

UNIWERSYTET GDAŃSKI
WYDZIAŁ MATEMATYKI, FIZYKI I INFORMATYKI

Krzysztof Wiśniewski
numer albumu: 274276

Kierunek studiów: Bioinformatyka
Specjalność: Ogólna

Optymalizacja oprogramowania do detekcji splątania kwantowego

Praca licencjacka
wykonana
pod kierunkiem
dr hab. Marcin Wieśniak, prof. UG

Gdańsk 2023

Spis treści

| | | |
|----------|--|-----------|
| 1 | Wstęp | 2 |
| 1.1 | Cele pracy | 2 |
| 1.2 | Działanie programu | 2 |
| 1.3 | Modularyzacja | 3 |
| 2 | Metody | 4 |
| 2.1 | Kompilacja AOT | 4 |
| 2.2 | Kompilacja JIT | 5 |
| 2.3 | Dane testowe | 7 |
| 2.4 | Środowisko testowe | 9 |
| 2.5 | Profilowanie | 9 |
| 2.6 | Narzędzia pomocnicze | 10 |
| 3 | Wyniki | 10 |
| 3.1 | Wstępne profilowanie | 10 |
| 3.2 | Wstępne pomiary wydajności | 12 |
| 3.3 | Pomiary czasu pracy re-implementacji | 13 |
| 3.3.1 | Python i NumPy | 13 |
| 3.3.2 | Python i NumPy z AOT | 14 |
| 3.3.3 | Python i NumPy z JIT | 15 |
| 3.3.4 | Rust i Ndaray | 16 |
| 3.3.5 | Rust i Ndaray z OpenBLAS | 18 |
| 3.4 | Precyzja obliczeń | 19 |
| 3.5 | Pomiary z pojedynczą precyzją | 20 |
| 3.5.1 | Python i NumPy | 20 |
| 3.5.2 | Python i NumPy z AOT | 21 |
| 3.5.3 | Python i NumPy z JIT | 22 |
| 3.5.4 | Rust i Ndaray | 22 |
| 3.5.5 | Rust i Ndaray z OpenBLAS | 23 |
| 4 | Dyskusja | 24 |
| | Odwołania | 25 |

1 Wstęp

Closest Separable State Finder (CSSFinder) jest programem pozwalającym na detekcję splątania kwantowego układu oraz określenie jak silnie owe splątanie jest. Bazuje on na dostosowanym algorytmie Elmera G. Gilberta, pozwalający na wyliczenie przybliżonej wartości odległości Hilberta-Schmidta (ang. Hilberta-Schmidta distance, HSD) pomiędzy stanem a zbiorem stanów separowanych. W literaturze pojawia się on pod nazwą ‘kwantowy algorytm Gilberta’(ang. quantum Gilbert algorithm, QGA)[9]. Wykorzystanie tego algorytmu zostało opisane w pracy ‘Hilbert-Schmidt distance and entanglement witnessing’ której autorami byli Palash Pandya, Omer Sakarya i Marcin Wieśniak[8]. Profesor Marcin Wieśniak utworzył implementację algorytmu QGA w języku Python, wykorzystując bibliotekę NumPy do przeprowadzania koniecznych obliczeń macierzowych. Wybór ten był podyktowany możliwościami oferowanymi przez taki zestaw narzędzi. Pozwalały one w szybki sposób stworzyć prosty kod, zdolny by relatywnie wydajnie przeprowadzać obliczenia na wszystkich najpopularniejszych systemach dla komputerów stacjonarnych. Alternatywy w postaci języków C, C++ czy Fortran wymagałyby większej ilości bardziej skomplikowanego kodu, jednocześnie zmuszając do ręcznego skompletowania systemu budowania, bibliotek oraz zastosowania dedykowanych rozwiązań dla każdego systemu operacyjnego.

Obecnie, posiadając sprawną implementację naturalnym jest rozpoczęcie eksploracji możliwości poprawy wydajności programu, jego dostępności i komfortu użytkowania.

1.1 Cele pracy

Celami tej pracy są:

- eksploracja dostępnych metod maksymalizacji wydajności algorytmu QGA,
- implementacja algorytmu QGA wybranymi z metod,
- weryfikacja poszczególnych z tych rozwiązań pod kątem zmian w czasie pracy programu,
- udoskonalenie interfejsu użytkownika programu CSSFinder,
- dostosowane postaci i sposób dystrybucji programu do standardów ekosystemu języka Python.

1.2 Działanie programu

Oryginalny program i jego re-implementacja posiadają praktycznie identyczną zasadę działania i tylko szczegóły dotyczące sposobu interakcji z nim zmieniły się. Z tego względu, w dalszej części tekstu, sposób działania programu będzie opisywany bez rozgraniczenia na wersję oryginalną i re-implementację.

Program jako dane wejściowe przyjmuje macierz gęstości opisującą pewien stan ρ_0 układu kwantowego. Następnie program w określonych wypadkach jest wydedukować wymiary

podukładów i ich ilość, lub można je podać jawnie. Następnie dobierany jest stan separowalny ρ_1 . Następnie program postępuje zgodnie z następującymi krokami:

1. Zwiększ licznik prób c_t o 1. Wylosuj czysty stan produktowy ρ_2 , zwany dalej stanem próbnym.
2. Uruchom preselekcję dla stanu próbnego poprzez sprawdzenie funkcji liniowej. Jeśli się nie powiedzie, wróć do punktu 1.
3. W przypadku udanej preselekcji symetryzujemy ρ_1 względem wszystkich symetrii przez ρ_0 , które respektują separowalność.
4. Znaleźć minimum $Tr(\rho_0 - p\rho_1 - (1-p)\rho_2)^2$ względem p .
5. Jeśli minimum występuje dla $0 \leq p \leq 1$, zaktualizuj $\rho_1 \leftarrow p\rho_1 + (1-p)\rho_2$, dodać nową wartość $D^2(\rho_0, \rho_1)$ do listy listy i zwiększyć wartość licznika sukcesu c_s o 1.
6. Przejdź do kroku 1, aż spełnione zostanie wybrane kryterium zatrzymania.

Jako dane wyjściowe program zapisuje do plików historię poprawek i stan ρ_1 . Dostępnymi kryteriami zatrzymania jest maksymalna ilość korelacji do uzyskania oraz maksymalna ilość iteracji do wykonania, ta z tych wartości która zostanie osiągnięta jako pierwsza decyduje o zatrzymaniu programu. Jeśli wyznaczona przez program odległość HSD jest odpowiednio niewielka, tj. mniejsza niż $1 \cdot 10^{-4}$, stan jest praktycznie separowalny, w przeciwnym wypadku może być uważany za splątany.

1.3 Modularyzacja

Re-implementując program CSSFinder planowałem wypróbować liczne rozwiązania, które wymagały zasadniczych zmian w kodzie algorytmu, w tym przepisania go w innym języku programowania. Jednocześnie część programu odpowiadająca za interakcję z użytkownikiem i ładowanie zasobów miała pozostawać taka sama. Zdecydowałem więc że tworzony przeze mnie kod musi być modularny, aby uniknąć duplikacji wspólnych elementów. Tak też program został podzielony na dwie części: główną ('core'), z interfejsem użytkowników i narzędziami pomocniczymi oraz część implementującą algorytm ('backend'). Korpus jest w całości napisany w języku Python i wykorzystuje wbudowany w ten język mechanizm importowania bibliotek w celu wykrywania i ładowania dostępnych implementacji algorytmu ('backendów'). Dane macierzowe w obrębie korpusu przechowywane są jako obiekty ndarray z biblioteki NumPy, ze względu na uniwersalność w świecie bibliotek do obliczeń tensorowych (wiele bibliotek w innych językach programowania oferuje gotowe narzędzia do transformacji obiektów ndarray na reprezentacje charakterystyczne dla tych bibliotek).

Pozwala to na proste podmiany implementacji o dowolnie różnym pochodzeniu, w tym implementacje w językach kompilowanych. Uprościło to znacznie proces weryfikacji zmian w zachowaniu programu i przyspieszyło proces tworzenia kolejnych implementacji, jako że

kod interfejsu programistycznego jest mniej pracochłonny niż kod pozwalający na interakcję z użytkownikiem. W przyszłości może to również pozwolić na łatwiejszy rozwój nowych implementacji oraz dodawanie nowych funkcjonalności do programu, ponieważ będą one mogły być implementowane stopniowo w różnych implementacjach.

2 Metody

2.1 Kompilacja AOT

Kompilacja AOT (Ahead Of Time) to proces tłumaczenia jednej reprezentacji programu (na przykład w języku programowania wysokiego poziomu) na inną (na przykład kod maszynowy) przed rozpoczęciem pracy kompilowanego programu.

Obecnie najpowszechniej używana implementacja języka Python, CPython, posiada możliwość korzystania z bibliotek współdzielonych (.so - Linux, .dll/.pyd - Windows) które powstały w skutek kompilacji kodu wysokiego poziomu. Dostęp do funkcji zawartych w takich bibliotekach można uzyskać na kilka sposobów:

1. Przy pomocy API modułu ctypes[18]. Pozwala ono opisać interfejs funkcji obcej (tj. takiej która została napisana w języku niższego poziomu i skompilowana do kodu maszynowego) i wywołać tak opisaną funkcję.
2. Poprzez zawarcie w bibliotece odpowiednio nazwanych symboli, automatycznie rozpoznawanych przez interpreter języka Python. Takie biblioteki określa się mianem modułów rozszerzeń [19]. W tym przypadku warto dodać, że pomimo, że oficjalna dokumentacja wspomina tylko o językach C i C++, natomiast powstały biblioteki które pozwalają wykorzystać w łatwy sposób wiele innych języków programowania, takich jak Rust przy pomocy Py03[16] lub GO z użyciem biblioteki gopy[12].
3. Wykorzystując bibliotekę Cython[13][3]. Oferuje ona dedykowany język, o tej samej nazwie, który jest nadzbiorem języka Python, który rozszerza jego składnię o możliwość statycznego typowania. Biblioteka zawiera transpiler, zdolny przetłumaczyć dedykowany język na C/C++, a następnie, wykorzystując osobno zainstalowany kompilator, skompilować do kodu maszynowego.
4. Kompilując kod pythona z użyciem biblioteki mypyc[23]. Ta, podobnie do biblioteki Cython, również zawiera transpiler, natomiast zamiast korzystać z dedykowanego języka, opiera się on na dodanych w Pythonie 3.5[4] (PEP 484[21] i PEP 483[20]), adnotacjach typów. Jest on rozwijany obok projektu mypy - pakietu do statycznej analizy typów dla języka Python, również opartej na adnotacjach typów[22].

