

Metody Programowania Równoległego

Raport I - OpenMPI

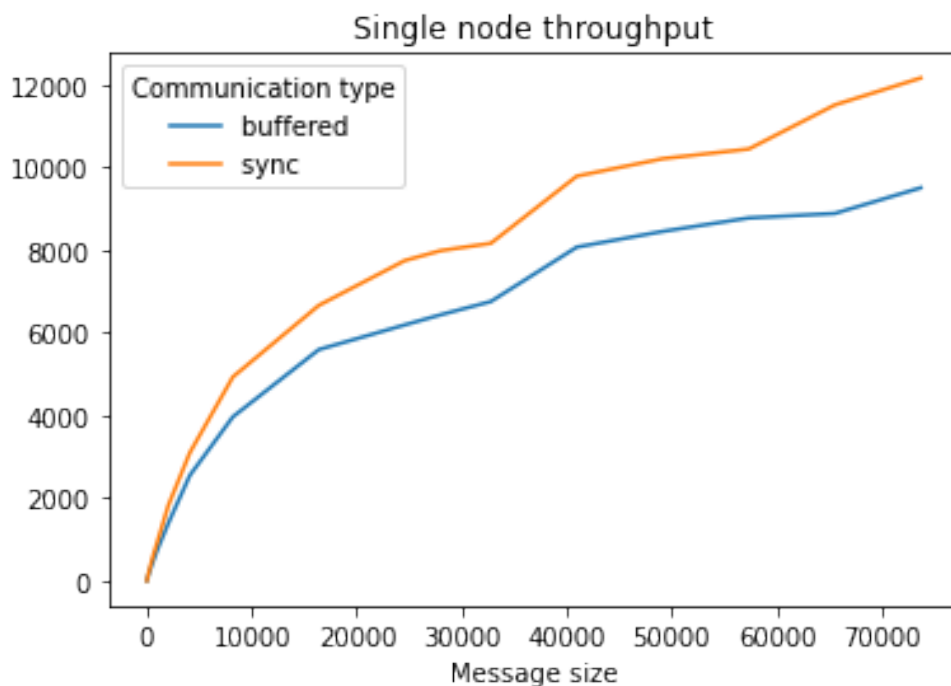
Maciej Trątnowiecki

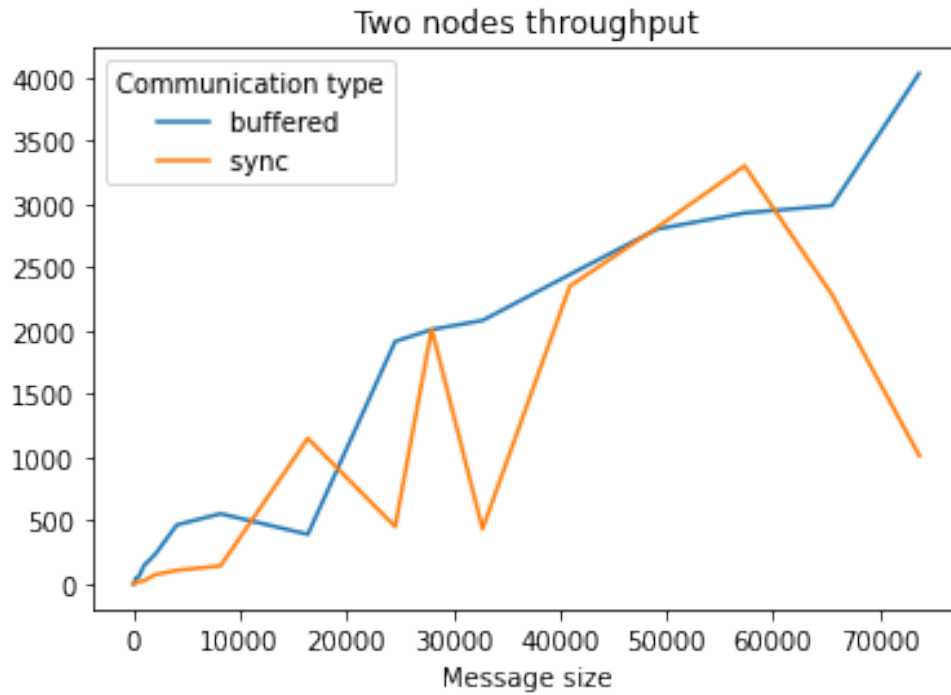
AGH, Semestr Letni, 2022

1 Pomiar przepustowości w komunikacji pomiędzy węzłami klastra

1.1 Pomiary i wyniki

W celu przeprowadzenia doświadczenia przygotowałem skrypt w języku Python, którego kod załączam poniżej. Program przesyła komunikaty pomiędzy dwoma procesami OpenMPI, predefiniowaną ilość razy, następnie liczy przepustowość połączenia. W celu pomiaru przepustowości wielokrotnie uruchomiłem program dla różnych rozmiarów wiadomości i różnych sposobów komunikacji - buforowanej i synchronicznej. Zebrane wyniki pomiarów przedstawiłem w poniższych wykresach.





