc používat nějaké komplikované stroj. instr., protože ten procesor si to pak stejně rozkládá na p-operace + víc menších instr. dává větší prostor pro optimalizace v rámci schedulingu překladačem

# Analýza rozsahu platnosti (live range analysis)

**DE:** Proměnná je v bodě a živá  $\equiv \exists$  cesta v CFG do bodu b t.j.:

- na cestě  $a \rightarrow b$  není žádný zápis do k' proměnné
- v bodě b je tato proměnná žena

jednoduché lot. bez aliasu

→ alokace registrů

**Alg:** CFG  $V=BB$ ,  $E=CF$  přechody, VAR... proměnné

1) lot. analýza - uvnitř BB

- W... do proměnné se zapisuje
  - R... k' proměnné se tam před 1. zápisem čte
- funkce  $BB \times VAR \rightarrow \{0,1\}$

2) glob. analýza

- L... proměnná je na začátku BB živá

1)  $\forall v \in VAR$ :

2)  $\forall b \in BB: L(b,v) \leftarrow R(b,v)$

3) Dohled dozadu ke zdrojům

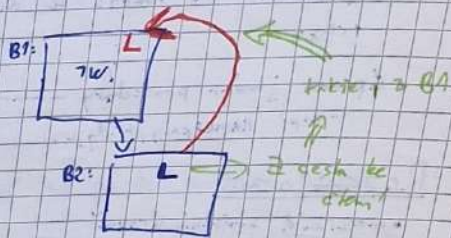
4)  $\forall (B1, B2) \in CF$ :

5)  $L(B1, v) \leftarrow v \cup W(B1, v) \& L(B2, v)$

(IBB) (CF) (VAR)

3) dopocít

- zač. a konce BB



tržba k' je zápis  
přítelk' se zápis  
nějak zápis  
nějak vybití  
k' významí, tak  
k' minimální  
příčet prakticky  
registru

## Alokace reg.

- vstupem je část. sek. kód - ekvivalenty stroj. instr. s virt. reg. a optimistický
- + máme spočítané rozsahy platnosti prom. - intervaly, kdy je prom. živá
- výstupem je přiřazení virt. reg. fyzickým

- Alg:**
- 1) vytvoření matice kódu  $C: C_{ij} = \begin{cases} 1 & \text{prom. } i \text{ a } j \text{ mají společnou platnost} \\ 0 & \text{jinak} \end{cases}$
  - 2) spočítání odhad ceny za spill kód - cena za přístup k proměnné, když nebude v reg.
  - 3) barvení grafu  $(VAR, C)$  n barvami - NP-úplná úloha  $\rightarrow$  heuristické řešení
  - 4) doplnění spill kódů

- ne je, ale ten píece kdy máte potřebovat nějak reg. a vyhradíte - to zvláštní reg. které se nečítají alq. aktivní v instr. pesimistický předpoklad, že bude spill kód a a to posílá

přítelk' L optimalizací  
kterou z  
deset-ánek  
binou používáme  
má záčet ze  
příjem obrnit

pak by s nejmenší  
cennou za spill kód  
někdy vřady s co  
největším/nejmenším  
shpřím

příjde obrnit  
oborování  
příběh do grafu  
registru do rebodu v reg.  
(a nezávislé by do grafu)

Optimalizace:

- minimalizace přesunových instr.
- zkratěvní vst. reg. spojené přes-instr. (pokud to jde) <sup>keřítak může pomoci</sup>
- úprava alg. - do grafu kolizi přidáme nové typ hrany, která bude říkat, že by bylo dobré, kdyby dané vrcholy měly stejnou barvu <sup>duplikace kódu</sup>
- řídíme se tím při volbě barvy

Vliv instr. sady

operand instr. má určenou množinu reg., v níž může být ortog. vzhledem

ortogonální  $\equiv$  dvě množiny povolených reg. jsou buď identické, nebo disj.