Ponieważ w każdym z wymienionych przypadków, kod niższego poziomu jest kompilowany przed dostarczeniem do użytkownika, pozwala to na wykorzystanie zaawansowanych możliwości automatycznej optymalizacji dostarczanych przez współczesne kompilatory, na przykład LLVM,

które jest sercem implementacji clang (język C++) oraz rustc (język Rust). Wiele bibliotek korzysta z mieszanek wymienionych powyżej metod, w tym cieszące się dużą popularnością NumPy, CuPy, Tensorflow, czy PyTorch. Dwie ostatnie biblioteki koncentrują się w głównej mierze na uczeniu maszynowym i głównie pod tym kontem są optymalizowane. Ich interfejsy są bardzo zbliżone do NumPy i CuPy, ale brakuje w nich niektórych narzędzi, które nie znajdują zastosowania w dziedzinie sztucznej inteligencji. W dalszej części pracy intensywnie wykorzystywana będzie biblioteka NumPy. Niestety, ze względu na ograniczenia czasowe oraz wstępne przewidywania dotyczące wydajności biblioteka CuPy nie wzięta pod uwagę¹.

Ponadto w zestawieniu pojawi się implementacja napisana w języku Rust, wykonująca operacje macierzowe w oparciu o bibliotekę Numpy[17]. Komunikacja pomiędzy interpreterem Pythona, a biblioteką oparta została o rozwiązanie opisane w punkcie drugim, dzięki wspomnianej tam bibliotece PyO3[16]. Jest to wymiennie reprezentatywny przykład tego jaką wydajność można uzyskać tworząc kod w typowym niskopoziomowym języku programowania, posiadającym relatywnie niskopoziomą kontrolę nad pamięcią. Dodatkowo w zestawieniach pojawi się wariant tej implementacji który dodatkowo będzie wykorzystywał bibliotekę OpenBLAS[25][17] do operacji mnożenia macierzowego.

2.2 Kompilacja JIT

Kompilacja JIT to proces tłumaczenia jednej reprezentacji programu (na przykład w języku programowania wysokiego poziomu) na inną (na przykład kod maszynowy) po rozpoczęciu pracy programu. Zazwyczaj wymaga to aby program rozpoczynał pracę w trybie interpretowanym, a następnie kompilował sam siebie i przechodził w tryb wykonywania skompilowanego kodu.

W momencie pisania tej pracy istnieją dwa szeroko dostępne i aktywnie utrzymywane narzędzia oferujące kompilację JIT dla języka Python.

Pierwszym z nich jest pełna alternatywna implementacja języka Python - PyPy[24]. Wykonywana przez nią kompilacja JIT działa on na podobnej zasadzie do uprzednio wymienionych - śledzi cały kod który wykonuje i automatycznie decyduje które fragmenty skompilować do kodu maszynowego[2]. Niestety, posiada ona zasadniczą wadę - jej interfejs binarny² oraz programistyczny³ różni się od CPythona, a większość pakietów które normalnie wykorzystują moduły rozszerzeń nie oferuje pre-kompilowanych pakietów dla PyPy. Powoduje to że instalacje takich pakietów są bardzo czasochłonne i obecności kompilatora na urządzeniu docelowym. Dodatkowo, pre-kompilowany kod nie czerpie żadnych korzyści z kompilatora JIT zawartego w PyPy. Problemy te powodują, że PyPy nadaje się głównie do wykonywania aplikacji napisanych w czystym języku Python.

¹CuPy jest odpowiednikiem NumPy który wykorzystuje do obliczeń GPU. Z tego względu radzi sobie wyśmienicie z operacjami na dużych macierzach, natomiast najprawdopodobniej macierze tutaj rozważane są zbyt małe aby uzyskać wzrost wydajności[6]. Jednocześnie pomimo podobieństwa do NumPy, biblioteka ta różni się i posiada problematyczne zależności (CUDA) co czyni adaptację kodu czasochłonną.

²ang. ABI - Application Binary Interface

³ang. API - Application Programming Interface.

Drugim narzędziem jest biblioteka Numba[5][7]. Ona, w przeciwieństwie do PyPy, wymaga aby fragmenty kodu, które mają być skompilowane, miały postać funkcji oznaczonych dedykowanymi dekoratorami⁴. Została ona również zaprojektowana aby dobrze współpracować z biblioteką NumPy. Jej zastosowanie z założenia ma generować wzrost wydajności nawet w sytuacjach gdy kod programu bardzo mocno eksploatuje możliwości biblioteki NumPy.

Z uprzednio wymienionych względów dotyczących preferowanych zastosowań powyższych rozwiązań w dalszej części będę próbował wykorzystać bibliotekę Numba, natomiast pominię możliwość skorzystania z PyPy.

⁴Obecnie dostępny jest też dekorator pozwalający na kompilację klas, niestety jest on niestabilny i nie radzi sobie w wielu sytuacjach.

Podczas pomiarów konsekwentnie wykorzystywałem ten sam zestaw macierzy gęstości, aby móc wygodnie porównywać wyniki wydajności poszczególnych implementacji. W dalszej części pracy będę wielokrotnie odnosił się do tych macierzy posługując się symbolem ρ z liczbą w indeksie dolnym. Liczba ta będzie wskazywać na konkretną z wymienionych poniżej macierzy.

[illegible]

Pierwsza wymieniana macierz opisuje układ 5 kubitów i posiada wymiary 32×32 . Pomimo że nie zawiera ona wartości, podczas analizy zawsze będzie reprezentowana przez macierze zawierające liczby zespolone, ponieważ szczególnie kosztowne obliczeniowo części algorytmu wymagają aby części urojone były obecne, co znaczy że usuwanie ich w wybranych miejscach nie niesie wymiernych zysków wydajnościowych.

$$\rho_n = \begin{bmatrix} 0.5 & 0 & \dots & 0 & 0.5 \\ 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \\ 0.5 & 0 & \dots & 0 & 0.5 \end{bmatrix}_{(2^n \times 2^n)}$$

7

W tekście macierze te będą oznaczane jako ρ_2 do ρ_6 , w zależności od reprezentowanej liczby kubitów⁵. Macierze te stanowią wygodny zestaw danych do weryfikacji ogólnej charakterystyki zachowania alternatywnych implementacji algorytmu, pomimo, że wyniki przy ich pomocy uzyskiwane tak bardzo odbiegają od tych uzyskiwanych przy pomocy ρ_0 .

⁵Tak więc macierz ρ_2 ma wymiary 4×4 i reprezentuje 2 kubity, macierz ρ_3 ma wymiary 8×8 i reprezentuje 3 kubity, macierz ρ_4 ma wymiary 16×16 i reprezentuje 4 kubity, itd. aż do ρ_6 , 64×64 .

2.4 Środowisko testowe

Podczas pomiarów wydajności wykorzystywałem każdorazowo to samo środowisko testowe. Do chłodzenia CPU wykorzystywane było chłodzenie wodne typu AIO, temperatura w pokoju oscylowała w okolicy 25°C, procesor podczas testów wydajności nie doświadczał temperatur powyżej 80°C.

| | |
|--------|--------------------------------|
| OS | Ubuntu 22.04.2 LTS 64-bit |
| Kernel | 5.19.0-42-generic |
| Python | 3.10.6 64-bit |
| NumPy | 1.23.5 |
| Numba | 0.56.4 |
| Cython | 3.0.0b1 |
| GCC | 11.3.0 64-bit |
| Rust | 1.68.2 64-bit |
| CPU | AMD Ryzen 9 7950X |
| RAM | 64GB DDR5 5600MHz CL40 |
| DRIVE | 512GB SSD GOODRAM CX400 (SATA) |

Tablica 1: Konfiguracja środowiska testowego.

2.5 Profilowanie

Podczas prac nad optymalizacją czasu pracy programu kluczowym było stałe zbieranie informacji na temat tego które fragmenty kodu pochłaniają najwięcej czasu. Standardowo proces zbierania takich danych określa się mianem profilowania i technologie po które sięgałem podczas re-implementacji algorytmu posiadają gotowe narzędzia pozwalające na skuteczne pozyskiwanie takich danych oraz ich wizualizację.

Dla kodu w języku Python, implementacja CPython tego języka posiada w bibliotece standardowej dwa dedykowane moduły oferujące funkcjonalność profilowania: ‘profile’ i ‘cProfile’. Pierwszy jest zaimplementowany w języku Python, drugi w C. Ponieważ drugi z nich posiada mniejszy dodatkowy narzut na procesor, zdecydowałem żeby to na nim oprzeć moje analizy. W celu wizualizacji uzyskanych wyników posłużyłem się otwartoźródłowym programem Snakeviz[14].

Do zbierania informacji na temat charakterystyki pracy kodu napisanego w języku Rust wykorzystałem narzędzie perf pochodzące z pakiety linux-tools-5.19.0-42-generic pobranego przy pomocy menadżera pakietów apt-get. Do wizualizacji uzyskanych wyników wykorzystałem jedno z otwartoźródłowych narzędzi funkcjonujące pod nazwą hotspot[11].

