

a SINT

### 3. Datový paralelizmus SIMD<sup>v</sup>, HW

### implementace a SW podpora. *na CPU a GPU*

V praxi sa vyskytuje rada výpočetných (pod)problémov, ktoré potrebujú vykonať rovnakú (nezávislú) transformáciu pre množinu prvkov = dátový paralelizmus. Pragmatický prístup výrobcov HW adresuje práve túto spoločnú vlastnosť výpočtov (rovnaká transformácia N prvkov) na dosiahnutie vyššieho výkonu (throughput), a to replikáciou funkčných jednotiek. Keďže operácia nad prvkami je rovnaká, netreba aby sa replikovali plnohodnotné procesory, ktoré obsahujú vlastnú riadiacu jednotku. Pridané procesory sú riadené jednou jednotkou, a teda dochádza k úspore miesta na čipe, resp. ceny. Zdieľanie riadiacej jednotky znamená, že všetky procesory pracujú synchronne a v danom čase vykonávajú tú istú inštrukciu - a teda sa takéto architektúra označuje Single Instruction, Multiple Data (SIMD).

Silné stránky:

- spracovanie vektorov, matíc, tenzorov

Slabé stránky:

- Transformácie, ktoré obsahujú veľa podmienok *, obsahujúce hodne vetvení*

Typy paralelizmu v CPU - kde sa zaraďuje SIMD?:

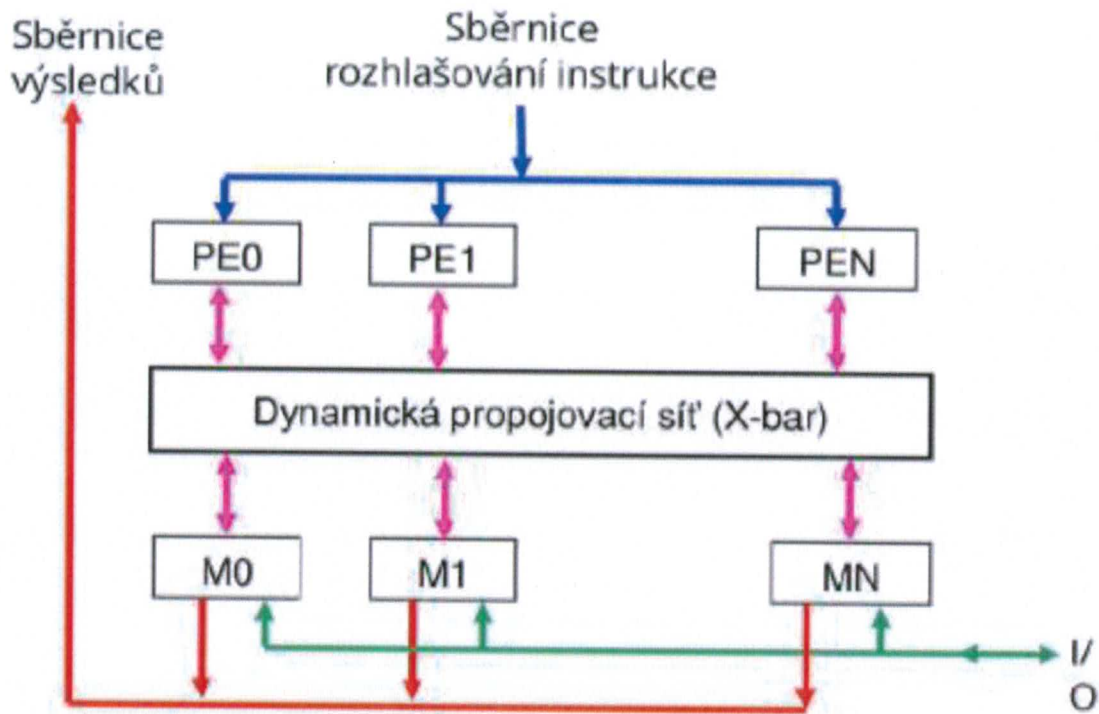
1. **Funkčný paralelizmus** - rôzne inštrukcie používajú rôzne časti CPU, a preto je ich možné vykonať paralelne (za predpokladu, že medzi nimi nie sú dátové závislosti). Príkladom sú superskalárne procesory (paralelizmus je riadený HW - OoO), a procesory VLIW (Very long instruction word; paralelizmus je riadený kompilátorom)
2. **Dátový paralelizmus**
  - a. **V čase** - zreťazené spracovanie
  - b. **V priestore** - replikácia funkčných jednotiek (SIMD)
  - c. **V oboch dimenziách** (paralelná vektorové linky, prípadne SINT - Single Instruction Multiple Threads)