→ abstrakce reg. pro množinu vstřet

neortogonální

→ do toho grafu kolizi přidáme úplný graf s vrcholy odpovídajícími všem fyz. reg. a spojíme ty VARS s těmi reg., kde nemohou být + by nové vrcholy dostanou nejvyšší prioritu (nikdy je nepřibíráme jako oběť)

+ to, který register odpovídá které barvě, je určeno tím, jakými barvami jsou obarveny ty fyz. reg.

- může ale být problém, že vyjde, že je nějaká proměnná spojena se všemi fyz. reg.  $\Rightarrow$  obarvená  $\Phi \Rightarrow$  stálo by se to do paměti  $\Rightarrow$  FUJ

- neortogonalita taky způsobuje volací konvence, která určuje, v jakých reg. co může být (i když by reg. byl dostupný)

(např. jeden se používá jako 2. parametr nějaké funkce a pak jako 3. parametr jiné funkce)

→ nepříjde přiřadit reg.

triede, když se má nějaká proměnná a vrcholy v grafu, tím množinami, jak počítat průběhův registry

překážka k optimalizaci, kterou z dosahových barev používáme

maže se z přístupu obarvit

pat by s nejvyšší cenou za spill kód nebo vrcholy s co největším / nejmenším stupněm

příjde obarvit  $\rightarrow$  obarvit a přikázat do grafu  $\rightarrow$  nepříjde  $\rightarrow$  nebude v reg. (a nevrátíme ho do grafu)

Scheduling - rozvrhování instrukcí

- typicky těsně před fyz. abstrak reg.

+ když abstrak reg. způsobí vytvoření spill kódu, tak se pak scheduling pustí znovu

- ILP

- Pipelining
- MIMD - superlativní prac.
- víc pipeline - v jednu chvíli se provádí stejné fáze několika instr.

• SIMD

- nečítá se při schedulingu, ale dřív - při generování instr. a často když je ten kód ještě ve Stromové formě

- provádí se na jednom BB a pak se rozšíří mírně dál:

• su pipelining - pro BB, které představují výšky

• trace scheduling - slíh' post. BB do jednoho - ale pořád v rámci jednoho procedury (pokud nebyla in-line) - ale pozor na závislosti, odskoky a příst.

kolem každého použití proměnné bude oblast deklaruji se reg-reg nové instrukcemi

tím se zúžností te' proměnné rozděl do víc, které se přetýkají rozumně

→ jde použít dokonce jednodušší alg. na barvení grafu

ale pozor na závislosti, odskoky a příst.

plivání instrukcí  $\rightarrow$  co nejvíce využít pipeline

- vstupem je DAG s uzly = instr. a hranami popisujícími závislosti

+ model procesoru, který řekne, jak je určité permunice instr. dobrá

inverze  
doba  
dok. času  
instr.

jaká ref. bod se bere

- výstup: uzlím toho DAGu přiřadíme časy (v jednotlivých tab. procesoru)  
(a hlavně určime topolog. pořadí tím grafem)

- latence - časová závislost dvojic instr.
- rezervacní tabulky - kdy která instr. obsazuje která výp. jednotky

+ kap. procesoru  
↓  
# výk. jednotek

Skutečné časování se ale může lišit - třeba kvůli rozpracování instr. z předchozího BB nebo kvůli nepřítomnosti dat v cache (problém předpokládá, že jsou v nejbližší cache) nebo kvůli době dok. dat instr.

- předávání dat v reg. :-)
- předávání dat přes paměť :-)
- musí se opatřit

Základní alg. - List scheduling

- jdeme postupně - vždy máme nějaké připravené instrukce, z nichž jednu rozvrhneme

$O(|E| \cdot (\log V + |E|))$

Krit. cesty = délka od té instr. do konce BB

čas na vst. tab. vs. čas na vst. tab. vs. čas na vst. tab.

instr. je to min. součet délek, kdy je s akt. rozvržením schopna být uložena - i s latencemi předchozích instr., které už jsou rozvrh.

• při schedulingu p. akt. reg. přibýdou anti-dependence na prom., které byly dány do stejného reg.

• kontrolní dependence v rámci Erase schedulingu a SW pipeliningu

• další technické závislosti - napr. manipulace se zás.