W pierwszym przypadku komunikacja odbywała się pomiędzy dwoma procesami uruchamianymi na tym samym węźle klastra, w drugim pomiędzy procesami z różnych węzłów. Zaobserwować możemy, że w przypadku gdy procesy wykonują się na różnych węzłach komunikacja buforowana pozwala uzyskać stabilniejszą przepustowość dla wiadomości o sporym rozmiarze. W przypadku gdy komunikują się procesy z tego samego węzła, komunikacja synchroniczna pozwala na uzyskanie mniejszego narzutu czasowego.

Opóźnienie w przesyłaniu wiadomości, zmierzone jako przepustowość dla wiadomości o rozmiarze jednego bajta wyniosło odpowiednio:

Rodzaj komunikacji	Typ wiadomości	Przepustowość [Mbit/s]
Jeden węzeł	buforowana	1.105
Jeden węzeł	synchroniczna	1.37
Różne węzły	buforowana	0.1468
Różne węzły	synchroniczna	0.0407

Tabela 1: Opóźnienie w przesłaniu wiadomości.

1.2 Kod programu

```

1000 #!/usr/bin/env python
1001 from mpi4py import MPI
1002 import numpy as np
1003 import os
1004 from sys import argv
1005
1006 MSG_N = 1000
1007
1008 if argv[1] == 'buff':
1009     BUFF_SIZE = 84000
1010 else:
1011     BUFF_SIZE = 0
1012
1013 MSG_SIZE = int(argv[2])
1014
1015 FIRST_NODE = 0
1016 SECOND_NODE = 1
1017
1018

```

```

def send(comm, *args, **kwargs):
1020     if BUFF.SIZE > 0:
        comm.bsend(*args, **kwargs)
1022     else:
        comm.ssend(*args, **kwargs)
1024

def attach_buffer():
1026     if BUFF.SIZE > 0:
        buf = np.empty((BUFF.SIZE * MSG.N,))
        MPI.Attach_buffer(buf)
1030

def detach_buffer():
1032     if BUFF.SIZE > 0:
        MPI.Detach_buffer()
1034

def get_message():
1036     return os.urandom(MSG.SIZE)
1038

def first_node(comm):
1040     data = get_message()
        for i in range(MSG.N):
1042         send(comm, data, dest=SECOND_NODE)
            data_rec = comm.recv(source=SECOND_NODE)
1044         assert data == data_rec
1046

def second_node(comm):
1048     for i in range(MSG.N):
        data = comm.recv(source=FIRST_NODE)
1050         send(comm, data, dest=FIRST_NODE)
1052

if __name__ == '__main__':
1054     comm_world = MPI.COMM_WORLD
        rank = comm_world.Get_rank()
1056
        attach_buffer()
        comm_world.Barrier()
        time=MPI.Wtime()
1058
        if rank == FIRST_NODE:
            first_node(comm_world)
1060         elif rank == SECOND_NODE:
            second_node(comm_world)
1062         else:
            print "Expected only two nodes"
1064
        time=MPI.Wtime()-time
        detach_buffer()
1066
        if rank == FIRST_NODE:
            capacity = float(MSG.SIZE*MSG.N*8*2)/time
1068             print 'buffered' if BUFF.SIZE > 0 else 'sync', ';', MSG.SIZE, ';', MSG.N, ';', capacity
1070
        MPI.Finalize()
1072
1074
1076