2.6 Narzędzia pomocnicze

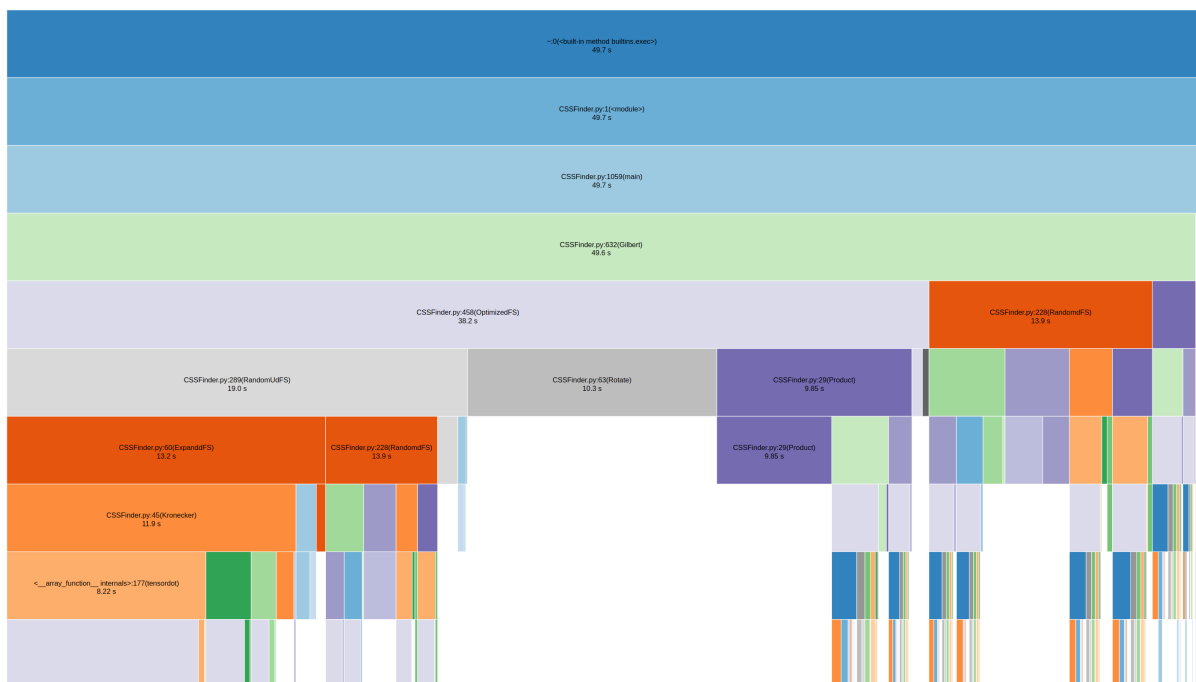
Wszystkie wykresy zamieszczone w tej pracy zostały utworzone przy pomocy skryptów w języku Python z wykorzystaniem biblioteki matplotlib[1].

3 Wyniki

3.1 Wstępne profilowanie

Prace nad optymalizacją kodu rozpocząłem od wstępnego profilowania pracy programu w trybie 1 (ang. full separability of an n-quDit state) przekazując do obliczeń układ 5 kubitów opisany macierzą ρ_0 (Rysunek 1).

Program wykonywał proces analizy stanu aż do uzyskania 1000 korekcji. Przekazany limit liczby iteracji wynosił 2.000.000 i nie został osiągnięty. Podczas pomiarów, program wykorzystywał domyślny globalny generator liczb losowych biblioteki NumPy (PCG64[15]) z ziarnem ustawionym na wartość 0.



Rysunek 3: Diagram podsumowujący pracę programu wygenerowany przez program Snakeviz.

Pozwoliło mi to wstępnie przyjrzeć się charakterystyce pracy programu i ocenić czy powszechnie dostępne narzędzia mogą zostać wykorzystane w tym wypadku. Rysunek 3 przedstawia diagram, typu Icicle, obrazujący udział czasu, pochłoniętego przez wykonywanie poszczególnych funkcji, w całkowitym czasie pracy programu. Pierwszy blok od góry (`:0(<built-in method builtins.exec>)`) to wywołanie funkcji wykonującej kod programu. Następne bloki, których opisy zaczynają się od `'CSSFinder.py'` to wywołania w kodzie programu. Bloki umieszczone najniżej, w większości pozbawione opisów, to wywołania do funkcji bibliotek, głównie NumPy, ale również modułów wbudowanych Pythona. Snakeviz automatycznie

podejmuje decyzję o nie adnotowaniu bloku gdy opis nie ma szansy zmieścić się w obrębie bloku. Aby usunąć z diagramu zbędny szum informacyjny, funkcje których wykonywanie zajęło mniej niż 1% czasu programu były pomijane.

| ncalls | totttime | percall | cumtime | percall | filename:lineno(function) |
|---------|-----------|-----------|-----------|-----------|-----------------------------------|
| 1 | 1.431e-05 | 1.431e-05 | 49.72 | 49.72 | CSSFinder.py:1(<module>) |
| 1 | 7.526e-05 | 7.526e-05 | 49.68 | 49.68 | CSSFinder.py:1059(main) |
| 1 | 0.3098 | 0.3098 | 49.63 | 49.63 | CSSFinder.py:632(Gilbert) |
| 1028 | 0.5381 | 0.0005234 | 38.2 | 0.03716 | CSSFinder.py:458(OptimizedFS) |
| 411200 | 0.8332 | 2.026e-06 | 19.03 | 4.627e-05 | CSSFinder.py:289(RandomUdFS) |
| 595516 | 0.67 | 1.125e-06 | 13.88 | 2.331e-05 | CSSFinder.py:228(RandomdFS) |
| 411200 | 0.384 | 9.338e-07 | 13.17 | 3.203e-05 | CSSFinder.py:60(ExpanddFS) |
| 822400 | 0.7256 | 8.823e-07 | 11.94 | 1.452e-05 | CSSFinder.py:45(Kronecker) |
| 849257 | 10.3 | 1.213e-05 | 10.3 | 1.213e-05 | CSSFinder.py:63(Rotate) |
| 1068026 | 6.535 | 6.118e-06 | 9.85 | 9.223e-06 | CSSFinder.py:29(Product) |
| 1332780 | 2.17 | 1.628e-06 | 4.502 | 3.378e-06 | CSSFinder.py:21(Normalize) |
| 1332780 | 2.247 | 1.686e-06 | 3.802 | 2.853e-06 | CSSFinder.py:33(Generate) |
| 737264 | 0.4225 | 5.73e-07 | 2.548 | 3.456e-06 | CSSFinder.py:18(Outer) |
| 595516 | 0.4642 | 7.794e-07 | 2.361 | 3.964e-06 | CSSFinder.py:26(Project) |
| 1233601 | 0.8998 | 7.294e-07 | 1.165 | 9.447e-07 | CSSFinder.py:39(IdMatrix) |
| 1 | 3.046e-06 | 3.046e-06 | 0.05277 | 0.05277 | CSSFinder.py:96(readmtx) |
| 1 | 1.752e-06 | 1.752e-06 | 0.05277 | 0.05277 | CSSFinder.py:552(Initrho0) |
| 1 | 4.597e-06 | 4.597e-06 | 0.002477 | 0.002477 | CSSFinder.py:1049(DisplayLogo) |
| 1 | 5.189e-06 | 5.189e-06 | 0.0004394 | 0.0004394 | CSSFinder.py:954(DetectDim0) |
| 1 | 1.628e-05 | 1.628e-05 | 2.526e-05 | 2.526e-05 | CSSFinder.py:556(Initrho1) |
| 1 | 1.903e-06 | 1.903e-06 | 5.671e-06 | 5.671e-06 | CSSFinder.py:599(DefineSym) |
| 40 | 3.038e-06 | 7.595e-08 | 3.038e-06 | 7.595e-08 | CSSFinder.py:192(writemtx) |
| 1 | 1.102e-06 | 1.102e-06 | 2.846e-06 | 2.846e-06 | CSSFinder.py:624(DefineProj) |
| 2 | 2.3e-07 | 1.15e-07 | 2.3e-07 | 1.15e-07 | CSSFinder.py:845(makeshortreport) |

Tablica 2: Dane dotyczące pracy oryginalnej implementacji programu CSSFinder uzyskane przy pomocy programu cProfile. Tabela posiada oryginalne nazwy kolumn, nadane przez program Snakeviz. Znaczenia kolumn, kolejno od lewej: **ncalls** - ilość wywołań funkcji. **totttime** - całkowity czas spędzony w ciele funkcji bez czasu spędzonego w wywołaniach do podfunkcji. **percall** - **totttime** dzielone przez **ncalls**. **cumtime** - całkowity czas spędzony wewnątrz funkcji i w wywołaniach podfunkcji. **percall** - **cumtime** dzielone przez **ncalls**. **filename:lineno(function)** - Plik, linia i nazwa funkcji.

Z uzyskanych danych wynika że znakomitą większość (77%⁶) czasu pracy programu zajmuje funkcja `OptimizedFS()`. W jej wnętrzu 38% czasu pochłania proces generowania losowych macierzy unitarnych, który w dużej mierze wykorzystuje mnożenia tensorowe (26%). Poza funkcją `OptimizedFS()`, znaczący wpływ na czas wykonywania ma też funkcja `rotate()`, która pochłania około 21% czasu działania programu. Kolejne 20% czasu zajmuje funkcja `product()`, obliczająca odległość Hilberta-Schmidta pomiędzy dwoma stanami. Pozostałe wywołania mają stosunkowo marginalny wpływ na czas pracy i ich analiza na tym etapie nie niesie za sobą znaczących korzyści.

Takie wyniki wskazują jednoznacznie że kluczowa dla czasu pracy programu jest tu maksymalizacja wydajności pętli optymalizacyjnej, w tym zawartych w niej operacji macierzowych. Najprostszym sposobem na uzyskanie takich efektów jest zastąpienie dynamicznego systemu typów i kodu bajtowego algorytmu wykonywanego przez interpretera pythona na statyczny system typów i kod maszynowy. Dodatkowo, niezastąpione są biblioteki

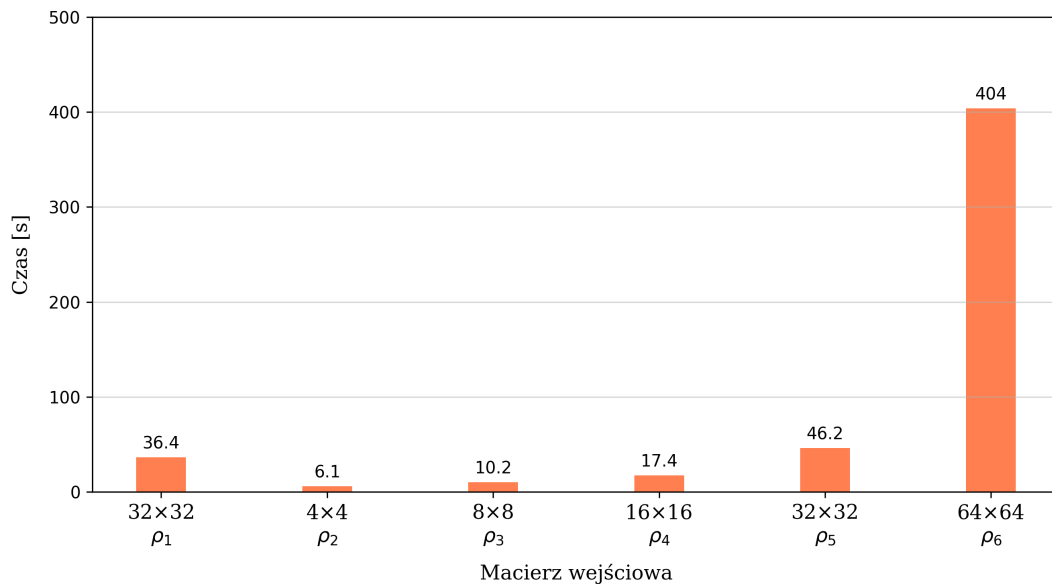
⁶Wartość 77% jak i wartości procentowe dalszej części tego akapitu zostały zaokrąglone do jedności, ze względu na małe znaczenie rzeczowe części ułamkowych.

zawierające wyspecjalizowane implementacje operacji macierzowych, takie jak OpenBLAS. Profilowanie pozwoliło również wykluczyć problemy z operacjami I/O⁷ oraz inne niespodziewane zjawiska.

