

# **Zpracování medicínských dat pro výpočet typického dne pacienta s diabetem**

Semestrální práce z předmětu KIV/PPR na ZČU v Plzni

Bc. Jiří Winter

31. ledna 2026

# Obsah

<b>1</b>	<b>Úvod a Struktura projektu</b>	<b>3</b>
1.1	Medicínská data . . . . .	3
1.2	Adresářová struktura . . . . .	3
1.3	Spuštění programu . . . . .	3
<b>2</b>	<b>Algoritmus zpracování dat</b>	<b>4</b>
2.1	Fáze 1: Medián pro každého pacienta . . . . .	4
2.2	Fáze 2: Medián napříč pacienty (Grand Median) . . . . .	4
2.3	Fáze 3: Downsampling (Redukce časových oken) . . . . .	5
<b>3</b>	<b>Načítání a Reprezentace dat</b>	<b>6</b>
3.1	Flat Data Layout . . . . .	6
3.2	Parsování časových razítek . . . . .	6
3.2.1	Frekvence 15ms . . . . .	6
3.3	Handling prázdných hodnot (Padding) . . . . .	7
<b>4</b>	<b>Sekvenční zpracování (CPU)</b>	<b>8</b>
4.1	Ukázka výpočtu mediánu . . . . .	8
4.2	Ukázka sekvenčního algoritmu . . . . .	9
4.3	Výběr algoritmu: std::nth_element . . . . .	10
<b>5</b>	<b>Paralelizace a SIMD Optimalizace</b>	<b>11</b>
5.1	Paralelizace algoritmu . . . . .	11
5.2	Varianta 1: Paralelizace bez SIMD . . . . .	11
5.3	Varianta 2: Paralelizace se SIMD . . . . .	11
5.3.1	Princip SIMD řazení - Sorting Network . . . . .	13
5.3.2	Extrakce mediánu ze seřazených registrů . . . . .	13
5.4	Medián napříč pacienty . . . . .	14
5.5	Downsampling časových slotů . . . . .	15
<b>6</b>	<b>GPU (OpenCL)</b>	<b>17</b>
6.1	Kernel 1: Median per Patient . . . . .	19
6.2	Kernel 2: Reduce Patients . . . . .	19
6.3	Kernel 3: Downsampling (Shell Sort) . . . . .	21
6.4	Řešení pro BVP data (740 milionů záznamů) . . . . .	21
<b>7</b>	<b>Naměřené výsledky a Vyhodnocení</b>	<b>23</b>
7.1	Hardwarová konfigurace . . . . .	23
7.2	Časy načítání dat . . . . .	23
7.3	Dexcom Data (Nízká zátěž) . . . . .	23
7.4	HR Data (Střední zátěž) . . . . .	24

7.5	BVP Data (Vysoká zátěž) . . . . .	24
7.6	Analýza fází GPU výpočtu . . . . .	25
<b>8</b>	<b>Analýza a Diskuze výsledků</b>	<b>26</b>
8.1	Zrychlení a efektivita (Speedup) . . . . .	26
8.2	Karp-Flattova metrika . . . . .	27
8.3	Škálovatelnost v čase . . . . .	27
<b>9</b>	<b>Diskuze nad implementací</b>	<b>29</b>
<b>10</b>	<b>Závěr</b>	<b>30</b>
	<b>Příloha A: Uživatelský manuál</b>	<b>31</b>

# Kapitola 1

## Úvod a Struktura projektu

Cílem této práce je efektivní zpracování velkých objemů medicínských dat (Dexcom, HR, BVP) s různou frekvencí vzorkování (od 5 minut až po 15 ms). Aplikace počítá mediány naměřených hodnot v čase, provádí redukci dat (downsampling) a využívá k tomu metody paralelizace na CPU (OpenMP, SIMD) a GPU (OpenCL).

### 1.1 Medicínská data

Práce umí zpracovat tři typy dat:

- **Dexcom Data:** Glukózová data s vzorkováním po 5 minutách (288 záznamů/den).
- **HR Data:** Srdeční frekvence s vzorkováním po 1 sekundě (86,400 záznamů/den).
- **BVP Data:** Biovoltážní signál s vzorkováním po 15 ms (5,760,000 záznamů/den).

### 1.2 Adresářová struktura

Projekt je rozdělen do modulů podle zodpovědnosti:

- **include/**
  - `DataReader.h` - Deklarace tříd pro načítání a konverzi dat.
  - `ReadDexcomData.h` - Datové struktury pro uložení záznamů (`DexcomData`).
- **src/**
  - `DataReader.cpp` - Implementace parseru CSV a převodu/normalizace časových razítek.
  - `CPUSequential.cpp` - Referenční sekvenční implementace.
  - `CPUParallel.cpp` - Implementace využívající OpenMP a vektorizaci.
  - `GPUParallel.cpp` - Host kód pro OpenCL (správa bufferů a kernelů).
  - `kernel.cl` - OpenCL kernely pro běh na grafické kartě.
  - `main.cpp` - Vstupní bod, benchmarky a validace výsledků.

### 1.3 Spuštění programu

Program je určen pro běh z příkazové řádky. Uživatelská příručka je uvedena v Příloze A.

# Kapitola 2

## Algoritmus zpracování dat

Algoritmus zpracování dat je identický pro všechny tři implementační metody (sekvenční, paralelní i GPU), liší se pouze způsobem, jakým jsou jednotlivé kroky vykonávány.

Algoritmus pracuje se třemi dimenzemi dat:

1. **Time Slots (Časová okna):** Diskrétní časové úseky během dne (např. 00:00:00, 00:00:15, ...).
2. **Patients (Pacienti):** Jednotliví měření lidé.
3. **Days (Dny):** Naměřené hodnoty pro daného pacienta a daný čas v různých dnech (6-8 dní měření).

Celý proces se skládá ze tří hlavních fází:

### 2.1 Fáze 1: Medián pro každého pacienta

V této fázi se pro každou kombinaci (**TimeSlot**, **Patient**) vezmou všechna naměřená data za všechny dny (6-8 hodnot) a vypočítá se z nich medián.

- **Vstup:** Matice  $[\text{TimeSlots} \times \text{Patients} \times \text{Days}]$
- **Výstup:** Matice  $[\text{TimeSlots} \times \text{Patients}]$
- **Operace:** Pro každého pacienta  $P$  v čase  $T$  se vezme vektor dní  $D_1 \dots D_N$ , seřadí se a vybere se střední hodnota.

### 2.2 Fáze 2: Medián napříč pacienty (Grand Median)

Následně se pro každý **TimeSlot** vezmou výsledné mediány všech pacientů z předchozí fáze a vypočítá se jeden reprezentativní medián pro tento časový okamžik.

- **Vstup:** Matice  $[\text{TimeSlots} \times \text{Patients}]$
- **Výstup:** Vektor  $[\text{TimeSlots}]$
- **Operace:** Pro každý čas  $T$  se vezme vektor pacientů  $P_1 \dots P_M$ , seřadí se a vybere se střední hodnota.

## 2.3 Fáze 3: Downsampling (Redukce časových oken)

V případě potřeby (např. pro vizualizaci nebo redukci šumu) se provádí snížení vzorkovací frekvence. Původní jemná časová okna se slučují do větších bloků.

- **Příklad:** Redukce z 15ms samplingu na 5 minut.
- **Vstup:** Vektor [TimeSlots] (např. 5 760 000 prvků)
- **Výstup:** Vektor [ReducedTimeSlots] (např. 288 prvků)
- **Operace:** Pro každé nové (větší) okno se vezme úsek původních hodnot, seřadí se a vypočítá medián.

# Kapitola 3

## Načítání a Reprezentace dat

Pro efektivní zpracování na GPU je klíčové převést data z objektové reprezentace (vektor objektů) do lineárního pole (Flat Data Layout).

### 3.1 Flat Data Layout

Data jsou uložena v jediném spojitém vektoru `std::vector<float> flat_data`. Indexace probíhá podle vzorce:

$$Index = TimeSlot \times (NumPatients \times NumDays) + Patient \times NumDays + Day \quad (3.1)$$

Tento layout zajišťuje, že data pro jednoho pacienta v rámci jednoho časového slotu (např. 8 dní měření) leží v paměti bezprostředně za sebou, což je optimální pro SIMD instrukce a cache procesoru.