Vylepšení

• branch-and-bound

- backtrackingem zkoušíme, jestli se nám nepovede najít lepší řešení
- po několika pokusech se zastavíme a lepší řešení už se přestane hledat

• SW pipelining

- pro cykly - přes hranice BB

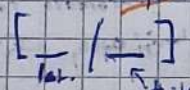
• Unroll-and-compact

- rozvineme ten cyklus a na to použijeme ten normální alg. => nejdřív krajem roli jen hraný s (a, v další iteraci i s / 1, ...)

- představujeme si, že je tam ten BB nekonečněkrát - řešíme asymptotické chování ->  $T = T_0 + k \cdot N$

- deklarijeme, kdy se už zacyklíme

- Kritická smyčka

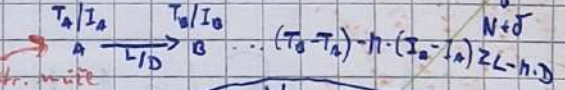


-> cyklus v grafu závis. s největším maxim. součtem délek instrukcí

• modulo scheduling

- snažíme se přímo najít rozvrh s periodou M (když tak pat. zkoušíme M1, M2, ...)

- roz. tab. používáme modulo M, instr. rozvrhujeme do prostoru čísel x ikence

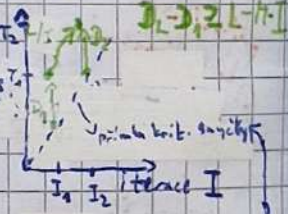
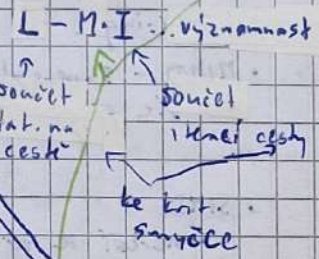


instr. může být v cache, když má být rozvržena

vedou z k' první

acyklické

- přivádíme instr. při rozvrhování



$T_1 - M \cdot I_1 = D_1$   
 $T_2 - M \cdot I_2 = D_2$

instr. rozvrhneme na co nejmenší čas, to pak opravíme přidáním virtuálnosti => nejvyšší prioritou má ta instr., s nejvyšším číselným číslem => rozvrheme

-> dolní odhad, jak to může být pomalejší (ale může to někdy být i horší) opravíme podle rezerv. tab.

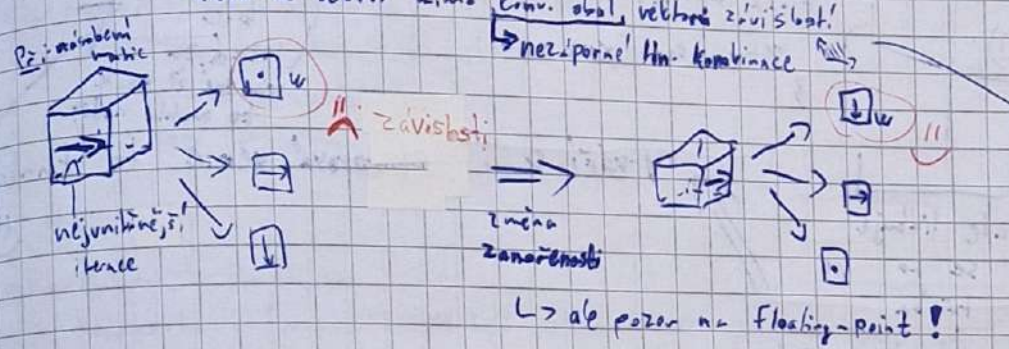
variable expansion

- pro eliminaci anhidependenci!
- pomocí duplikace kódů a přeměnění proměnných
- největší duplikace a pak rozvahování
- největší eliminace anhidependenci (přeměněním) a pak duplikace a pak rozvahování
- ale pozor na spill kód

Vektorizace - jemnozrnný paralelismus

- překladač si musí být jistý, že ta transformace bude korektní!
- cykly musí mít předvídatelný # iterací
- cykly s afinním chováním - +/\*/cmp
- > afinní transformace:

- o loop reversal - výměna zanoření
- o loop skewing - změna souřadného systému v iter. prost.