3.2 Wstępne pomiary wydajności

Aby uzyskać dobrą bazę porównawczą, wykonałem serię pomiarów czasu pracy programu na macierzach $\rho_1, \rho_2 - \rho_6$, przedstawionych na rysunkach 1 i 2.

Dane przekazywałem kolejno do programu z poleceniem działania w trybie 1 (full separability of an n-quDit state) do osiągnięcia 1000 korekcy lub do 2.000.000 iteracji algorytmu, w zależności od tego co nastąpi szybciej. Dla wszystkich macierzy algorytm uzyskał 1000 korekcy i w żadnym przypadku nie osiągnął maksymalnej liczby iteracji. Dla każdej macierzy pomiar był powtarzany pięciokrotnie, a wyniki z pomiarów zostały uśrednione. Podczas obliczeń ziarno globalnego generatora liczb losowych biblioteki NumPy było ustawione na 0. Pomiary czasu pracy dotyczyły wyłącznie samego algorytmu⁸.

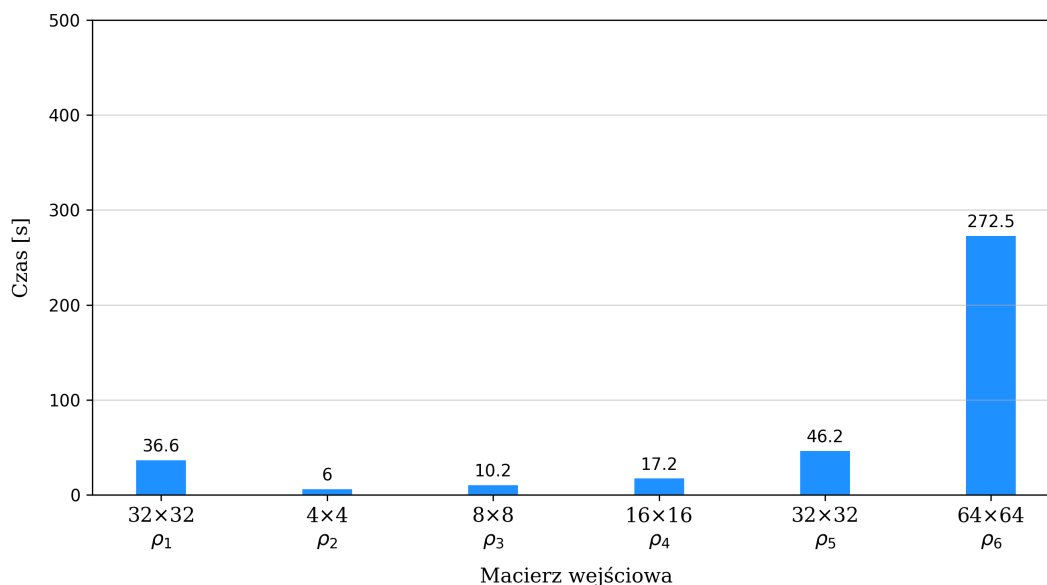


Rysunek 4: Wyniki wstępnych testów wydajności oryginalnego kodu dla macierzy $\rho_2 - \rho_6$.

Podczas testów zaobserwowałem interesujące zjawisko dotyczące wydajności dla macierzy 64×64 . W przypadku takich rozmiarów danych biblioteka NumPy automatycznie decyduje o wykorzystaniu wielowątkowej implementacji mnożenia macierzowego. Niestety, daje to efekt odwrotny do zamierzonego - obliczenia zamiast przyspieszać zwalniają. Na rysunku 4 zostały przedstawione czasy obliczeń dla macierzy $\rho_2 - \rho_6$ z domyślnym zachowaniem biblioteki.

⁷I/O - operacje wejścia wyjścia, w tym wypadku odczyt z i pisanie do plików.

⁸tj. funkcji 'Gilbert()', nie biorą więc pod uwagę czasu pochłoniętego przez importowanie modułów, ładowanie danych itp. natomiast operacje pisania do plików które były wykonywane w obrębie tej funkcji są wliczane w czas pracy.



Rysunek 5: Wyniki wstępnych testów wydajności oryginalnego kodu z zablokowaną liczbą wątków obliczeniowych dla macierzy $\rho_2 - \rho_6$.

Jeśli przy pomocy zmiennych środowiskowych ustawimy ilość wątków wykorzystywanych do obliczeń na 1 uzyskujemy znaczące skrócenie czasu obliczeń dla macierzy 64×64 . Wyniki testów w takich warunkach zostały przedstawione na rysunku 5. Dla macierzy w mniejszych rozmiarach nie odnotowałem różnicy w wydajności pomiędzy konfiguracją domyślną, a manualnie dostosowywaną. Warto dodać że ilość iteracji wykonywanych przez program nie zmienia się, różnica wynika wyłącznie z czasu trwania operacji arytmetycznych. Taki stan rzeczy najprawdopodobniej jest wynikiem dodatkowego obciążenia ze strony komunikacji i/lub synchronizacji między wątkami.

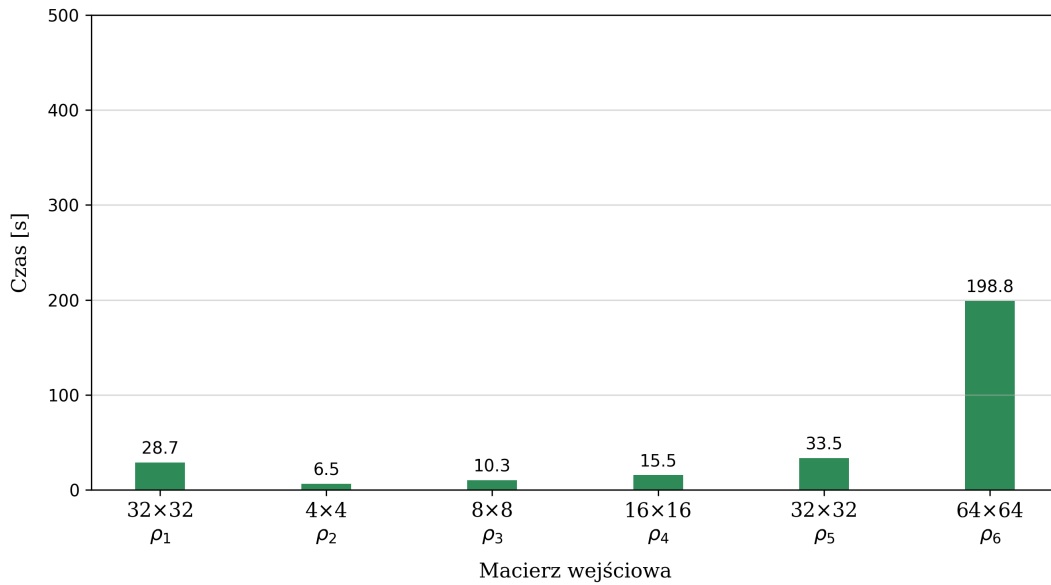
3.3 Pomiary czasu pracy re-implementacji

3.3.1 Python i NumPy

Pierwsza wykonana przeze mnie re-implementacja algorytmu, została napisana w języku Python, a do realizowania obliczeń na macierzach liczb zespolonych wykorzystywała bibliotekę NumPy. Podczas przepisywania podjąłem jednak dodatkowe wysiłki aby zastępować kod Pythona wywołaniami do funkcji zawartych w bibliotece NumPy. Ponieważ kluczowe dla wydajności fragmenty kodu tego pakietu są zaimplementowane w języku niższego poziomu, a następnie skompilowane kompilatorem optymalizującym, oferują znacznie wyższą wydajność niż analogiczny kod napisany w języku Python. Proces ten pozwolił mi również zapoznać się lepiej z charakterystyką programu i udoskonalić interfejs służący do komunikacji pomiędzy częścią główną, a samą implementacją (backend'em). Sam algorytm pozostał taki sam, natomiast konstrukcja kodu zmieniła się diametralnie, więc dogłębne analizy różnic byłyby nieczytelne, dlatego nie zostaną tutaj zawarte. W dalszej części pojawiają się wyniki pomiarów wydajności.

Pomiary czasu pracy były wykonywane przy użyciu macierzy $\rho_2 - \rho_6$. Dane przekazywałem

kolejno do programu z poleceniem działania w trybie FSnQd⁹ do osiągnięcia co najmniej 1000 korekcji lub do 2.000.000 iteracji algorytmu, w zależności od tego co nastąpi szybciej. Dla wszystkich macierzy algorytm uzyskał co najmniej 1000 korekcji i w żadnym przypadku nie osiągnął maksymalnej liczby iteracji. Dla każdej macierzy pomiar był powtarzany pięciokrotnie a wyniki zostały uśrednione. Podczas obliczeń ziarno domyślnego globalnego generatora liczb losowych biblioteki NumPy było ustawione na 0. Program działał z zablokowaną liczbą wątków obliczeniowych. Pomiary czasu pracy dotyczyły przede wszystkim samego algorytmu¹⁰.



Rysunek 6: Wyniki testów wydajności alternatywnej implementacji Python z użyciem biblioteki NumPy dla macierzy $\rho_2 - \rho_6$.

Uzyskane wyniki zostały przedstawione na rysunku 6.