### 3.2 Parsování časových razítek

Pro každý typ dat je nutné převést textové časové razítko na index časového slotu pro zarovnání dat. Pro každou frekvenci vzorkování je implementována samostatná funkce spolu s funkcí pro normalizaci času:

- **Frekvence 5 minut:** `minutes_since_midnight()`
- **Frekvence 1 sekunda:** `seconds_since_midnight()`
- **Frekvence 15 ms:** `units_15ms_since_midnight()`

#### 3.2.1 Frekvence 15ms

U vstupních dat je nutné převést textový čas na index slotu. Implementace využívá `std::from_chars` pro maximální rychlost a řeší i normalizaci milisekund.

```
1 std::optional<uint32_t> DataConverter::units_15ms_since_midnight(const std::string_view& ts) {
2     // Rychlá cesta: předpokládáme formát "YYYY-MM-DD HH:MM:SS.ms"
3     if (ts.size() >= 23) {
4         // ... parsování HH, MM, SS ...
5
6         // Helper funkce pro parsování milisekund
7         auto to_int = [](std::string_view s) -> std::optional<int> {
8             int out{};
```

```

9         auto res = std::from_chars(s.data(), s.data() + s.size(), out);
10        if (res.ec == std::errc{}) return out;
11        return std::nullopt;
12    };
13
14    int digit_count = 0;
15    // ... spočítání číslic v milisekundách -> digit_count
16
17    // Parsování milisekund a normalizace
18    if (digit_count > 0) {
19        auto sv_ms = ts.substr(20, digit_count);
20        auto oms = to_int(sv_ms);
21        if (oms) {
22            ms = *oms;
23            // ISO 8601 pravidlo: .5 je 500ms, .05 je 50ms
24            if (digit_count == 1) ms *= 100;
25            else if (digit_count == 2) ms *= 10;
26        }
27    }
28
29    // Výpočet celkových milisekund a slot indexu
30    uint32_t total_ms = (*oh * 3600 + *om * 60 + *os) * 1000 + ms;
31    return total_ms / 15;
32 }
33 return std::nullopt;
34 }

```

Listing 3.1: Metoda pro výpočet indexu slotu z času (DataReader.cpp)

### 3.3 Handling prázdných hodnot (Padding)

Jelikož jsou v měřeních prázdné úseky, kdy data chybí, používáme padding a vyplníme tyto prázdná data konstantou. Pro BVP data, která mohou nabývat záporných hodnot, nelze použít `-1.0f` a je tedy použito nejmenší možný float.

```

1 // Používáme velmi malé číslo, aby se nepletlo s reálnými daty
2 constexpr float EMPTY_VALUE = -1.0e30f;
3 constexpr float EPSILON = 100.0f; // Tolerance pro porovnání

```

Listing 3.2: Definice paddingu



# Kapitola 4

## Sekvenční zpracování (CPU)

Referenční implementace slouží pro ověření správnosti paralelních algoritmů. Iteruje přes časové sloty, pacienty a dny sekvenčně s potlačením automatické vektorizace/paralelizace kompilátorem.

```
1 // Při implemetaci a benchmarku byl použit clang-kompilátor
2 #pragma clang loop vectorize(enable) interleave(disable)
3 for (...) {}
```

Listing 4.1: Definice paddingu

### 4.1 Ukázka výpočtu mediánu

Následující kód ukazuje výpočet mediánu, který filtruje padding a řeší případ, kdy se sudým počtem prvků. V takovém případě je medián průměrem dvou prostředních hodnot.

```
1 inline float calculateMedian(const std::vector<float>& data) {
2     if (data.empty()) return EMPTY_VALUE;
3
4     // Filtrace: Vytvoříme buffer jen pro platná data
5     // Ignorujeme EMPTY_VALUE padding, ale bereme i záporná čísla (BVP)
6     std::vector<float> valid_data;
7     valid_data.reserve(data.size());
8
9     for (float val : data) {
10         // Podmínka pokud je hodnota větší než EMPTY_VALUE + odchylka floatu
11         if (val > EMPTY_VALUE + EPSILON) {
12             valid_data.push_back(val);
13         }
14     }
15
16     if (valid_data.empty()) return EMPTY_VALUE;
17
18     // Samotný výpočet mediánu (nth_element)
19     size_t n = valid_data.size() / 2;
20     std::nth_element(valid_data.begin(), valid_data.begin() + n, valid_data.
21 end());
22     float median = valid_data[n];
23
24     // Průměrování pro sudý počet prvků
25     // Pokud je počet prvků sudý, medián je průměr n-tého a (n-1)-tého prvku
26     // nth_element nám zaručil n-tý, pro (n-1)-tý musíme najít maximum v lev
27 é části.
28     if (valid_data.size() % 2 == 0) {
```

```

27     auto max_it = std::max_element(valid_data.begin(), valid_data.begin
    () + n);
28     median = (*max_it + median) / 2.0f;
29 }
30 return median;
31 }

```

Listing 4.2: Sekvenční výpočet mediánu s filtrací (CPUSquential.cpp)

## 4.2 Ukázka sekvenčního algoritmu

Následující kód ukazuje hlavní smyčku sekvenčního výpočtu mediánů a redukce časových slotů. Nejprve se spočítají mediány pro všechny pacienty napříč dny, poté se spočítá medián pro každé časové okno napříč pacienty a poté se provede downsampling. Výsledky jsou uloženy do `result_medians_seq` a `updated_timeslots_seq`.

```

1 void DexcomData::processSequential(int32_t num_wanted_time_slots) {
2 #pragma clang loop vectorize(enable) interleave(disable)
3     // smyčka napříč časovými sloty
4     for (auto i = 0; i < num_time_slots; i++) {
5         std::vector<float> patients_medians;
6         patients_medians.reserve(num_patients);
7         // smyčka napříč pacienty
8         for (auto j = 0; j < num_patients; j++) {
9             std::vector<float> vectorData;
10            vectorData.reserve(num_days);
11            auto data = getPatientDataPtr(i, j);
12            for (auto k = 0; k < num_days; k++) {
13                vectorData.push_back(*(data+k));
14            }
15            // medián pro pacienta j v časovém slotu i napříč dny
16            float median = calculateMedian(vectorData);
17            patients_medians.push_back(median);
18        }
19        // medián pro časový slot i napříč pacienty
20        float median = calculateMedian(patients_medians);
21        result_medians_seq.push_back(median);
22    }
23    // downsampling časových slotů, pokud je požadováno
24    // num_wanted_time_slots je počet cílových slotů po redukci
25    if (num_wanted_time_slots > 0) {
26        updated_timeslots_seq.reserve(num_wanted_time_slots);
27        std::vector<float> vectorData;
28        for (int i = 0; i < num_wanted_time_slots; ++i) {
29            // Výpočet rozsahu pro aktuální cílový slot
30            // Rozdíl v end_idx a start_idx určuje počet prvků ke zpracování
31            size_t start_idx = (static_cast<size_t>(i) * num_time_slots) /
num_wanted_time_slots;
32            size_t end_idx = (static_cast<size_t>(i + 1) * num_time_slots) /
num_wanted_time_slots;
33
34            if (end_idx > num_time_slots) end_idx = num_time_slots;
35            if (start_idx >= end_idx) {
36                updated_timeslots_seq.push_back(EMPTY_VALUE);
37                continue;
38            }
39
40            for (size_t k = start_idx; k < end_idx; ++k) {
41                vectorData.push_back(result_medians_seq[k]);

```

```

42     }
43     // výpočet mediánu pro aktuální cílový slot
44     float median = calculateMedian(vectorData);
45     updated_timeslots_seq.push_back(median);
46 }
47 }
48 }

```

Listing 4.3: Sekvenční výpočet mediánu s filtrací (CPUSequential.cpp)

## 4.3 Výběr algoritmu: `std::nth_element`

Pro výpočet mediánu byla zvolena funkce `std::nth_element` namísto klasického řazení `std::sort`. Důvodem je algoritmická složitost:

- `std::sort`: Složitost  $\mathcal{O}(N \log N)$ , seřadí celý vektor.
- `std::nth_element`: Složitost  $\mathcal{O}(N)$  v průměrném případě. Používá algoritmus Quickselect, který pouze přehází prvky tak, aby na  $n$ -té pozici byl prvek, který by tam patřil v seřazené posloupnosti.

Jelikož nás zajímá pouze jedna hodnota uprostřed (medián), je plné řazení zbytečné a `std::nth_element` je výrazně rychlejší.