### SIMD Architektúry

- **Obsahujú veľký počet výpočetných jednotiek (processing units PE) a jednu riadiacu jednotku**
- Počet PE na CPU a GPU je v rozsahu 4-32
- pri návrhu SIMD architektúry je potreba riešiť dilema ako moc zložitých budú PE (čo všetko budú vedieť spočítať). Čím zložitejšie, tým viac miesta zaberú a tým menej sa ich na plochu o daných rozmerov zmestí.
  - Typicky umožňujú veľký počet jednotiek, ktoré umožňujú základné výpočty (ADD, SUB, CMP, ...) a malý počet jednotiek pre náročnejšie operácie (DIV, SIN, SQRT, ...)

- Typicky realizácie týchto architektúr vystupujú v úlohe koprocesorov (GPU, AVX), prípadne HW akceleratorov – ich použitie ako general purpose CPU by bolo problematické, keďže pre niektoré druhy výpočtov sú pomalé.

### SIMD so zdieľanou pamäťou:

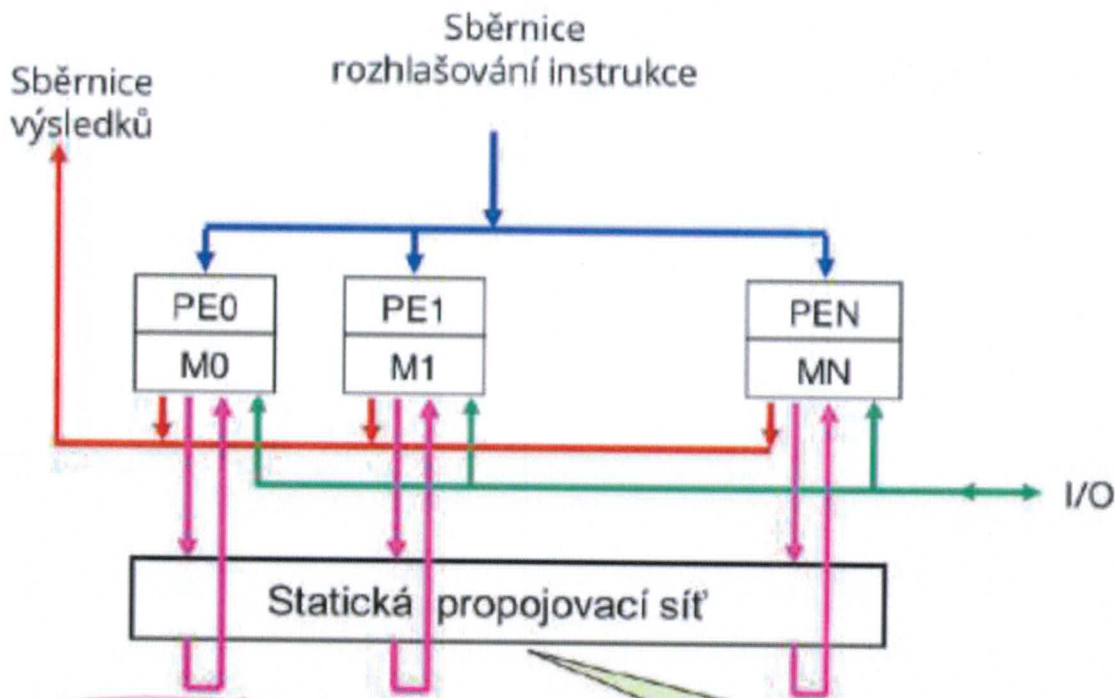


- Každý PE má rovnaký prístup k pamäťovým modulom, ktorý prebieha prostredníctvom dynamickej prepojovacej siete (napr. crossbar - križový prepínač)
- Komunikácie medzi PE prebieha cez zdieľanú pamäť
- Príklady realizácie: AVX jednotky v CPU, SM procesory v GPU

- halle dnes využívajú procesory a GPU
- dá sa to modelovať PRAM modelom a CREW

Kdy má vektorizace smysl? - Aritmetická intenzita

- výronnost řádkového algoritmu je slova omezena buď propustností ALU nebo paměti
  - pokud nedosáhneme rozumné aritmetické intenzity (FLOP/B), nemá smysl se do vektorizace vůbec pouštět
- SIMD s distribuovanou pamětí:



- Každý PE má vlastní paměťový modul - výpočty vykonává nad lokálními daty (data sú blízko pri PE = rýchlosť prístupu)
- Komunikácia s ostatnými PE prebieha cez statickú prepojavaciu sieť (mriežka/torus a pod.)
- Príklady realizácie: IBM Cell v Sony playstation 3

Orientácia v SIMD v moderných CPU - inštrukčné sady

- X86