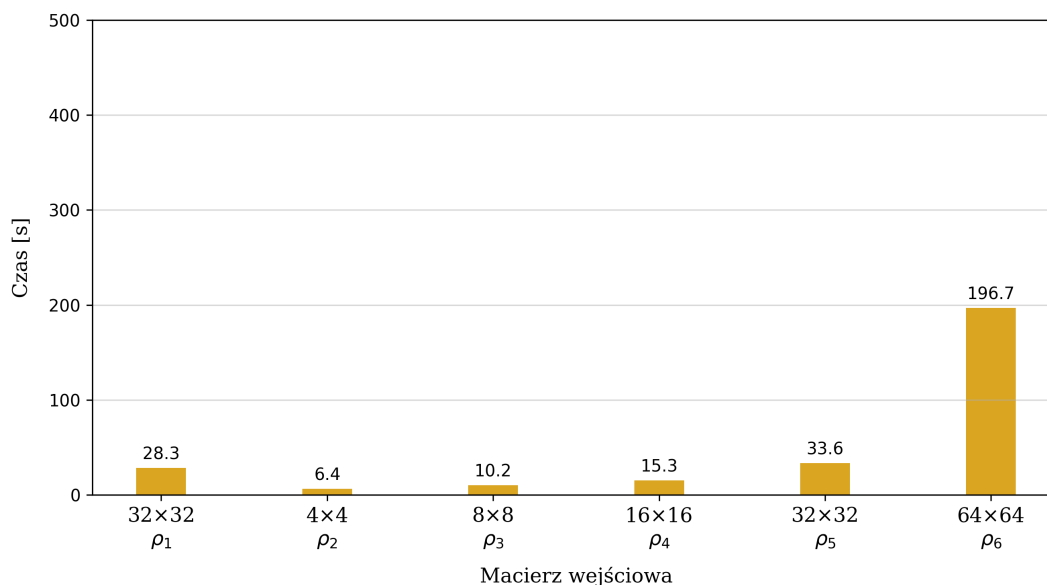
3.3.2 Python i NumPy z AOT

Następnym wykonanym przeze mnie krokiem było skompilowanie mojej implementacji korzystającej z NumPy do kodu maszynowego przy pomocy biblioteki Cython. Kod przeznaczony do takiej kompilacji nie musi być adnotowany dedykowanymi informacjami o typach. Zostanie on w tedy przetłumaczony na odpowiednie operacje w języku C/C++, a potem skompilowany do kodu maszynowego. Brak adnotacji powoduje niestety, że program zachowuje swoją dynamiczną naturę, charakterystyczną dla języka Python. Kompilacja pozwala jednak usunąć dodatkowy narzut na procesor ze strony interpretera. W takim scenariuszu spodziewać należy się, że zyski z kompilacji będą niewielkie, ale mogą wystąpić.

Pomiary czasu pracy były wykonywane w taki sam sposób jak dla implementacji bez AOT.

⁹Tryb FSnQd jest odpowiednikiem trybu 1 (full separability of an n-quDit state) z oryginalnego kodu.

¹⁰Pomiary nie biorą więc pod uwagę czasu pochłoniętego przez importowanie modułów itp., natomiast operacje wczytywania danych i pisania do plików są wliczane w czas pracy, ponieważ wbudowany w program mechanizm pomiaru czasu pracy rozpoczyna pomiar zanim dane zostaną załadowane.



Rysunek 7: Wyniki testów wydajności implementacji Python z użyciem biblioteki NumPy oraz pakietu Cython do kompilacji AOT dla macierzy $\rho_2 - \rho_6$.

Na rysunku 7 przedstawione zostały wyniki pomiarów czasu pracy skompilowanej wersji w języku Python bazującej na bibliotece NumPy wykorzystujące macierze $\rho_2 - \rho_6$. Dodatkowa kompilacja nie poskutkowała widocznym skróceniem czasu pracy programu, jedynie wynik dla macierzy 64×64 różni się nieznacznie. Może to być spowodowane usunięciem szczątkowego obciążenia ze strony interpretera, które nie jest mierzone podczas krótszych testów z mniejszymi macierzami. Natomiast możliwe jest również że ta różnica wynika z korzystniejszych warunków losowo zapewnionych przez system operacyjny.

3.3.3 Python i NumPy z JIT

Ostatnia stworzona przeze mnie re-implementacja w języku Python bazująca na bibliotece NumPy dodatkowo korzysta z kompilacji JIT. Pakiet Numba, który został wykorzystany do zrealizowania kompilacji JIT, posiada dwa tryby pracy. Pierwszy wykonuje kompilację na podstawie specjalnie dostarczonych przez programistę deklaracji typów dla funkcji podlegających kompilacji i jest wykonywany zaraz po rozpoczęciu pracy programu¹¹. Drugi polega na śledzeniu typów wejściowych i wyjściowych funkcji i automatycznie kompiluje funkcję dla tych typów danych które są odpowiednio często używane¹².

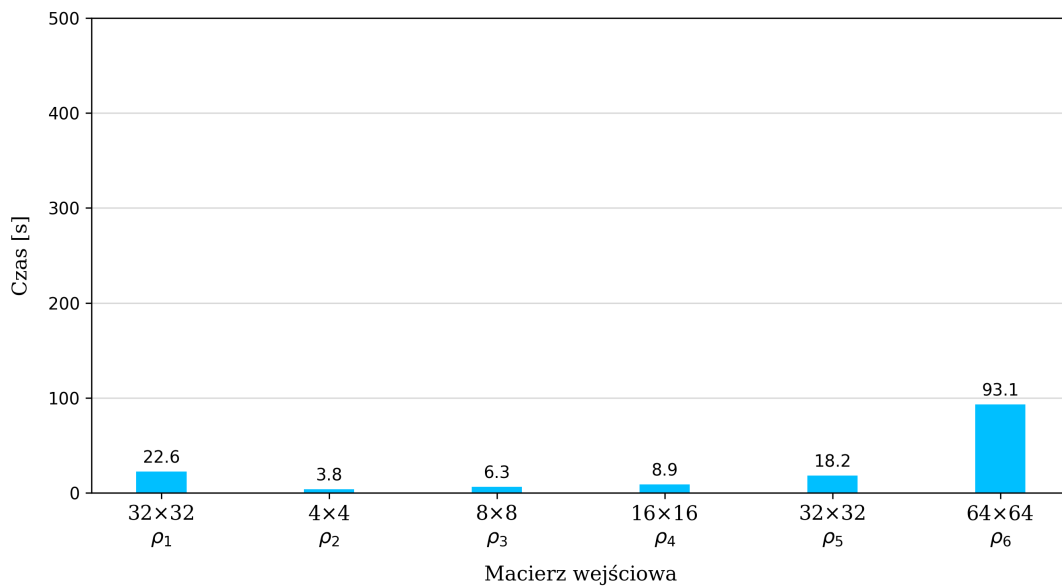
Ponadto, Numba posiada dodatkowe parametry kompilacji, które można przekazać do funkcji `numba.jit`. Jednym z nich, posiadającym szczególnie duży wpływ na wydajność, flaga `nopython`. Tryb `nopython=True` oferuje znacznie większe możliwości optymalizacji i potencjalnie lepszą wydajność. Niestety nie wszystkie funkcje dostępne w bibliotece NumPy są akceptowane przez kompilator JIT pakietu Numba w trybie `nopython=True`. Do niekompatybilnych należy

¹¹ang. eager (compilation) - niecierpliwa (kompilacja).

¹²ang. lazy (compilation) - leniwa (kompilacja).

między innymi funkcja `tensor.dot` która implementuje mnożenie tensorowe. Wspomniana funkcja może zostać skompilowana tylko w trybie obiekowym (`nopython=False`), który po kompilacji zachowuje dynamiczną naturę Pythona. Niestety, brak możliwości skompilowania funkcji używającej `tensor.dot` powoduje również brak możliwości skompilowania funkcji wyżej w drzewie wywołań. W efekcie znacząca część implementacji używającej JIT musi używać trybu obiekowego.

Pomiary czasu pracy były wykonywane w taki sam sposób jak dla implementacji bez JIT.



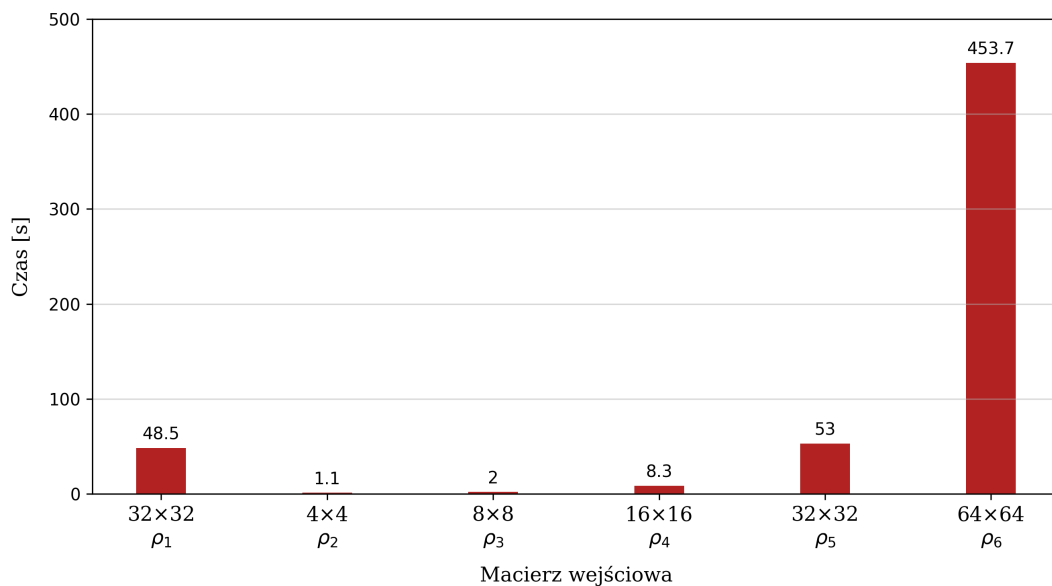
Rysunek 8: Wyniki testów wydajności implementacji w języku Python z użyciem biblioteki NumPy i pakietu Numba do kompilacji JIT dla macierzy $\rho_2 - \rho_6$.