# Kapitola 5

## Paralelizace a SIMD Optimalizace

Pro zvýšení výkonu na moderních procesorech byla využita kombinace vláknové paralelizace (OpenMP) a vektorizace (SIMD). V rámci projektu byly implementovány a změřeny dvě varianty paralelního zpracování na CPU. Obě využívají knihovnu OpenMP pro distribuci práce mezi vlákna, ale liší se v úrovni optimalizace výpočtu uvnitř jednotlivých vláken. Jedna metoda používá běžné sekvenční výpočty, zatímco druhá využívá ručně optimalizovaný kód pro výpočet mediánu z pevného počtu 8 hodnot pomocí SIMD instrukcí.

1. **Parallel (Bez SIMD):** Jednoduchá paralelizace smyčky, která pro výpočet mediánu používá alokaci `std::vector` a funkci `std::nth_element`.
2. **Parallel + SIMD:** Ručně optimalizovaná varianta pomocí SIMD instrukcí, která používá řazení přímo v registrech procesoru - Sorting Network.

### 5.1 Paralelizace algoritmu

Data pro různé časové sloty jsou na sobě nezávislá, takže může být snadno paralelizována. Každý časový slot tak může být vypočítán paralelně. Paralelizace probíhá ve dvou fázích:

1. **Fáze 1:** Výpočet mediánů pro jednotlivé TimeSloty. Smyčka napříč časovými sloty je rozdělena mezi vlákna pomocí `#pragma omp parallel for`. Každé vlákno zpracovává svůj blok časových slotů.
2. **Fáze 2:** Redukce (Downsampling) časových oken. Paralelizace přes cílové (redukované) sloty.

### 5.2 Varianta 1: Paralelizace bez SIMD

Tato metoda (reprezentovaná funkcí `processParallelCPUNonVectorized`) paralelizuje vnější smyčku přes časové sloty pomocí `#pragma omp parallel for`. Uvnitř vlákna se však chová stejně jako sekvenční kód: pro každého pacienta dynamicky alokuje vektor, naplní jej daty a zavolá obecnou funkci pro medián.

### 5.3 Varianta 2: Paralelizace se SIMD

Tato varianta (funkce `processParallelCPU`) se snaží o ruční optimalizaci pomocí vektorových instrukcí SIMD. Využívá faktu, že pro daný dataset máme maximálně 8 dní měření. To umožňuje řazení pomocí Sorting Network nad vektory o 4 floatech - `float32x4_t`. Pro každý den

se vytvoří jeden vektor a do jednoho vektoru se načtou hodnoty pro 4 pacienty. Jedna operace řazení tak zpracuje 4 pacienty najednou.

```

1 // Funkce pro transpozici 4x4 matice uložené ve 4 vektorech
2 inline void transpose_matrix(float32x4_t& r0, float32x4_t& r1, float32x4_t&
  r2, float32x4_t& r3) {
3     float32x4_t t0 = vtrn1q_f32(r0, r1);    // bere všechny sudé indexy
4     float32x4_t t1 = vtrn2q_f32(r0, r1);    // bere všechny liché indexy
5     float32x4_t t2 = vtrn1q_f32(r2, r3);
6     float32x4_t t3 = vtrn2q_f32(r2, r3);
7     r0 = vcombine_f32(vget_low_f32(t0), vget_low_f32(t2));
8     r1 = vcombine_f32(vget_low_f32(t1), vget_low_f32(t3));
9     r2 = vcombine_f32(vget_high_f32(t0), vget_high_f32(t2));
10    r3 = vcombine_f32(vget_high_f32(t1), vget_high_f32(t3));
11 }
12
13 void DexcomData::processParallelCPU(int32_t num_wanted_time_slots) {
14     // ... inicializace promenných a vektoru
15     const uint32_t aligned_patients = (num_patients / 4) * 4;
16     const uint32_t remainder = num_patients % 4;
17     const uint32_t padding_patients = (remainder == 0) ? 0 : (4 - remainder)
18     ;
19
20     for (auto i = 0; i < padding_patients; i++) {
21         for (auto j = 0; j < num_time_slots; j++) {
22             flat_data.push_back(-1);
23         }
24     }
25     const uint32_t patients_with_padding = aligned_patients +
padding_patients;
26     #pragma omp parallel for
27     {
28         for (uint32_t ts = 0; ts < num_time_slots; ++ts) {
29             for (uint32_t p = 0; p < patients_with_padding; p += 4) {
30                 float* ptr = getPatientDataPtr(ts, p);
31
32                 // Načtení dat pro 4 pacienty a 4 dny měření najednou do
vektorů
33                 // d0 obsahuje první 4 dny pro pacienta p
34                 float32x4_t d0 = vld1q_f32(ptr);
35                 // d1 obsahuje první 4 dny pro pacienta p+1
36                 float32x4_t d1 = vld1q_f32(ptr + 1*8);
37                 float32x4_t d2 = vld1q_f32(ptr + 2*8);
38                 float32x4_t d3 = vld1q_f32(ptr + 3*8);
39                 // Transpozice matic pro správné zarovnání
40                 // d0 obsahuje den 1 pro pacienty p, p+1, p+2, p+3
41                 transpose_matrix(d0, d1, d2, d3);
42
43                 // Načtení dat pro dny 5-8
44                 float32x4_t d4 = vld1q_f32(ptr + 4);
45                 float32x4_t d5 = vld1q_f32(ptr + 1*8 + 4);
46                 float32x4_t d6 = vld1q_f32(ptr + 2*8 + 4);
47                 float32x4_t d7 = vld1q_f32(ptr + 3*8 + 4);
48                 // Transpozice matic pro správné zarovnání
49                 transpose_matrix(d4, d5, d6, d7);
50
51                 // Výpočet mediánu pro 4 pacienty najednou ...
52             }
53         }
54     }
55 }

```

54 }

Listing 5.1: Inicializace SIMD vektorů pro paralelní výpočet mediánu (CPUParallel.cpp)

### 5.3.1 Princip SIMD řazení - Sorting Network

Namísto obecného řazení používáme tzv. **Sorting Network** (řadicí síť). Jde o algoritmus s pevně daným počtem porovnání a prohození, který nezávisí na hodnotách dat (data-independent). Hledání využívá paralelní porovnání a prohození hodnot v registrech procesoru "**Compare and Swap**" pomocí instrukcí `vminq_f32` a `vmaxq_f32`.

```
1 #define SORT_VEC(A, B) { \
2     float32x4_t min_v = vminq_f32(A, B); \
3     float32x4_t max_v = vmaxq_f32(A, B); \
4     A = min_v; \
5     B = max_v; \
6 }
```

Listing 5.2: Makro pro porovnání a prohození dvou SIMD vektorů (CPUParallel.cpp)

Seřazení 8 hodnot pro 4 pacienty probíhá ve 21 krocích (Compare and Swap).

```
1 void DexcomData::processParallelCPU(int32_t num_wanted_time_slots) {
2     for (uint32_t ts = 0; ts < num_time_slots; ++ts) {
3         for (uint32_t p = 0; p < patients_with_padding; p += 4) {
4             // Načtení a transpozice dat (viz předchozí kód) ...
5
6             // Ruční řazení 8 hodnot pomocí Sorting Network
7             SORT_VEC(d0, d1); SORT_VEC(d2, d3); SORT_VEC(d4, d5); SORT_VEC(
8             d6, d7);
9             SORT_VEC(d0, d2); SORT_VEC(d1, d3); SORT_VEC(d4, d6); SORT_VEC(
10            d5, d7);
11            SORT_VEC(d1, d2); SORT_VEC(d5, d6); SORT_VEC(d0, d4); SORT_VEC(
12            d3, d7);
13            SORT_VEC(d1, d5); SORT_VEC(d2, d6);
14            SORT_VEC(d1, d4); SORT_VEC(d3, d6);
15            SORT_VEC(d2, d4); SORT_VEC(d3, d5);
16            SORT_VEC(d3, d4); SORT_VEC(d1, d2); SORT_VEC(d5, d6);
17
18            // Filtrace paddingu a extrakce mediánu z registrových vektorů
19            ...
20        }
21    }
```

Listing 5.3: Makro pro porovnání a prohození dvou SIMD vektorů (CPUParallel.cpp)