```

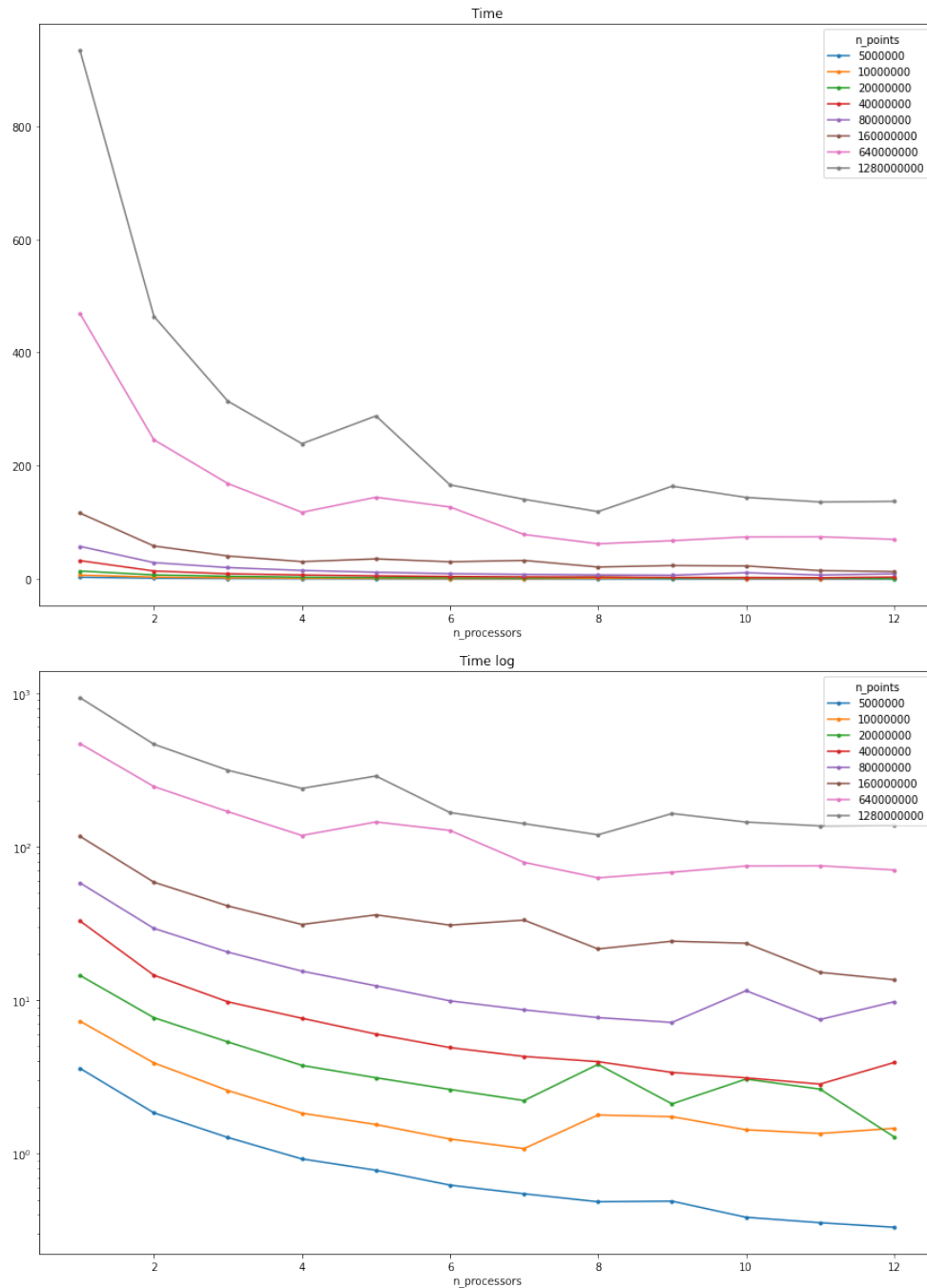
../ex1/ping_pong.py

2 Przybliżenie liczby PI metodą Monte Carlo

W ramach eksperymentu przygotowałem skrypt obliczający przybliżenie liczby PI metodą Monte Carlo w sposób równoległy.