Na rysunku 8 przedstawione zostały wyniki uzyskane podczas pomiarów czasu pracy implementacji z JIT, w zależności od rozmiaru macierzy gęstości. Kod który wykorzystywał kompilację JIT oferował podczas testów dwu-czterokrotnie (w zależności od rozmiaru macierzy testowej) większą wydajność niż kod bez niej. Tak znaczącą poprawę implementacja zawdzięcza prawdopodobnie temu, że kompilator JIT może specjalizować kod dla dokładnie jednej platformy, korzystając z całego spektrum jej możliwości. Dotyczy to na przykład instrukcji SIMD, takich jak AVX512, które są dostępne w procesorze użytym do testów, ale wiele wciąż popularnych procesorów ich nie posiada. Wymusza to, przy kompilacji AOT, zastąpienie tych instrukcji innymi szerzej dostępnymi, aby zmaksymalizować przenośność kodu. Dodatkowo kompilator może brać pod uwagę inne charakterystyczne cechy konkretnych architektur. Te dodatkowe informacje i możliwość dodatkowej specjalizacji kodu czynią kompilację JIT bardzo potężnym narzędziem

3.3.4 Rust i Ndaray

Aby uczynić to porównanie jak najpełniejszym, podjąłem również wysiłek zaimplementowania części obliczeniowej programu w języku Rust. Język ten wybrałem z kilku względów.

Posiada on pełną infrastrukturę pozwalającą w łatwy sposób kompilować programy i biblioteki wykorzystujące stworzone przez innych programistów rozwiązania. Daje mu to znaczącą przewagę nad językami takimi jak C/C++ które wymagają aby bardziej skomplikowane projekty samodzielnie skompletowały systemy budowania opierającego się na rozwiązaniach podmiotów trzecich, takich jak CMake, szczególnie jeśli chcą być dostępne na wielu platformach. Ponadto konkurencja nie posiada ujednoliconego standardu pozwalającego na łatwe uzyskanie odstępu do bibliotek otwartoźródłowych, Rust natomiast taki system posiada. W efekcie w łatwy sposób mogłem dołączyć gotowe rozwiązania pozwalające na prowadzenie obliczeń na macierzach liczb zespolonych. W efekcie cały proces wstępnej konfiguracji sprowadził się do około godziny, co stanowi wyśmienity wynik, biorąc pod uwagę, że przed podejściem do tego projektu nie miałem żadnej praktycznej styczności z tym językiem programowania. Dodatkową zaletą tego języka jest automatyczny system zarządzania pamięcią oparty na koncepcji własności (ang. *ownership*), który usuwa konieczność manualnego zarządzania pamięcią, jednocześnie bez konieczności wprowadzania mechanizmu liczenia referencji i dedykowanego automatycznego ‘odśmiecacza’ (ang. *garbage collector*) które to są częstym źródłem problemów z wydajnością i użyciem pamięci.



Rysunek 9: Wyniki testów wydajności implementacji w języku Rust dla macierzy ρ_2 - ρ_6 .

Pomiary czasu pracy implementacji w języku Rust były wykonywane przy użyciu macierzy ρ_2 - ρ_6 . Dane przekazywałem kolejno do programu z poleceniem działania w trybie FSnQd do osiągnięcia co najmniej 1000 korekcy lub do 2.000.000 iteracji algorytmu, w zależności od tego co nastąpi szybciej. Dla wszystkich macierzy algorytm uzyskał co najmniej 1000 korekcy i w żadnym przypadku nie osiągnął maksymalnej liczby iteracji. Dla każdej macierzy pomiar był powtarzany pięciokrotnie a wyniki zostały uśrednione. Pomiary czasu pracy dotyczyły przede wszystkim samego algorytmu¹³.

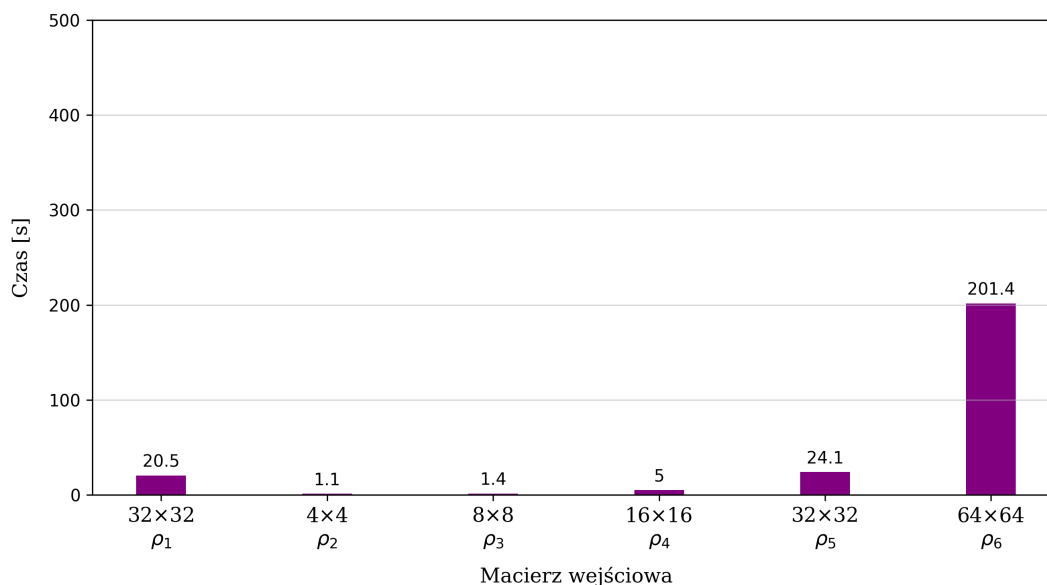
¹³Nie biorą więc pod uwagę czasu pochłoniętego przez importowanie modułów itp., natomiast operacje

Na rysunku 9 zaprezentowane zostały wyniki pomiarów czasu pracy implementacji w języku Rust. Prezentuje ona znacząco lepsze wyniki dla małych macierzy oraz znacznie gorsze wyniki dla dużych macierzy. Jest to prawdopodobnie spowodowane tym, że sama implementacja mnożenia macierzowego nie jest wyspecjalizowana, aby wykorzystywać maksimum możliwości procesora na którym jest wykonywana, w przeciwieństwie do na przykład biblioteki NumPy, które wewnętrznie wykorzystuje bibliotekę OpenBLAS[10]. W efekcie nie czerpie ona korzyści z instrukcji SIMD, takich jak AVX512.

3.3.5 Rust i Ndaray z OpenBLAS

Biblioteka Ndaray, która jest sercem implementacji w języku Rust, posiada przełącznik funkcjonalności¹⁴ który pozwala wykorzystać funkcje zawarte w bibliotece OpenBLAS jako implementację mnożenia macierzowego. O ile kompilacja dla wszystkich platform które ma wspierać CSSFinder (Windows, Linux i MacOS) jest poza moim zasięgiem, to uznałem, że warto zweryfikować jakie efekty daje wykorzystanie tej funkcjonalności w środowisku laboratoryjnym.

Pomiary czasu pracy były wykonywane w taki sam sposób jak dla implementacji która nie korzystała z OpenBLAS.



Rysunek 10: Wyniki testów wydajności implementacji w języku Rust z użyciem biblioteki OpenBLAS dla macierzy $\rho_2 - \rho_6$.

Wykorzystanie biblioteki OpenBLAS poskutkowało znaczącym wzrostem wydajności, przekraczającym możliwości oryginalnej implementacji. Wyniki te zostały przedstawione na rysunku 10. Kod tutaj omawiany ustępuje jedynie implementacji z sekcji 3.3.3, wykorzystującej JIT. Jednocześnie sprawuje się on gorzej dla małych macierzy niż wariant bez OpenBLAS.

wczytywania danych i pisanie do plików są wliczane w czas pracy.

¹⁴ang. feature switch

3.4 Precyzja obliczeń

Oryginalny program, jak i re-implementacje które pojawiły się powyżej, posługiwały się liczbami zespolonymi na bazie liczb zmiennoprzecinkowych podwójnej precyzji. Jedna taka liczba zajmuje 64 bity. Jednak w wielu przypadkach taka precyzja obliczeń nie jest konieczna do uzyskania poprawnych wyników. Podstawową zaletą wykorzystania liczb zmiennoprzecinkowych pojedynczej precyzji, czyli 32 bitowych, jest zmniejszenie rozmiaru macierzy. Pozwala na umieszczenie większej części macierzy w pamięci podręcznej procesora. Dodatkowo zwiększa to przepustowość obliczeń wykorzystujących instrukcje SIMD, ponieważ wykorzystują one rejestry o stałych rozmiarach (128, 256, 512 bitów) które mogą na ogół pomieścić dwukrotnie więcej liczb 32 bitowych niż 64 bitowych. Pozwala to oczekiwać że obliczenia wykorzystujące liczby zmiennoprzecinkowe pojedynczej precyzji będą trwały krócej.

Tworzony przeze mnie kod od początku powstawał z zamysłem umożliwienia wykorzystania liczb zmiennoprzecinkowych o różnych precyzjach, dlatego transformacja ta była dość prosta. W języku Python, wykorzystując bibliotekę NumPy przejście na liczby pojedynczej precyzji wymagało prawie każdorazowego deklarowania że wynik operacji ma posiadać typ `complex64` (cały czas mówimy o liczbach zespolonych które składają się z dwóch wartości zmiennoprzecinkowych). Nie wszystkie operacje które przyjmują parametr określający typ wejściowy są akceptowane przez kompilator JIT biblioteki Numba gdy jest on przekazywany. To ograniczenie można obejść wykonując zmianę typu jako osobną operację przy pomocy metody `astype()`.