### 5.3.2 Extrakce mediánu ze seřazených registrů

Po seřazení pomocí Sorting Network jsou hodnoty v registrech uspořádány od nejmenšího k největšímu. Následuje filtrace paddingu (`EMPTY_VALUE`) a výpočet mediánu. Počet platných hodnot se spočítá pomocí porovnání s `EMPTY_VALUE + EPSILON` a na základě toho se určí pozice mediánu.

```
1 // Filtrace paddingu a počítání platných hodnot po Sorting Network - viz
2 // výše
3 float32x4_t limit = vdupq_n_f32(EMPTY_VALUE + EPSILON);
4 uint32x4_t counts = vdupq_n_u32(0);
5
6 // Sčítáme masky pro nalezení neplatných hodnot
```

```

6   counts = vsubq_u32(counts, vcgtq_f32(d0, limit));
7   counts = vsubq_u32(counts, vcgtq_f32(d1, limit));
8   counts = vsubq_u32(counts, vcgtq_f32(d2, limit));
9   counts = vsubq_u32(counts, vcgtq_f32(d3, limit));
10  counts = vsubq_u32(counts, vcgtq_f32(d4, limit));
11  counts = vsubq_u32(counts, vcgtq_f32(d5, limit));
12  counts = vsubq_u32(counts, vcgtq_f32(d6, limit));
13  counts = vsubq_u32(counts, vcgtq_f32(d7, limit));
14
15  uint32_t cnt[4];
16  vst1q_u32(cnt, counts);
17
18  // extrakce seřazených hodnot do pole
19  float sorted_cols[32];
20  vst1q_f32(sorted_cols + 0, d0);
21  vst1q_f32(sorted_cols + 4, d1);
22  vst1q_f32(sorted_cols + 8, d2);
23  vst1q_f32(sorted_cols + 12, d3);
24  vst1q_f32(sorted_cols + 16, d4);
25  vst1q_f32(sorted_cols + 20, d5);
26  vst1q_f32(sorted_cols + 24, d6);
27  vst1q_f32(sorted_cols + 28, d7);
28
29  // Výpočet mediánu pro každého z 4 pacientů
30  for(int k=0; k<4; ++k) {
31      int n = cnt[k];
32      float median = EMPTY_VALUE;
33      if(n > 0) {
34          int start = 8 - n;
35          int mid = n / 2;
36
37          int idx1 = (start + mid) * 4 + k;
38          median = sorted_cols[idx1];
39
40          // Sudý počet prvků - průměr dvou prostředních hodnot
41          if(n % 2 == 0) {
42              int idx2 = (start + mid - 1) * 4 + k;
43              median = (sorted_cols[idx2] + median) / 2.0f;
44          }
45          grand_median_buffer.push_back(median);
46      }
47      result_medians_par_per_pat[(ts * num_patients) + p + k] = median;
48  }

```

Listing 5.4: Extrakce mediánu ze seřazených SIMD registrů (CPUParallel.cpp)

## 5.4 Medián napříč pacienty

Optimální algoritmus pro vektorové řazení se mi nepodařilo implementovat. Níže je ukázka implementace bitonic sort pro 16 prvků, ale to neseřadilo hodnoty napříč všemi vektory - jen v rámci jednoho vektoru o 4 prvcích. Z maticového pohledu to seřadilo sloupce i řádky.

```

1  inline void bitonic_sort_16(float32x4_t& q0, float32x4_t& q1, float32x4_t&
    q2, float32x4_t& q3) {
2
3      SORT_VEC(q0, q1); SORT_VEC(q2, q3); SORT_VEC(q0, q2); SORT_VEC(q1, q3);
4      SORT_VEC(q0, q3); SORT_VEC(q1, q2);
5
6      transpose_matrix(q0, q1, q2, q3);

```

```

7
8     SORT_VEC(q0, q1); SORT_VEC(q2, q3); SORT_VEC(q0, q2); SORT_VEC(q1, q3);
9     SORT_VEC(q0, q3); SORT_VEC(q1, q2);
10
11     transpose_matrix(q0, q1, q2, q3);
12 }

```

Listing 5.5: Bitonic Sort pro 16 prvků ve 4 vektorech (CPUParallel.cpp)

Proto je v této fázi využita standardní funkce `std::nth_element`, která sice není vektorizovaná, ale celá smyčka přes časové sloty je stále paralelizována pomocí OpenMP.

```

1     // Výpočet Grand Median pro každý časový slot
2     float gm = EMPTY_VALUE;
3     if (!grand_median_buffer.empty()) {
4         size_t n = grand_median_buffer.size() / 2;
5         std::nth_element(grand_median_buffer.begin(), grand_median_buffer.
6         begin() + n, grand_median_buffer.end());
7         gm = grand_median_buffer[n];
8
9         if (grand_median_buffer.size() % 2 == 0) {
10             auto max_it = std::max_element(grand_median_buffer.begin(),
11             grand_median_buffer.begin() + n);
12             gm = (*max_it + gm) / 2.0f;
13         }
14     }
15 }

```

Listing 5.6: Bitonic Sort pro 16 prvků ve 4 vektorech (CPUParallel.cpp)

## 5.5 Downsampling časových slotů

Downsampling časových slotů je implementován paralelně pomocí OpenMP. Každé vlákno zpracovává jeden cílový (redukovaný) časový slot. Pro každý cílový slot se vypočítá rozsah původních slotů, které do něj spadají, a z těchto hodnot se spočítá medián pomocí `std::nth_element`. Implementace je bez manuální SIMD optimalizace.

```

1     #pragma omp for schedule(static)
2     for (int i = 0; i < num_wanted_time_slots; ++i) {
3         size_t start_idx = (static_cast<size_t>(i) * num_time_slots) /
4         num_wanted_time_slots;
5         size_t end_idx = (static_cast<size_t>(i + 1) * num_time_slots) /
6         num_wanted_time_slots;
7
8         if (end_idx > num_time_slots) end_idx = num_time_slots;
9         size_t count = end_idx - start_idx;
10        if (count <= 0) {
11            updated_timeslots_par[i] = (start_idx < num_time_slots) ?
12            result_medians_par[start_idx] : EMPTY_VALUE;
13            continue;
14        }
15
16        // Lokální buffer pro hodnoty v aktuálním rozsahu
17        std::vector<float> local_buffer;
18        // ... naplnění local_buffer hodnotami z result_medians_par ...
19
20        float median = EMPTY_VALUE;
21        if (!local_buffer.empty()) {
22            auto n = local_buffer.size() / 2;
23            std::nth_element(local_buffer.begin(), local_buffer.begin() + n,
24            local_buffer.end());
25        }
26    }

```



```

21         median = local_buffer[n];
22
23         if (n % 2 == 0) {
24             auto max_it = std::max_element(local_buffer.begin(),
local_buffer.begin() + n);
25             median = (*max_it + median) / 2.0f;
26         }
27     }
28     updated_timeslots_par[i] = median;
29 }

```

Listing 5.7: Paralelní downsampling časových slotů (CPUParallel.cpp)

# Kapitola 6

## GPU (OpenCL)

Pro zpracování větších dat (zejména BVP s **740 milionů záznamů** (5.7 milionu slotů  $\times$  16 pacientů  $\times$  8 dní)) byla implementována pipeline skládající se ze 3 kernelů. Rozdělení do tří kernelů bylo zvoleno z důvodu synchronizace paralelních operací. Implementace je realizována v OpenCL.