- meze musí být konstanty (ne nutně zadané za překladač, ale proči se nesmí měnit)
- polyhedrální kompilace - ten prostor se natáhne na n-dim. mnohostrany
- musí se zapnout!
- např. POLLY v LLVM
- loop skewing je speciální případ rotace

Jemnozrnný paralelismus

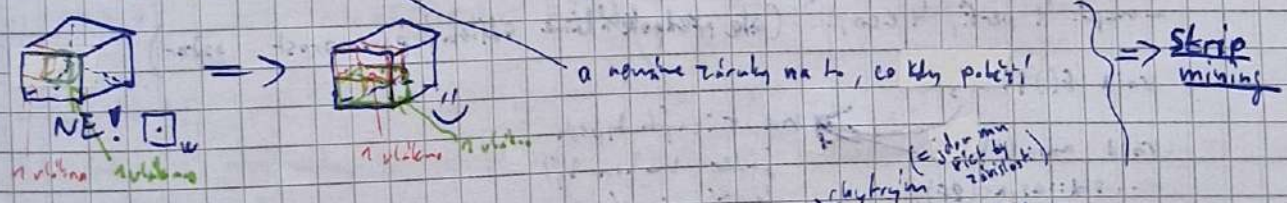
- o ILP -> scheduling
- o SIMD

Hrubozrnný paralelismus

- > SVE
- o SMP
- hyper-threading
- multicore, multi-socket
- o MHA
- o Cluster
- Strip mining - dělení na větší bloky

Hrubozrnná paralelizace

- typicky ne u nejvnitřnější smyčky, ale naopak u nejvnější => nezávislé bloky
- > dostatečně velké, aby se to vyplatilo
- krmíme nějakou knihovnu - OS, framework, ...
- > neumíme řídit závislosti => potřebujeme mít zcela nezáv. bloky



- > lepší je, když máme k dispozici taskový paralelismus s vrstevně schedulingem (a task stealingem)
- => řídíme si to, jak to krmíme => synchronizace v tom scheduleru k řešení závislosti

- chceme nahraze mít ty nezáv. bloky
- a dale tak, aby chom mohli dělat jemnozrnnou paralelizaci, závislosti mohou být v pohledu naší
- => typicky ale máme malý prostor (treba zaručený cyklus tabulky neodtáhne)
- => rozdělíme na víc bloků a už to vyjde (a může to pomoci cachím)
- > a nad nimi děláme strip mining

možná třeba vzájemně inkonzistentní nejmenších funkcí!

# Optimalizace pro cache

→ prostorová lokalita

→ časová lokalita

→ ideálně přístupy hned vedle sebe, ale tím zase rozbijeme přístupy jiné  
 ⇒ z-křivka → nepravidelné přístupy ☹️

⇒ Kompromis



číslo je u toho nějaký #pragma

## Problémy s linkováním

```
for(...)
    a[c] = f(b[i], c[i]);
```

```
float f(float x, float y) {
    return sqrt(x*x + y*y);
}
```

ale když se to linkuje dynamicky, tak to neprojde inlinem a vektorizací

vytváří se i vektorová verze té funkce f ☺️

kdysi překladač toho b.c. věděl, že vhodně, že by se to mohlo hodit

kdysi ne, tak to linker udělá v rámci link-time optimizations

## Heterogenní jádra

nebo dokonce GPU

- např. 4 perf., 4 eco, ... (ale předpokládáme sdílený paměťový prostor)

```
void f() { ... }
void main() {
    ... std::async(f);
}
// thread pool
```

na různých typech jader mohou být lepší různé optimalizace a nebo taky dokonce mohou mít jinou instr. sadu

tabulka

je potřeba mít víc verzí f() a rozšířit ten koncept pointeru na funkci

víc pointerů v jednom

# Analýza aliasů (points-to analyza) Kam pointer míří (v Javě escape analyza)

↳ jestli 2 výrazy označují stejnou paměťovou lokaci

- indexy !!
- překryvy pointerů !!