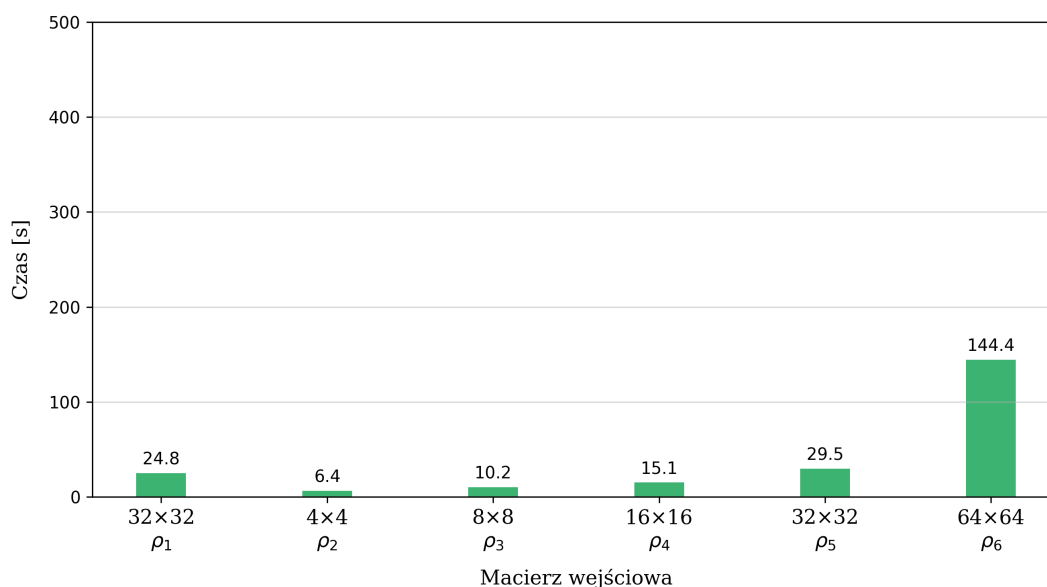
Warto tutaj zaznaczyć że wszystkie implementacje w języku Python powstają ze wspólnego szablonu który był ewaluowany przez bibliotekę Jinja2 do różnych wariantów kodu, w zależności od tego jakie parametry były do niego przekazywane. Pozwoliło to uniknąć wielokrotnego pisania wspólnych fragmentów kodu, a elementy unikalne są dodawane warunkowo. Zastosowanie introspekcji do konstruowania odpowiedniego kodu w trakcie wykonywania programu mogłoby w znaczący sposób obniżyć wydajność, dlatego zdecydowałem się sięgnąć po system bardziej statyczny, który na pewno nie wpływał na czas pracy programu.

W przypadku języka Rust, posiada on dedykowany konstrukt składniowy pozwalający na deklarowanie funkcji w oparciu o symbole zastępcze wobec których stawia się zbiór wymagań dotyczących wspieranych interfejsów. W efekcie funkcja może zostać wyspecjalizowana żeby akceptować zarówno liczby zespolone skonstruowane z liczb zmiennoprzecinkowych pojedynczej jak i podwójnej precyzji. Pozwoliło to uniknąć sięgania po zewnętrzne mechanizmy do tworzenia szablonów, tak jak było to konieczne w języku Python.

3.5 Pomiary z pojedynczą precyzją

Wszystkie testy wydajności z dla obliczeń wykorzystujących liczby zmiennoprzecinkowe pojedynczej precyzji były przeprowadzane w taki sam sposób jak odpowiadające testy z podwójną precyzją.

3.5.1 Python i NumPy

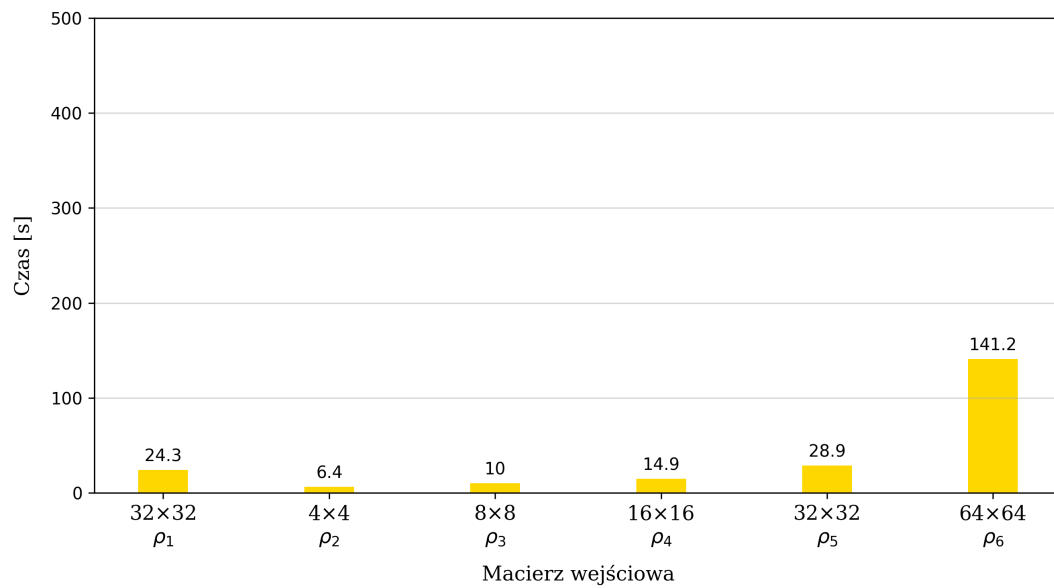


Rysunek 11: Wyniki testów wydajności implementacji w języku Python z użyciem biblioteki NumPy dla macierzy $\rho_2 - \rho_6$ i obliczeniami pojedynczej precyzji.

Na wykresie 11 przedstawiłem wyniki wydajności dla implementacji napisanej w języku Python wykorzystującej bibliotekę NumPy do przeprowadzania obliczeń na macierzach liczb zespolonych pojedynczej precyzji. W przypadku mniejszych macierzy (4×4 , 8×8 , 16×16) różnice w czasie pracy, względem wariantu opartego na liczbach podwójnej precyzji, są minimalne. Dzieje się tak prawdopodobnie dlatego, że macierze te są na tyle niewielkie (do 4KB) że mieszczą się w pamięci cache L1 procesora¹⁵, więc wyznaczanie ich jest procesem bardzo szybkim. W momencie kiedy docieramy do macierzy 32×32 wzrost wydajności staje się zauważalny, co również można wytłumaczyć odwołując się do pojemności pamięci cache procesora. Macierze podwójnej precyzji zajmują dokładnie 16KB ($32 \times 32 \times 2 \times 8$), natomiast dostęp do tej pamięci nie jest ekskluzywny dla jednego procesu, nie może on więc korzystać z całych 16KB. W efekcie część danych przebywa poza pamięcią cache. Natomiast macierze wykorzystujące liczby pojedynczej precyzji zajmują tylko 8KB. Można się więc spodziewać że większość czasu spędzają one w pamięciach L1 i L2, co pozwala przyspieszyć obliczenia. Dodatkowo mniejszy rozmiar liczb pojedynczej precyzji pozwala dwukrotnie zwiększyć przepustowość instrukcji SIMD, co prawdopodobnie odgrywa również bardzo istotną rolę, szczególnie w przypadku macierzy 64×64 , dla których obliczenia przyspieszają znacznie bardziej niż w przypadku mniejszych macierzy.

¹⁵Wykorzystywany tutaj Ryzen 9 7950X posiada 32×16 KB cache L1

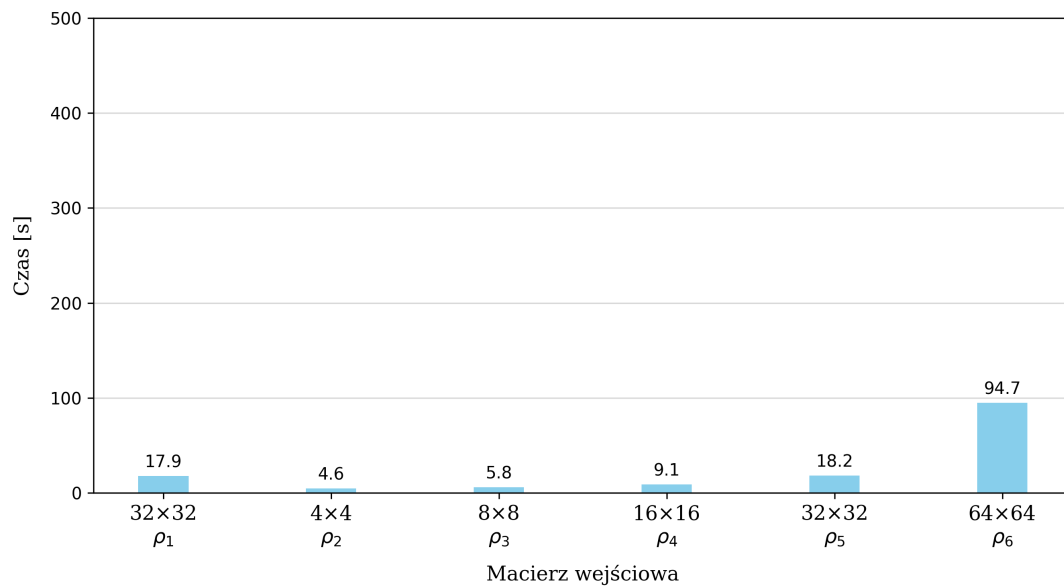
3.5.2 Python i NumPy z AOT



Rysunek 12: Wyniki testów wydajności implementacji w języku Python z użyciem biblioteki NumPy i pakietu Cython do kompilacji AOT dla macierzy ρ_2 - ρ_6 i obliczeniami pojedynczej precyzji.

Wyniki dla wariantu pre-kompilowanego przy pomocy biblioteki Cython nie różnią się znacząco od wariantu nie pre-kompilowanego, podobnie jak w przypadku obliczeń podwójnej precyzji, zostały one zaprezentowane na rysunku 12.

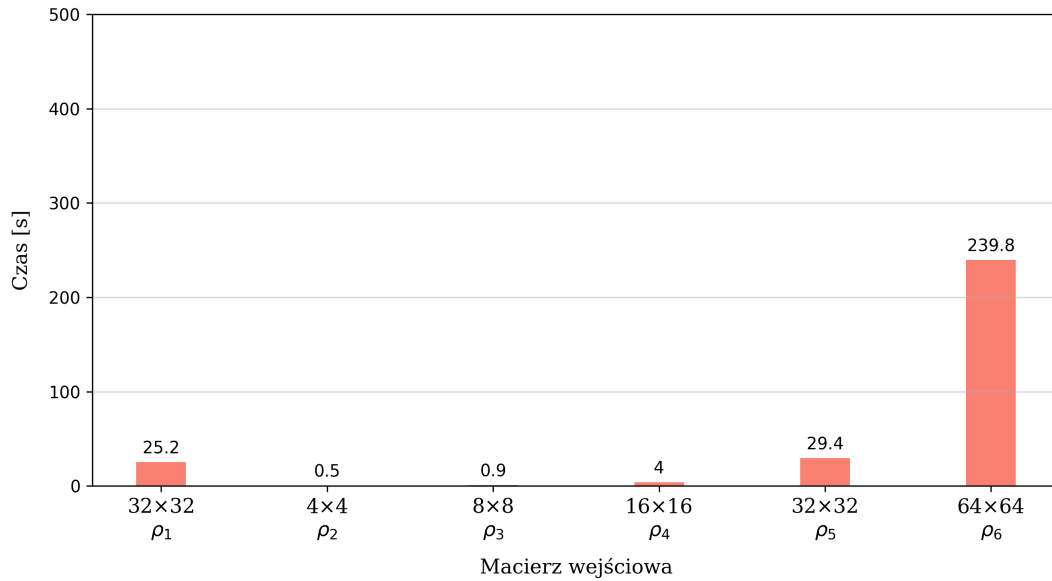
3.5.3 Python i NumPy z JIT



Rysunek 13: Wyniki testów wydajności implementacji w języku Python z użyciem biblioteki NumPy i pakietu Numba do kompilacji JIT dla macierzy ρ_2 - ρ_6 i obliczeniami pojedynczej precyzji.