  - SSE (SIMD Streaming Extensions! - 128b registre. Pôvodne 4x single precision float, s neskoršími vydaniaми (SSE2, SSE3, ...) sa pridáva možnosť interpretácií: register = 2x double | 16x byte | 4x 32 bit integer | ...
  - AVX - rozširuje šírku registra na 256b
  - AVX512 - rozširuje šírku registra na 512b
- ARM
  - ARM Neon - inštrukčná sada kombinujúca 64b a 128b operácie

SW podpora

Ako môžem použiť SIMD na urýchlenie výpočtu?

- Vektorizované knižnice
  - Intel MKL, TensorFlow

- Automatická vektorizace kompilátorem (je nutné uistiť sa, že kompilátor bol schopný kód automaticky vektorizovať)
- Pragma hinty kompilátoru:
 

```
#pragma omp simd
for (int i = 0; i < K; ++i) { a[i] = 2*a[i] + b[i]; }
```
- Vektorové intrinsic funkcie:
  - Hlavičkový súbor, ktorý pridáva C funkcie, ktoré sa svojim menom podobajú assembleru (napr. `__mm_add_ps`). Tento hlavičkový súbor rovnako pridáva dátové typy, ktoré odrážajú SIMD registre (`__m512`).
- Asembler kód

### Automatická vektorizácia kompilátorom - čo dokáže vektorizovať?

- Smyčky, kde je známy a počas iterácie nemenný počet iterácií (*Kde je predem známy počet iterácií*)
- Smyčky s jedným vstupom a výstupom
- Smyčky, ktoré neobsahujú zložité vetvenie. Jednoduché podmienky je možné vektorizovať s využitím maskovania, napr. v nasledujúcom kóde sa vyhodnotí podmienka pre celý SIMD register obsahujúci napr. 8 float čísiel pričom výsledok je binárna maska, ktorá hovorí, či táto podmienka platí alebo pre daný element registra. Následne sa táto maska použije a do elementov, ktoré ju spĺňajú sa presunie výsledok then vetvy, a do časti ktoré ju nespĺňajú sa presunie výsledok else vetvy.
 

```
for (int i = 0; i < length; i++) {
    if (s >= 0) x[i] = b[i] * b[i] - 4 * a[i] * c[i];
    else x[i] = 0.;
}
```
- Iba najvnútornejšie smyčky (*pouze najvnútornejšie smyčky*)
- Smyčky bez volania funkcií až na intrinsic matematiku, inline funkcie, OMP SIMD funkcie...