- myšlenkové dynamiky ⇒ obravský graf  $V = \text{prom.} + \text{dyn. abt. bloky}$   
 → staticky  $\Leftrightarrow \text{points-to} + \text{čas. rozložení, kdy to platilo}$

- grafy variabil různými běhy složený do jednoho (ekvivalence nad hodnoty jako  $\text{map}$ )  
 - pokusíme umístit dokázat pře nějaké pointery ne míří nikdy (vždy vytvořit na zásobníku)  
 - norm. odhad dyn. analyzy

- každý proměnná má svou skupinu, dyn. abt. objekty mají skupinu = Kontext (troub. i celý stack trace)

## • bezkontextová analyza

- neřešíme ten kontext
- informace společně pro celý program

## • kontextová analyza

- points-to informace může být relativní vůči pozici v programu (IP), nebo s call stackem (např. číselným)

## Andersenova metoda (bezkontextová)

- u  $\forall$  prom.  $p$  si představíme množinu míst, kam může ukazovat

⇒ soustava množinových nerovnic: ⇒ vyřeší se metoda nejím. prvního bodu

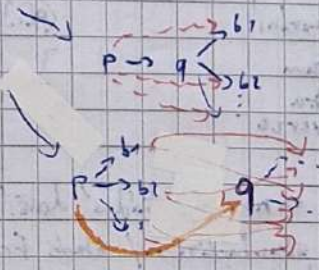
$$p = \&x; \Rightarrow \{x\} \in \sigma(p)$$

$$p = q; \Rightarrow \sigma(q) \subseteq \sigma(p)$$

$$p = *q; \Rightarrow \forall b \in \sigma(q): \sigma(b) \subseteq \sigma(p)$$

$$*p = q; \Rightarrow \forall b \in \sigma(p): \sigma(q) \subseteq \sigma(b)$$

2) a opakujeme, dokud se to nezastaví ⇒ nejím. první bod



ale mohou to být

struktury a prvky mohou být uloženy na další prvky... - ale my to máme celé v jednom

⇒  $\sigma(p, m) :=$  místa, kam může ukazovat pointer v objektu  $p$

$$\hookrightarrow p.m = \&a;$$

kvadratická složitost !!

krádež derivance  
 rovnost (X Andersen - nerovnost)

## Steensgaardova metoda (bezkontextová; méně přesná, ale rychlejší)

$T(p)$ ... "typ" proměnné  $p$ , na začátku  $\forall p: T(p)$  nazýváme všemi, přířazení ⇒ koloniální typy

- zase můžeme rozšířit by pointerky ( $T(x, m) + \text{pravidlo } T(x) = T(y) \Rightarrow \forall m: T(x, m) = T(y, m)$ )

# Částečné vyhodnocení

- konst. podvýrazy, sparse conditional const. prop. - sledování, kdy jsme do proměnné uložili konstantu

## Algebraické úpravy

- machine idioms - párkrát instrukce, která udělají víc věcí najednou
- odstranění zbytečných kontrolních vět
- distributivita
- conditional move → užitečný podm. skok

napiš FMA

→ různé věci, kdy tam přijde, to mohou být různé konstanty nebo runtimové hodnoty

## Odstranění redundance

ale můžeme něco dělat zbytečně !!

může pomoci profiling  
může spojit více operací do jedné

- Copy Propagation - paralyzujeme si, že proměnná je kopie jiné, napojeno na analyzu rozsahu platnosti prom.
- Elimination spol. podvýrazů - velmi dobře se detekuje v SSA mezibodu
  - lokální - umí BB, třeba pomocí hashování
  - globální - umí procedury
- přesun invariantního kódu z cyklu - n.př.  $for(i < k.size()) \Rightarrow S = k.size()$   
 $for(i < S)$

- partial-redundancy elimination
  - ↳ redundance, ale jen někdy (v ifu)
  - ⇒ kopie duplikace BB s tím vyřadím bez

## lazy code motion

→ něco spočítáme, až když je to potřeba, aby bylo to potřeba  
- bariéry: • návraty z funkcí  
• podmínky v control-flow  
- rovnou eliminuje mrtvý kód • ne-loop-invariant výrazy