W przypadku wariantu wykorzystującego kompilację JIT, zysk czasowy jest minimalny lub wręcz nie występuje. Ciężko mi wytłumaczyć dlaczego tak się dzieje, możliwe, że coś powoduje że obliczenia korzystają z tej samej implementacji.

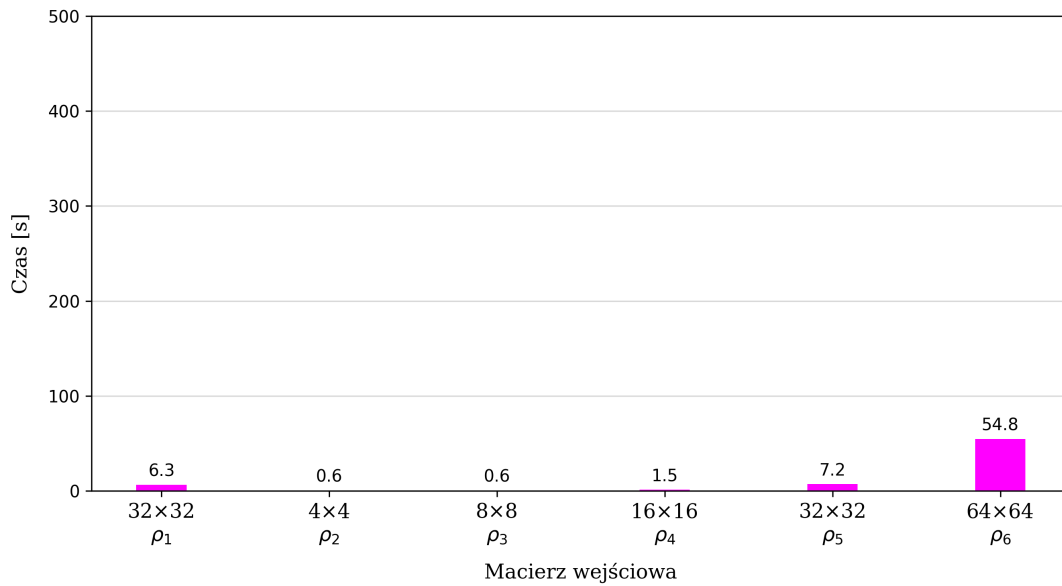
3.5.4 Rust i Ndaray



Rysunek 14: Wyniki testów wydajności implementacji w języku Rust z użyciem biblioteki Ndaray dla macierzy $\rho_2 - \rho_6$ i obliczeniami pojedynczej precyzji.

Implementacja w języku Rust wykorzystująca bibliotekę Ndaray prezentuje najlepszą wydajność podczas obliczeń na małych macierzach, do 16×16 włącznie. Wynika to najprawdopodobniej z braku dodatkowego obciążenia ze strony interpretera, którego wewnętrzne operacje wprowadzają dodatkowe dane do pamięci cache, tym samym wypierając z niej macierze na których są prowadzone obliczenia. W przypadku języka interpretowanego, CPU musi wykonywać o znacznie więcej instrukcji, ponieważ każda operacja w języku wysokiego poziomu musi zostać załadowana, po czym odpowiednia akcja musi zostać wybrana i dopiero wykonana. Wzrost wydajności występuje również dla macierzy 32×32 i 64×64 , natomiast w ich przypadku kluczowa dla wydajności staje się wielomianowa złożoność obliczeniowa algorytmu, którą znacznie lepiej rekompensują wyspecjalizowane implementacje operacji macierzowych zawarte w bibliotece OpenBLAS.

3.5.5 Rust i Ndaray z OpenBLAS



Rysunek 15: Wyniki testów wydajności implementacji w języku Rust z użyciem biblioteki Ndaray dla macierzy ρ_2 - ρ_6 i obliczeniami pojedynczej precyzji.

Zestawienie języka Rust i biblioteki Ndaray z pakietem OpenBLAS i liczbami zmiennoprzecinkowymi pojedynczej precyzji poskutkowało uzyskaniem zaskakująco dobrych wyników wydajności, które zostały przedstawione na rysunku 15. W przypadku macierzy 4×4 i 8×8 wydajność jest bardzo zbliżona do wariantu nie korzystającego z OpenBLAS, natomiast wraz ze wzrostem rozmiaru macierzy, skrócenie czasu pracy staje się coraz bardziej widoczne, osiągając $4.3\times$ przyspieszenie dla macierzy 64×64 .

4 Dyskusja

Odwołania

- [1] J. D. Hunter. „Matplotlib: A 2D graphics environment”. W: *Computing in Science & Engineering* 9.3 (2007), s. 90–95. DOI: 10.1109/MCSE.2007.55.
- [2] Carl Friedrich Bolz i in. „Tracing the meta-level: PyPy’s tracing JIT compiler”. W: *Proceedings of the 4th workshop on the Implementation, Compilation, Optimization of Object-Oriented Languages and Programming Systems*. 2009, s. 18–25.
- [3] Stefan Behnel i in. „Cython: The best of both worlds”. W: *Computing in Science & Engineering* 13.2 (2010), s. 31–39.
- [4] Yury Selivanov Elvis Pranskevichus. *What’s New In Python 3.5*. 2015. URL: <https://docs.python.org/3/whatsnew/3.5.html> (term. wiz. 14.05.2023).
- [5] Siu Kwan Lam, Antoine Pitrou i Stanley Seibert. „Numba”. W: *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*. ACM, list. 2015. DOI: 10.1145/2833157.2833162. URL: <https://doi.org/10.1145/2833157.2833162>.
- [6] Feng Li i in. „CPU versus GPU: which can perform matrix computation faster—performance comparison for basic linear algebra subprograms”. W: *Neural Computing and Applications* 31 (2019), s. 4353–4365.
- [7] Inc. Anaconda i in. *Numba documentation*. 2020. URL: <https://numba.readthedocs.io/en/stable/user/index.html> (term. wiz. 12.05.2023).
- [8] Palash Pandya, Omer Sakarya i Marcin Wieśniak. „Hilbert-Schmidt distance and entanglement witnessing”. W: *Physical Review A* 102.1 (2020), s. 012409.
- [9] Mirko Consiglio, Tony JG Apollaro i Marcin Wieśniak. „Variational approach to the quantum separability problem”. W: *Physical Review A* 106.6 (2022), s. 062413.
- [10] NumPy Developers. *NumPy documentation*. 2022. URL: <https://numpy.org/doc/stable/> (term. wiz. 12.05.2023).
- [11] Klaralvdalens Datakonsult AB. *GitHub - KDAB/hotspot: The Linux perf GUI for performance analysis*. 2023. URL: <https://github.com/KDAB/hotspot> (term. wiz. 31.05.2023).
- [12] The go-python Authors. *go-python/gopy: gopy generates a CPython extension module from a go package*. 2023. URL: <https://github.com/go-python/gopy> (term. wiz. 14.05.2023).
- [13] *Cython C-Extensions for Python*. 2023. URL: <https://cython.org/> (term. wiz. 14.05.2023).
- [14] Matt Davis. *snakeviz · PyPI*. 2023. URL: <https://pypi.org/project/snakeviz/> (term. wiz. 22.05.2023).
- [15] NumPy Developers. *Random Generator — NumPy v1.24 Manual*. 2023. URL: <https://numpy.org/doc/1.24/reference/random/generator.html> (term. wiz. 22.05.2023).

- [16] The PyO3 developers. *PyO3 user guide*. 2023. URL: <https://pyo3.rs/v0.18.3/> (term. wiz. 14.05.2023).
- [17] Agner Fog. *Lists of instruction latencies, throughputs and micro-operation breakdowns for Intel, AMD, and VIA CPUs*. 2023. URL: <https://docs.rs/ndarray/latest/ndarray/index.html> (term. wiz. 22.05.2023).
- [18] Python Software Foundation. *ctypes — A foreign function library for Python*. 2023. URL: <https://docs.python.org/3/library/ctypes.html> (term. wiz. 14.05.2023).
- [19] Python Software Foundation. *Extending Python with C or C++*. 2023. URL: <https://docs.python.org/3/extending/extending.html> (term. wiz. 14.05.2023).
- [20] Ivan Levkivskiy Guido van Rossum. *PEP 483 - The Theory of Type Hints*. 2023. URL: <https://peps.python.org/pep-0483/> (term. wiz. 14.05.2023).
- [21] Łukasz Langa Guido van Rossum Jukka Lehtosalo. *PEP 484 - Type Hints*. 2023. URL: <https://peps.python.org/pep-0484/> (term. wiz. 14.05.2023).
- [22] Jukka Lehtosalo i mypy contributors. *mypy 1.2.0 documentation*. 2023. URL: <https://mypy.readthedocs.io/en/stable/> (term. wiz. 14.05.2023).
- [23] mypyc team. *mypyc 1.2.0 documentation*. 2023. URL: <https://mypyc.readthedocs.io/en/stable/> (term. wiz. 14.05.2023).
- [24] The PyPy Team. *PyPy Home Page*. 2023. URL: <https://www.pypy.org/> (term. wiz. 14.05.2023).
- [25] Zhang Xianyi. *OpenBLAS : An optimized BLAS library*. 2023. URL: <https://www.openblas.net/> (term. wiz. 31.05.2023).