Celá pipeline je kontrolována z C++ hostitelského kódu, který připraví data, nastaví OpenCL prostředí, nahraje data na GPU, spustí jednotlivé kernely a stáhne výsledky zpět na CPU - `GPUParallel.cpp`. Mezi jednotlivými kernely zůstávají data uložena v globální paměti GPU a zpět do paměti programu jsou přenesena jen finální výsledky po skončení celého výpočtu.

```
1 void DexcomData::processGPU(int32_t num_wanted_time_slots, bool
    read_all_outputs, int8_t num_kernels_use, const std::string& kernel_file)
    {
2     // nastavení OpenCL prostředí, načtení kernelů z kernel_file ...
3
4     // alokace bufferů v globální paměti GPU
5     // input data
6     cl_mem d_data = clCreateBuffer(context, CL_MEM_READ_ONLY |
7         CL_MEM_COPY_HOST_PTR, flat_data.size() * sizeof(float),
8         flat_data.data(), &err);
9     checkErr(err, "CreateBuffer Input Data");
10
11    // medians per patients
12    cl_mem d_medians_per_patients_results = clCreateBuffer(context,
13        CL_MEM_WRITE_ONLY, total_items * sizeof(float), NULL, &err);
14    checkErr(err, "CreateBuffer medians per patients");
15
16    // medians per time_slots
17    cl_mem d_time_slots_results = clCreateBuffer(context, CL_MEM_READ_WRITE,
18        num_time_slots * sizeof(float), NULL, &err);
19    checkErr(err, "CreateBuffer medians per time_slots");
20
21    size_t total_items = num_time_slots * num_patients;
22
23    // kernel 1: Výpočet mediánu napříč dny pro každého pacienta
24    cl_kernel kernel = clCreateKernel(program, "median_kernel", &err);
25
26    // nastavení argumentů kernelu 1 ...
27    clSetKernelArg(kernel, 0, sizeof(cl_mem), &d_data);
28    clSetKernelArg(kernel, 1, sizeof(cl_mem),
29        &d_medians_per_patients_results);
30    clSetKernelArg(kernel, 2, sizeof(int), &n_days);
31    clSetKernelArg(kernel, 3, sizeof(int), &n_total);
32    size_t global_work_size = total_items;
33
```

```

34 err = clEnqueueNDRangeKernel(queue, kernel, 1, NULL, &global_work_size,
35 NULL, 0, NULL, NULL);
36
37 // kernel 2: Výpočet mediánu napříč pacienty pro každý časový slot
38 if (num_kernels_use >= 2) {
39     cl_kernel timeslots_kernel = clCreateKernel(program,
40         "reduce_patients_kernel", &err);
41
42     // nastavení argumentů kernelu 2 ...
43     clSetKernelArg(timeslots_kernel, 0, sizeof(cl_mem),
44         &d_medians_per_patients_results);
45     clSetKernelArg(timeslots_kernel, 1, sizeof(cl_mem),
46         &d_time_slots_results);
47     clSetKernelArg(timeslots_kernel, 2, sizeof(int), &num_patients);
48
49     size_t global_work_size2 = num_time_slots;
50
51     err = clEnqueueNDRangeKernel(queue, timeslots_kernel,
52         1, NULL, &global_work_size2, NULL, 0, NULL, NULL);
53
54     // release kernel 2
55 }
56 // kernel 3: Downsampling časových slotů
57 if (num_kernels_use >= 3 && num_wanted_time_slots > 0) {
58
59     // buffer pro výsledky downsamplingu
60     cl_mem d_wanted_time_slots_results = clCreateBuffer(context,
61         CL_MEM_READ_WRITE, num_wanted_time_slots * sizeof(float),
62         NULL, &err);
63     checkErr(err, "CreateBuffer medians per wanted_time_slots");
64
65     cl_kernel reduce_slots_kernel = clCreateKernel(program,
66         "reduce_slots_kernel", &err);
67
68     // nastavení argumentů kernelu 3 ...
69     clSetKernelArg(reduce_slots_kernel, 0, sizeof(cl_mem),
70         &d_time_slots_results);
71     clSetKernelArg(reduce_slots_kernel, 1, sizeof(cl_mem),
72         &d_wanted_time_slots_results);
73     clSetKernelArg(reduce_slots_kernel, 2, sizeof(int),
74         &num_time_slots);
75     clSetKernelArg(reduce_slots_kernel, 3, sizeof(int),
76         &num_wanted_time_slots);
77
78     size_t global_work_size3 = num_wanted_time_slots;
79
80     err = clEnqueueNDRangeKernel(queue, reduce_slots_kernel,
81         1, NULL, &global_work_size3, NULL, 0, NULL, NULL);
82
83     // stáhnutí výsledků z GPU
84     rr = clEnqueueReadBuffer(queue, d_wanted_time_slots_results,
85         CL_TRUE, 0, num_wanted_time_slots * sizeof(float),
86         updated_timeslots_gpu.data(), 0, NULL, NULL);
87
88     // release kernel 3
89 }
90 // release ostatních zdrojů GPU ...
91
92 }

```

## 6.1 Kernel 1: Median per Patient

Tento kernel počítá medián napříč 8 dny pro každého pacienta v každém časovém okně.

- **Vstup:** Raw data [TimeSlots  $\times$  Patients  $\times$  Days]
- **Implementace:** Každé vlákno (Work-item) zpracovává jednoho pacienta v jednom čase. Protože počet dnů je fixní (8), data se načtou do privátních registrů a seřadí pomocí **Sorting Network** (série instrukcí CAS).

```
1 __kernel void median_kernel(  
2     __global const float* data,  
3     __global float* result_medians,  
4     const int num_days,  
5     const int total_patients_slots  
6 ) {  
7     int gid = get_global_id(0);  
8     long offset = (long)gid * (long)num_days;  
9  
10    // Načtení 8 dnů do registrů  
11    float v0 = data[offset + 0];  
12    float v1 = data[offset + 1];  
13    float v2 = data[offset + 2];  
14    float v3 = data[offset + 3];  
15    float v4 = data[offset + 4];  
16    float v5 = data[offset + 5];  
17    float v6 = data[offset + 6];  
18    float v7 = data[offset + 7];  
19  
20    // Sorting Network (Compare-And-Swap)  
21    CAS(v0, v1); CAS(v2, v3); CAS(v4, v5); CAS(v6, v7);  
22    // ... další kroky sítě ...  
23  
24    // Filtrace a výpočet mediánu  
25    int valid_count = 0;  
26    valid_count += (v0 > VALID_THRESHOLD);  
27    // ...  
28  
29    if (valid_count > 0) {  
30        // ... výběr mediánu z pole 'sorted' a uložení do result_medians ...  
31    }  
32 }
```

Listing 6.1: Median Kernel s CAS sítí (kernel.cl)

## 6.2 Kernel 2: Reduce Patients

Tento kernel počítá medián napříč všemi pacienty pro každý časový slot na základě výstupu z prvního kernelu.

- **Vstup:** Výstup z Kernelu 1 (mediány jednotlivých pacientů).
- **Úkol:** Pro každý časový slot spočítat medián ze všech pacientů.
- **Důvod oddělení:** Výpočet mediánu přes pacienty vyžaduje přístup k datům, která byla vyprodukována různými vlákny v Kernelu 1. Spojení do jednoho kernelu by vyžadovalo globální bariéry nebo atomické operace, což by degradovalo výkon. Spuštěním nového kernelu je synchronizace zaručena globálně.

- **Implementace:** Každé vlákno načte mediány všech pacientů pro daný čas do lokálního bufferu a provede řazení.

```

1 __kernel void reduce_patients_kernel(
2     __global const float* input_matrix,
3     __global float* output_medians,
4     const int num_patients
5 ) {
6     int ts = get_global_id(0); // time_slot
7
8     #define MAX_BUFFER_SIZE 16
9     float buffer[MAX_BUFFER_SIZE];
10
11     int valid_count = 0;
12     long offset = (long)ts * (long)num_patients;
13
14     for (int p = 0; p < num_patients; ++p) {
15         float val = input_matrix[offset + p];
16         if (val > VALID_THRESHOLD && valid_count < MAX_BUFFER_SIZE) {
17             buffer[valid_count++] = val;
18         }
19     }
20
21     float result = EMPTY_VALUE;
22     if (valid_count > 0) {
23         sort(buffer, valid_count);
24
25         int mid = valid_count / 2;
26         result = buffer[mid];
27
28         // sudý počet prvků - průměr dvou prostředních hodnot
29         if (valid_count % 2 == 0) {
30             result = (buffer[mid - 1] + result) / 2.0f;
31         }
32     }
33
34     output_medians[ts] = result;
35 }

```

Listing 6.2: reduce\_patients\_kernel pro medián napříč pacienty (kernel.cl)

Pro řazení jsou implementované 2 funkce - **Bubble Sort** a **Shell Sort**. Bubble Sort by měl být vhodnější pro malý počet pacientů (16), zatímco Shell Sort je efektivnější pro větší množiny dat. V další kapitole jsou změřeny výkony obou přístupů. Shell sort provádí řazení přímo v globální paměti bez alokace další paměti, zatímco Bubble Sort využívá lokální buffer.

```

1 void shell_sort_global(__global float* data, int start_idx, int count) {
2     for (int gap = count / 2; gap > 0; gap /= 2) {
3         for (int i = gap; i < count; i += 1) {
4             float temp = data[start_idx + i];
5             int j;
6             for (j = i; j >= gap && data[start_idx + j - gap] > temp; j -=
gap) {
7                 data[start_idx + j] = data[start_idx + j - gap];
8             }
9             data[start_idx + j] = temp;
10        }
11    }
12 }
13
14 void sort(float* arr, int n) {

```

```

15     for (int i = 0; i < n - 1; i++) {
16         for (int j = 0; j < n - i - 1; j++) {
17             if (arr[j] > arr[j + 1]) {
18                 float temp = arr[j];
19                 arr[j] = arr[j + 1];
20                 arr[j + 1] = temp;
21             }
22         }
23     }
24 }

```

Listing 6.3: Řazení: Bubble Sort a Shell Sort (kernel.cl)

## 6.3 Kernel 3: Downsampling (Shell Sort)

Tento kernel provádí downsampling časových slotů na základě výstupu z druhého kernelu. Jako parametr přijímá počet cílových časových slotů a pro každý cílový slot spočítá medián z příslušných původních slotů.