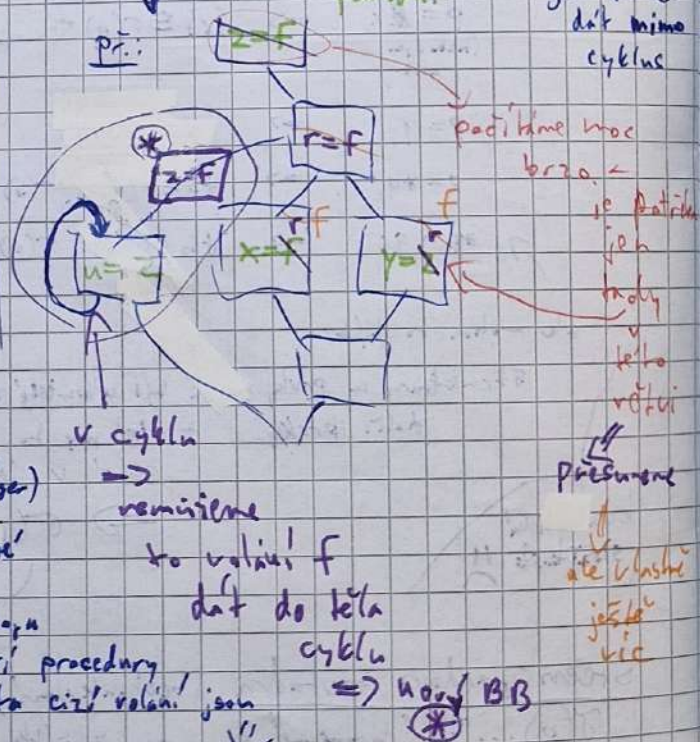
## Odstranění neúčinného kódu

- mrtvý kód - něco dělá, ale zbytečně
- nedosažit. kód (nejčastěji v důsledku předchozích optimalizací)
- optimalizační skoky - skoky na skoky ⇒ skok

## Integrace procedur (inline expansion)

- ⇒ duplikace kódu
- problémy v jazycích, kde je oddělený příklad od linkerů
- ↳ musí se tam nechat i ta původní verze

```
n.př.
x = a + b
for(i...) {
    p[i] = x + i
}
```



- listové procedury = nevolají žádné další
  - není potřeba standardizovaný stack frame (třeba pro debugger)
  - lze zjednodušit
    - např. vypustit ukládání některých registrů v příloh a obnovení v epilogu
- shrink wrapping - ve větších, které nevolají další procedury
  - prolog a epilog dáme jen do větvi, kde ta cizí volání jsou
- tail merging - u procedur se stejným koncem (detekce např. pomocí hashování)
  - spojíme přes skoky ⇒ mrtvý kód (ale problémy s debuggerem - není, když se tam spojí)

## Specializace procedur

- když se procedura volá s konst. parametry
  - pro ty varianty natlačujeme a přizpůsobíme je okolnostem (→ např. eliminace podmínek)
- nebo taky pro řešení aliasingu → verze pro ja i ne
  - ↳ optimalizovanější verze
  - ↳ opatrnější verze
- nebo třeba  $F(x \neq 1)$ 
  - $F(x)$  a to  $+1$
  - + if (jo) → jo() else → ne()
  - az někde v kódu (může jít slepit s jiným  $+1$ )

## Interproc. alokace reg.

- úprava volací konvence - např. nemusíme ukládat nějaké registry
- oproti integraci procedur tu úpravou funkci můžeme volat z víc míst

## Návrhy pro procesor

- branch prediction
  - prefetching do cache (ignoruje ujimky ☹️ ~~xx~~ "hard load předem")
  - instruction prefetching
    - překladač se by instrukce (BB) snažil
    - procesory instrukce načítají po balíčcích (např. 16B u Intelu)
    - rychlejší (lip kódu dekodují)
- na table zaznamenal  
- ale tím se zase plyná kódová cache  
→ znamená se třeba jin načítaly BB v cyklech