*keďže bude  
ve smyčce  
break, bude  
to problém*

*jednoduchá  
vetven  
pres blonde*

*#omp declare simd*

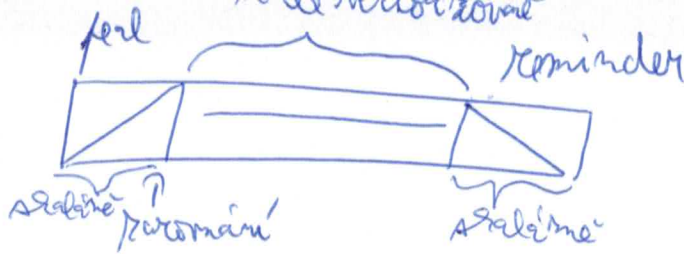
### Automatická vektorizácia kompilátorom - kde zlyháva? *Kdy to nejde:*

- Nejednotkový rozostup - je nutné, aby elementy boli za sebou v pamäti s konštantným rozostupom. Ak je rozostup  $> 1$ , je možné použiť inštrukcie gather/scatter (typicky obsahujúce penalizáciu).
- Nezarovnané štruktúry
- Dátové závislosti medzi iteráciami
- Pointer aliasing - je nutné testovať za chodu, či sa použité pamäťové lokality môžu prekryvať

*- musíš říct, že je vektorizace možná, ale dle revidovaný kompilátoru nepředimá*

- Zlé zarovnanie pamäte je jeden z hlavných dôvodov zlého výkonu. Kompilátor preto vkladá do kódu test na zarovnanie, a pokiaľ nie je spracovávané pole korektné zarovnané je vygenerovaný peal, body, reminder - body obsahuje vektorizovaný výpočet na

peal a reminder  
 se chce zarovnat  
 => pecho zarovnání



zarovnanaj pamäti, peal obsahuje nevektorizovaný kód spracúvajúci časť poľa pred zarovnanou časťou, obdobne reminder je nevektorizovaný kód spracúvajúci časť za zarovnanou časťou.

- Pre predídanie generovania nevektorizovaných častí (peal, reminder) je nutné explicitne deklarovať zarovnanie, a to pri:
  - Alokácii dát (použiť napr. `aligned_alloc` namiesto `malloc`)
  - Deklarácii nového ukazateľa
  - Deklarácii novej funkcie ktorej parameter je ukazateľ

- Keďže je vstupem funkcie ukazateľ, tak kompilátor nebude predpokladať zarovnanie

```
#pragma omp simd align(a:64)
```

říkáme kompilátoru, že a je zarovnáno

## SIMT (Single Instruction Multiple Threads)

- typ SIMD

- většinou GPU

- např: ~~matice~~ ~~matice~~ ~~matice~~ ~~matice~~

smyčky úplně rozbíjíme a řádky jeden prvek v matici bude počítat jedno vláknem (musíme jich mít až miliardu)

Programyernal pře kód pro jedno vláknem a pracujeme s jeho indexem pro určení jeho "práce"

- Strady se sbíhují do WARPů, což je vlastně skupina vláken, které mají společný programový článek (obdobna AVX), WARP pře do bloků, které se přidávají <sup>na jednotlivé</sup> procesory
- GPU pře vždy zpracovává celý WARP a switchuje mezi nimi pro udržování celku
- GPU jsou většinou skalární in-order

# Programování pro GPU (SIMT)

- mikrokontroléry - CUDA
- využití OpenMP direktiv (OpenMP 4.5)
  - využívá <sup>na pozadí</sup> řady, protože GPU nedovoluje bariéry

```
# pragma omp target
# pragma omp teams
# pragma omp distribute
```