- **Vstup:** Výstup z Kernelu 2 (jeden medián pro každý časový slot).
- **Úkol:** Spočítat medián z  $K$  po sobě jdoucích časových slotů a tím provést downsampling.
- **Implementace:** Protože  $K$  může být velké (tisíce hodnot), nelze použít registry. Kernel využívá **In-Place Shell Sort** přímo v globální paměti grafické karty.

```

1  __kernel void reduce_slots_kernel(
2      __global float* input_slots,
3      __global float* output_final,
4      const int num_total_slots,
5      const int num_wanted_slots
6  ) {
7      // 1. Komprese dat (přesun validních hodnot na začátek úseku)
8      int valid_count = 0;
9      int write_ptr = start_idx;
10     for (int read_ptr = start_idx; read_ptr < end_idx; ++read_ptr) {
11         if (input_slots[read_ptr] > VALID_THRESHOLD) {
12             if (read_ptr != write_ptr)
13                 input_slots[write_ptr] = input_slots[read_ptr];
14             write_ptr++;
15             valid_count++;
16         }
17     }
18
19     // 2. Řazení přímo v globální paměti (bez alokace)
20     if (valid_count > 0) {
21         shell_sort_global(input_slots, start_idx, valid_count);
22         // ... výběr mediánu ...
23     }
24 }

```

Listing 6.4: Downsampling s Shell Sortem

## 6.4 Řešení pro BVP data (740 milionů záznamů)

Při zpracování BVP dat se objevily specifické problémy spojené s jejich objemem:

- **Přetečení indexů:** Počet prvků ( $5.7 \times 10^6 \times 16 \times 8 \approx 740 \times 10^6$ ) se vejde do 32-bitového `int`, ale při výpočtu indexů (např. násobení) může dojít k přetečení mezivýsledku. V kernelech byla proto důsledně zavedena 64-bitová aritmetika (`long`) pro adresaci paměti.
- **Filtrace:** BVP data mohou být záporná. Byla zavedena konstanta `EMPTY_VALUE = -1e30` pro padding, aby se odlišila chybějící data od platných záporných hodnot.

# Kapitola 7

## Naměřené výsledky a Vyhodnocení

V této kapitole jsou prezentovány výsledky benchmarků pro různé sady dat. Měření byla prováděna opakovaně a výsledky byly validovány křížovou kontrolou (Sequential vs. Parallel vs. GPU), aby bylo zaručeno, že jsou výsledky shodné (s definovanou tolerancí).

### 7.1 Hardwarová konfigurace

Měření probíhalo na následující konfiguraci:

- **CPU:** M1 Pro (8 jader)
- **GPU:** M1 Pro (14 jader)
- **RAM:** 16GB unifikované paměti
- **OS:** macOS Tahoe 26.2

### 7.2 Časy načítání dat

Vstupem pro algoritmy jsou CSV soubory. Měření zahrnuje čas potřebný pro přečtení souborů z disku a jejich konverzi do `flat_data` formátu. Kategorie dat DEXCOM, HR a BVP mají různé velikosti souborů kvůli rozdílným vzorkovacím intervalům.

Tabulka 7.1: Časy načítání vstupních souborů (15 souborů pro typ dat DEXCOM, 16 souborů pro typy dat HR a BVP)

Typ Dat (vzorkovací interval)	Průměrná velikost jednoho souboru	Čas načítání [ms]
Dexcom (5 min)	139 KB	166 ms
HR (1 sec)	15.8 MB	18 435 ms (18s)
BVP (15 ms)	1.28 GB	1 568 830 ms (26m 9s)

### 7.3 Dexcom Data (Nízká zátěž)

- **Charakteristika:** Vzorkování po 5 minutách (288 slotů/den).
- **Velikost:** Malá data, minimální paměťová náročnost.



Tabulka 7.2: Výsledky pro Dexcom Data

Metoda	Čas [ms]
CPU Sequential	5.45483 ms
CPU Parallel + SIMD	0.8895 ms
CPU Parallel (No Vectorization)	4.79196 ms
GPU (Full Pipeline)	59.51 ms

**Zhodnocení:** U malých dat je režie spojená s přenosem dat na GPU (PCIe latence) a spouštěním kernelů větší než samotný výpočet. Zde dominuje CPU Parallel verze díky vektorovým instrukcím, efektivnímu využití cache a absence overheadu.

## 7.4 HR Data (Střední zátěž)

- **Charakteristika:** Vzorkování po 1 sekundě (86 400 slotů/den).
- **Velikost:** Středně velká data.

Tabulka 7.3: Výsledky pro HR Data

Metoda	Čas [ms]
CPU Sequential	1544.08 ms
CPU Parallel + Vectorized	226.49 ms
CPU Parallel (No Vectorization)	1375.42 ms
GPU (Full Pipeline)	19.66 ms

**Zhodnocení:** S rostoucím objemem dat je rychlost GPU znatelnější. Paralelní CPU se SIMD instrukcemi stále nabízí dobrý výkon, ale pravděpodobně je limitována implementací, která vektorizuje jen první část výpočtu (medián napříč dny). GPU verze těží z paralelizace napříč všemi částmi výpočtu.

Zajímavé je, že GPU zde běží kratší dobu než GPU u Dexcom dat. Důvodem je, že GPU běží první nad Dexcom daty a následně na HR daty. Režie inicializace OpenCL a JIT kompilace kernelu se započítá do prvního běhu a při druhém běhu již není potřeba.

## 7.5 BVP Data (Vysoká zátěž)

- **Charakteristika:** Vzorkování po 15 ms (cca 5 760 000 slotů/den).
- **Velikost:** Velká data, miliony záznamů. Zde se projevují limity paměti a cache.
- **Specifika:** Data obsahují záporné hodnoty, což vyžadovalo robustní handling paddingu (EMPTY\_VALUE = -1e30).

Tabulka 7.4: Výsledky pro BVP Data

Metoda	Čas [ms]
CPU Sequential	80440.40 ms
CPU Parallel + Vectorized	16506.80 ms
CPU Parallel (No Vectorization)	68084.00 ms
GPU (Full Pipeline)	1583.48 ms

**Zhodnocení:** Toto je scénář, pro který byla GPU optimalizace navržena. GPU je více než 36x rychlejší než sekvenční CPU a více než 7.5x rychlejší než paralelní CPU s vektorizací. Velká data plně využívají paralelní architekturu GPU a minimalizují režii spojenou s přenosem dat díky pipeline.

## 7.6 Analýza fází GPU výpočtu

Pro detailní pochopení výkonu GPU bylo provedeno měření po jednotlivých krocích (kernelech). To umožňuje identifikovat úzká hrdla pipeline.

- **1 Kernel:** Spuštění pouze `median_kernel` (výpočet mediánů pro každého pacienta a časový slot). Zpět do RAM se přenáší kompletní matice výsledků (časové sloty  $\times$  pacienti).
- **2 Kernely:** Navíc spuštění `reduce_patients_kernel`. Redukce pacientů na GPU, do RAM se přenáší jen jeden medián pro každý časový slot.
- **3 Kernely (Full):** Kompletní pipeline včetně `reduce_slots_kernel` (downsampling). Zpět do RAM se přenáší pouze finální downsampled výsledky.

Tabulka 7.5: Časy GPU podle počtu aktivních kernelů (BVP data)

Konfigurace GPU Pipeline	Čas [ms]
Pouze Kernel 1 (Median per Patient)	1094.36
Kernel 1 + Kernel 2 (Reduce Patients)	534.659
Kernel 1 + 2 + 3 (Full Pipeline)	1024.54

**Analýza:** Čas pro "1 Kernel"(1094 ms) je dvojnásobný oproti "2 Kernelům"(535 ms). To indikuje **Memory Bottleneck** v přenosu dat z GPU. Výstupem prvního kernelu je cca 368 MB dat (92 milionů hodnot), které se musí přenést přes sběrnici do RAM. Při použití dvou kernelů zůstávají tato data v paměti GPU a přenáší se pouze výsledek po 2. kernelu (cca 23 MB). Časová úspora na přenosu dat zde výrazně převyšuje výpočetní čas druhého kernelu.

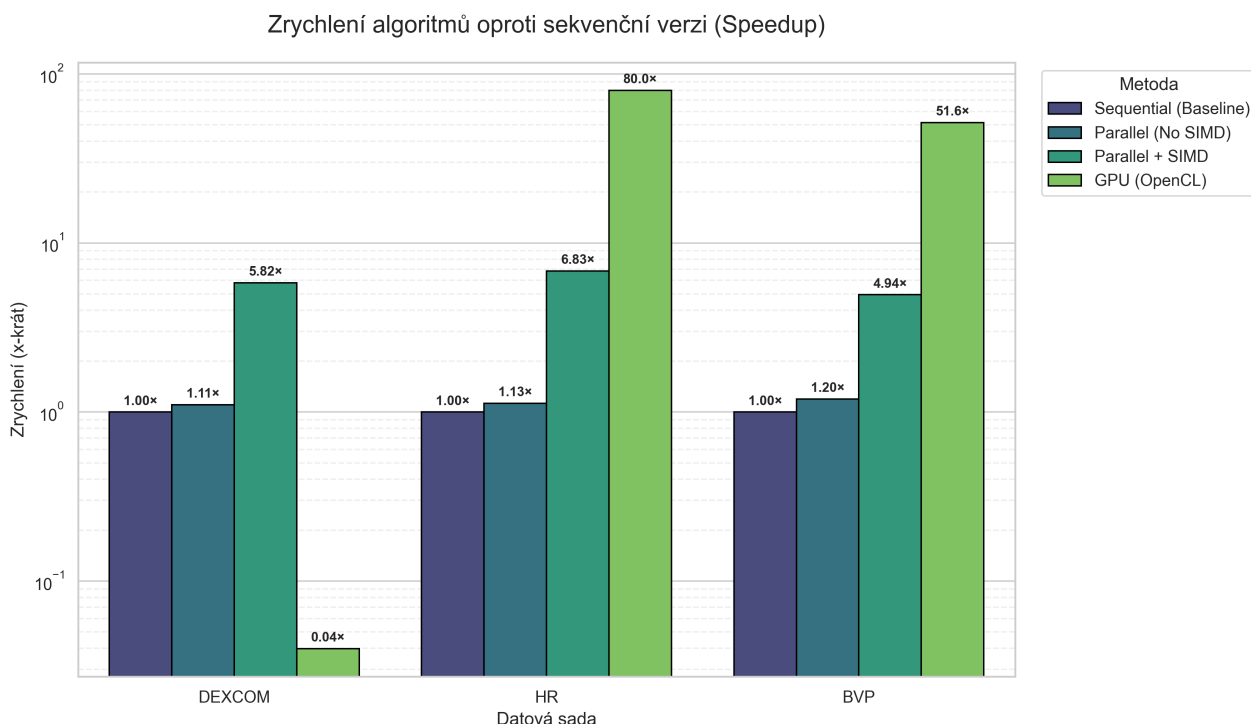
# Kapitola 8

## Analýza a Diskuze výsledků

Cílem této kapitoly je vyhodnotit efektivitu implementovaných algoritmů na základě naměřených časů. Zaměřujeme se na porovnání absolutního výkonu (škálovatelnost) a efektivity paralelizace (Karp-Flattova metrika).

### 8.1 Zrychlení a efektivita (Speedup)

Zrychlení ( $S$ ) je definováno jako poměr času sekvenčního algoritmu k času paralelního algoritmu:  $S = T_{seq}/T_{par}$ . Z naměřených dat je patrné, že charakteristika zrychlení se dramaticky mění s velikostí dat.



Obrázek 8.1: Graf zrychlení (Speedup) jednotlivých metod oproti sekvenční verzi (osa Y je logaritmická). U malých dat (Dexcom) je GPU neefektivní, zatímco u velkých (BVP) dominuje se zrychlením přes 50×.

Z grafu 8.1 vyplývají tři klíčové poznatky:

- Režie GPU u malých dat:** U sady Dexcom dosahuje GPU zrychlení  $S < 1$  (zpomalení). Režie spojená s přenosem dat a spouštěním kernelu zde převyšuje samotný výpočet.

2. **Dominance GPU u větších dat:** U sady HR a BVP je GPU výrazně efektivnější. Pro BVP dosahuje GPU zrychlení  $51,6\times$ .
3. **Role SIMD instrukcí:** U CPU implementace je viditelný rozdíl mezi čistou paralelizací (NoVect) a vektorizací. Zatímco ‘Parallel (No SIMD)’ dosahuje u BVP zrychlení pouze  $1,2\times$ , optimalizovaná verze ‘Parallel + SIMD’ dosahuje zrychlení  $4,9\times$ .

## 8.2 Karp-Flattova metrika

Pro hlubší analýzu efektivity paralelizace na CPU byla použita Karp-Flattova metrika ( $e$ ), která odhaduje sériový podíl kódu včetně režie paralelizace. Výpočet byl proveden pro procesor Apple M1 Pro s počtem jader  $p = 8$ .

$$e = \frac{\frac{1}{S} - \frac{1}{p}}{1 - \frac{1}{p}}$$

Nízká hodnota  $e$  značí efektivní paralelizaci. Hodnoty v Tabulce 8.1 ukazují, jak efektivně dokáže algoritmus využít dostupná jádra.

Tabulka 8.1: Karp-Flattova metrika pro CPU implementace ( $p = 8$ )

Dataset	Metoda	Speedup ( $S$ )	Karp-Flatt ( $e$ )	Efektivita
Dexcom	CPU Par + SIMD	5,82	0,053	Vysoká
HR	CPU Par + SIMD	6,83	0,024	Velmi vysoká
BVP	CPU Par + SIMD	4,94	0,089	Střední
<b>BVP</b>	<b>CPU Par (No SIMD)</b>	<b>1,19</b>	<b>0,817</b>	<b>Kritická</b>

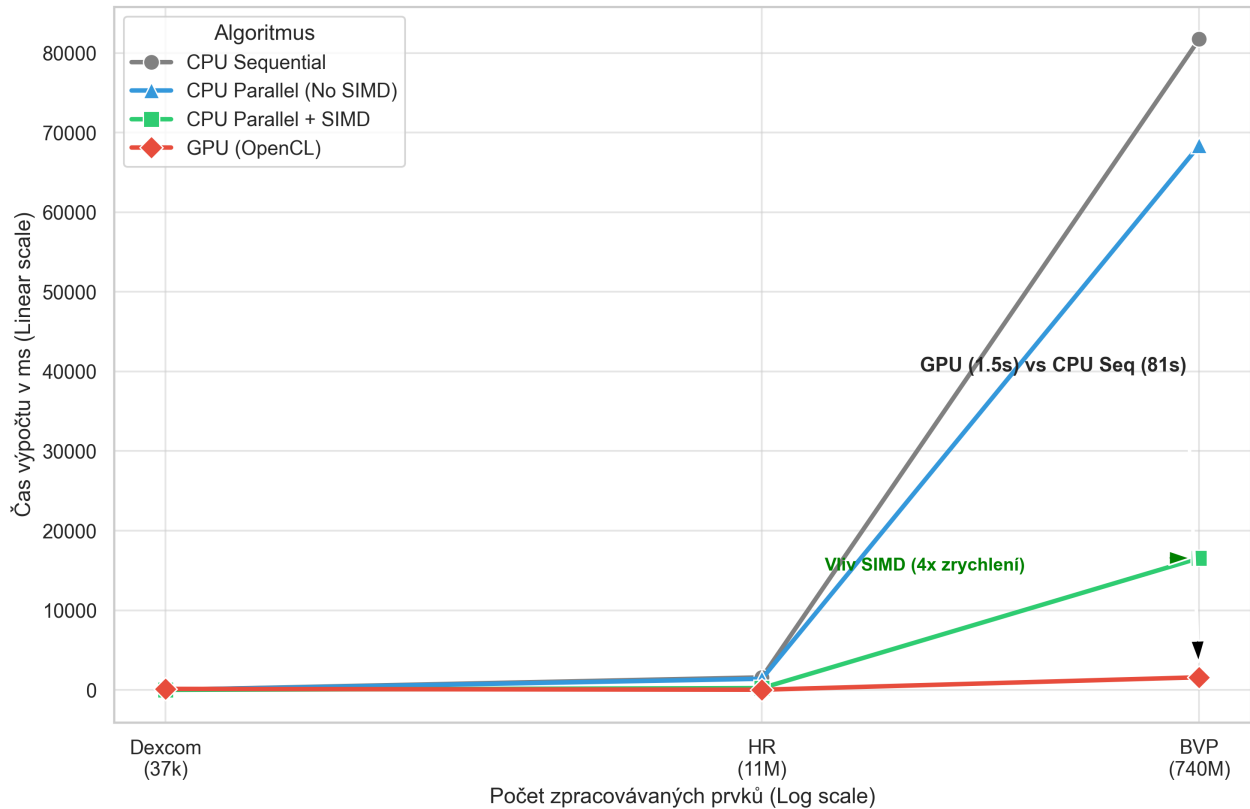
### Interpretace:

- U sady **HR** dosahuje metrika hodnoty 0,024, což znamená, že pouze 2,4 % času tvoří sériová část nebo režie. Paralelizace je zde téměř ideální.
- U sady **BVP** (s vektorizací) hodnota stoupá na 0,089 (8,9 % režie). To je pravděpodobně způsobeno saturací paměťové sběrnice (Memory Bound), kdy 8 jader soutěží o přístup k 1,2 GB dat.
- Metoda **BVP bez SIMD** vykazuje  $e = 0,817$ . To znamená, že 81,7 % výpočtu se chová sériově. Bez SIMD instrukcí, které zpracují 4 hodnoty najednou a efektivně využívají registry, je CPU zahlceno režii správy vláken a čekáním na data, což degraduje výkon téměř na úroveň sekvenčního kódu.

## 8.3 Škálovatelnost v čase

Poslední graf zobrazuje absolutní čas běhu v závislosti na velikosti vstupu.

## Kompletní porovnání škálovatelnosti



Obrázek 8.2: Závislost doby výpočtu na velikosti vstupních dat (log-linear měřítko).

Graf 8.2 ilustruje, že zatímco sekvenční a čistě paralelní řešení rostou lineárně s velikostí dat (strmý nárůst u BVP), GPU řešení (červená křivka) si zachovává nízký čas zpracování i pro 740 milionů záznamů. Zelená křivka (SIMD) představuje ideální volbu pro středně velká data, kde je rychlejší než GPU (díky absenci přenosu dat), ale u větších dat již na GPU ztrácí.

# Kapitola 9

## Diskuze nad implementací

Projekt demonstruje efektivní zpracování velkých datových sad pomocí různých přístupů - sekvenčního CPU, paralelního CPU s vektorizací a GPU akcelerační. Každý přístup má své výhody a nevýhody v závislosti na velikosti dat a dostupném hardwaru. Při větších datových sadách (BVP) je GPU akcelerační jasně nejefektivnější, zatímco pro menší datové sady (Dexcom) může být optimalizovaná paralelní verze na CPU výhodnější díky nižší režii.

Práce úspěšně řeší několik technických výzev, včetně efektivního využití SIMD instrukcí na CPU a implementace robustní pipeline na GPU. Nicméně, existují oblasti pro další zlepšení:

- **Vektorizace redukční fáze na CPU:** Současná implementace využívá SIMD pouze pro výpočet mediánu napříč dny. Vektorizace dalších částí výpočtu (medián napříč pacienty, medián při downsamplingu) by mohla dále zlepšit výkon paralelní CPU verze. Efektivní implementace řazení dynamicky dlouhého pole pomocí vektorových instrukcí je však náročná. Já jsem se o to pokoušel (viz. kapitola 5.4), ale bez úspěchu.
- **Dynamická volba strategie:** Vývoj adaptivního systému, který by na základě velikosti dat a dostupného hardwaru automaticky zvolil nejefektivnější zpracovatelský přístup (CPU vs. GPU).

Vzhledem k velikosti některých datových sad (BVP) trvá načítání dat z disku značný čas (více než 26 minut). Pro reálné aplikace by bylo vhodné zvážit optimalizace v této oblasti, například paralelním načítáním dat.

# Kapitola 10

## Závěr

Projekt úspěšně implementuje zpracování velkých medicínských datových sad 3 způsoby - sekvencí CPU, paralelní CPU+SIMD a GPU. Během vývoje byly vyřešeny klíčové technické výzvy:

- **Filtrace dat:** Specifika BVP signálu (záporné hodnoty) si vynutila změnu logiky pro detekci prázdných míst. Zavedení konstanty `EMPTY_VALUE = -1e30` a sjednocení podmínek ve všech metodách zajistilo konzistenci výsledků.
- **Optimalizace:**
  - **GPU:** Využití Sorting Networks v prvním kernelu a in-place Shell Sortu ve třetím kernelu minimalizovalo nároky na paměť.

Výsledná aplikace prokazuje, že pro masivní data (BVP) je GPU akcelerace vysoce efektivní, zatímco pro menší datové sady (Dexcom) je vhodnější optimalizovaná paralelní verze na CPU.

# Příloha A: Uživatelský manuál

Aplikace je navržena jako nástroj příkazové řádky (CLI), který umožňuje flexibilní konfiguraci benchmarků. Uživatel může volit, které datové sady se mají zpracovat, jaké výpočetní metody (sekvenční, paralelní, GPU) se mají použít a kam se mají uložit výsledky. **V aplikaci jsou použity vektorové instrukce NEON SIMD specifické pro ARM architekturu.**

## Překlad a Sestavení

Projekt využívá sestavovací systém CMake. Pro sestavení aplikace je nutné mít nainstalovaný kompilátor podporující C++23 a knihovnu OpenCL.

```
mkdir build
cd build
cmake ..
make
```

Po úspěšném sestavení vznikne v adresáři `build` spustitelný soubor `ppr_semestralka`. Příkazem `./ppr_semestralka -h` lze zobrazit nápovědu k ovládání.

## Přidání složky s daty

Vstupní data by měla být umístěna ve složce `data` ve kořenovém adresáři projektu (ne ve složce `build`). Struktura složky `data` by měla být následující:

```
data/
  001/
    DEXCOM_001.csv
    HR_001.csv
    BVP_001.csv
  002/
    DEXCOM_002.csv
    HR_002.csv
    BVP_002.csv
  003/
    DEXCOM_003.csv
    HR_003.csv
    BVP_003.csv
  ...
```

## Parametry příkazové řádky

Program přijímá následující přepínače pro konfiguraci běhu:



**-d, -datasets <list>** Specifikuje čárkou oddělený seznam datových sad ke zpracování.

- Možnosti: `dexcom`, `hr`, `bvp`
- Výchozí: Všechny (`dexcom,hr,bvp`)

**-m, -mode <list>** Specifikuje čárkou oddělený seznam výpočetních metod.

- Možnosti: `seq` (Sekvenční), `par` (Paralelní CPU), `gpu` (OpenCL), `all` (Všechny)
- Výchozí: `all`

**-k, -kernels <num>** Nastavuje počet aktivních fází (kernelů) v GPU pipeline. Slouží primárně pro analýzu úzkých hrdel.

- 1: Pouze výpočet mediánů pro pacienta (masivní přenos dat).
- 2: Přidána redukce pacientů na GPU (minimalizace přenosu).
- 3: Kompletní pipeline včetně downsamplingu (vysoká výpočetní zátěž).
- Výchozí: 3

**-o, -output <file>** Cesta k souboru pro logování výsledků benchmarku.

- Výchozí: `benchmark_logs.txt`

**-h, -help** Vypíše nápovědu k ovládání.

## Příklady použití

**1. Základní spuštění (všechny testy):** Spustí všechny algoritmy nad všemi datovými sadami s výchozím nastavením.

```
./ppr_semestralka
```

**2. Benchmark pouze GPU na velkých datech:** Spustí pouze GPU implementaci nad BVP daty a výsledky uloží do specifického souboru.

```
./ppr_semestralka -d bvp -m gpu -o gpu_results.txt
```

**3. Analýza GPU pipeline (Memory Bound test):** Spustí GPU výpočet nad BVP daty pouze s prvním kernelem pro otestování propustnosti paměti.

```
./ppr_semestralka -d bvp -m gpu -k 1
```

**4. Porovnání CPU metod:** Spustí sekvenční a paralelní verzi nad HR a Dexcom daty pro ověření zrychlení na procesoru.

```
./ppr_semestralka -d hr,dexcom -m seq,par
```