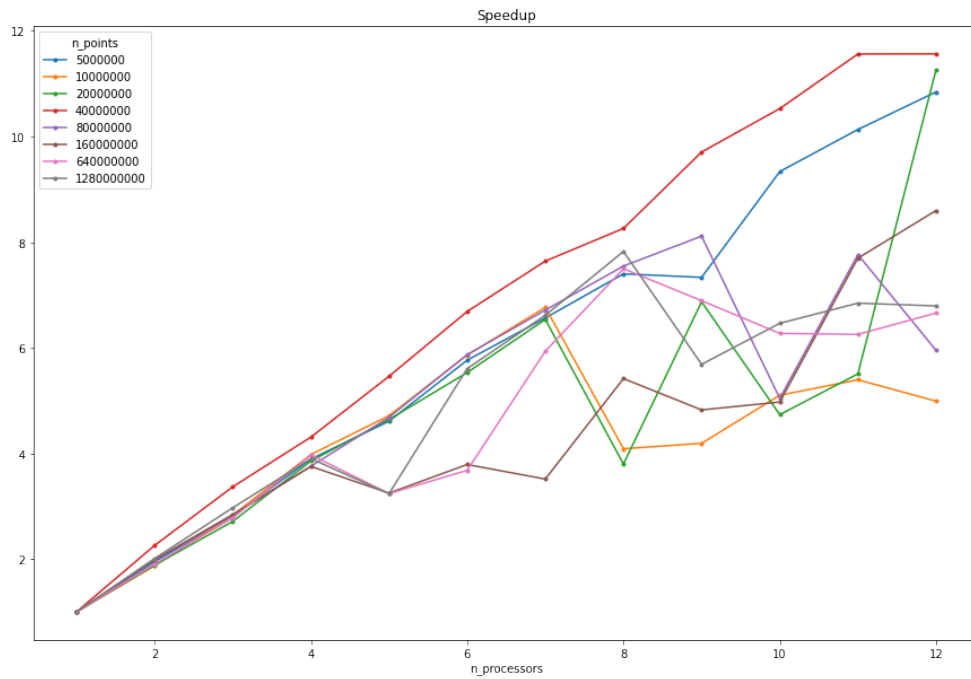
2.1 Zależność czasu od liczby procesorów

Skrypt uruchomiłem dla problemów o rosnącej złożoności dla liczby procesorów z zakresu 1 do 12 i zmierzyłem czasy wykonania. Poniższe wykresy przedstawiają zmianę czasu wykonania programu w zależności od liczby procesorów. Drugi z wykresów przedstawiony został w skali logarytmicznej.



2.2 Przyspieszenie

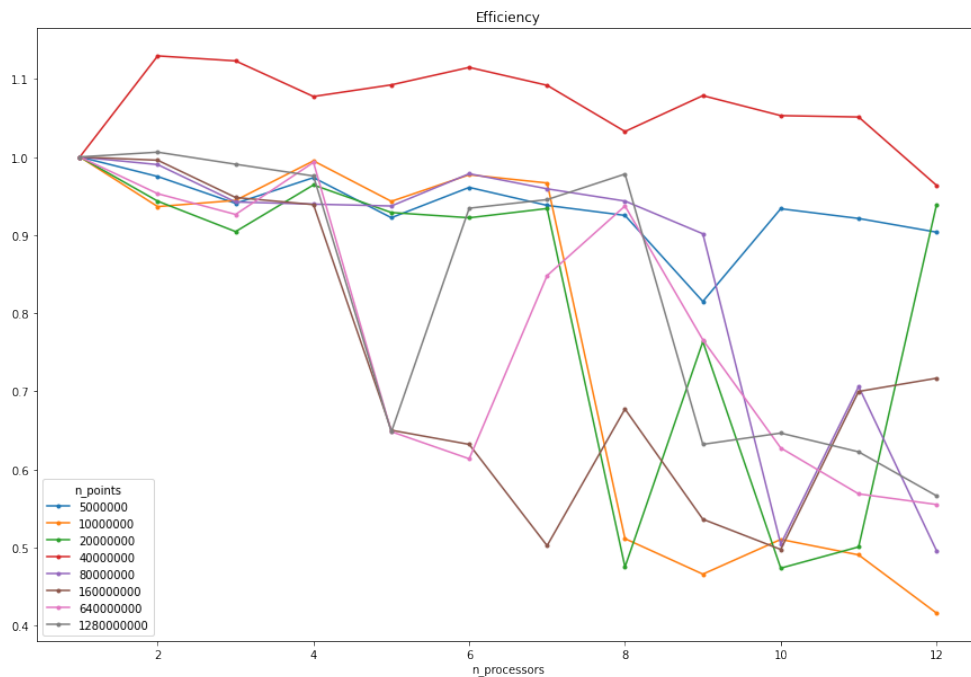
Na podstawie otrzymanych wcześniej pomiarów czasu wykonania programu, przygotowałem wykres przyspieszenia w zależności od liczby procesorów dla problemów o różnym rozmiarze. Przyspieszenie obliczyłem jako stosunek czasu wykonania programu bez współbieżności dla otrzymanych wyników dla wyższej liczby procesorów.



Zauważyć możemy, że dla większości rozmiarów problemu wykres przyspieszenia odbiega od teoretycznie idealnego liniowego wzrostu $y=x$.

2.3 Efektywność

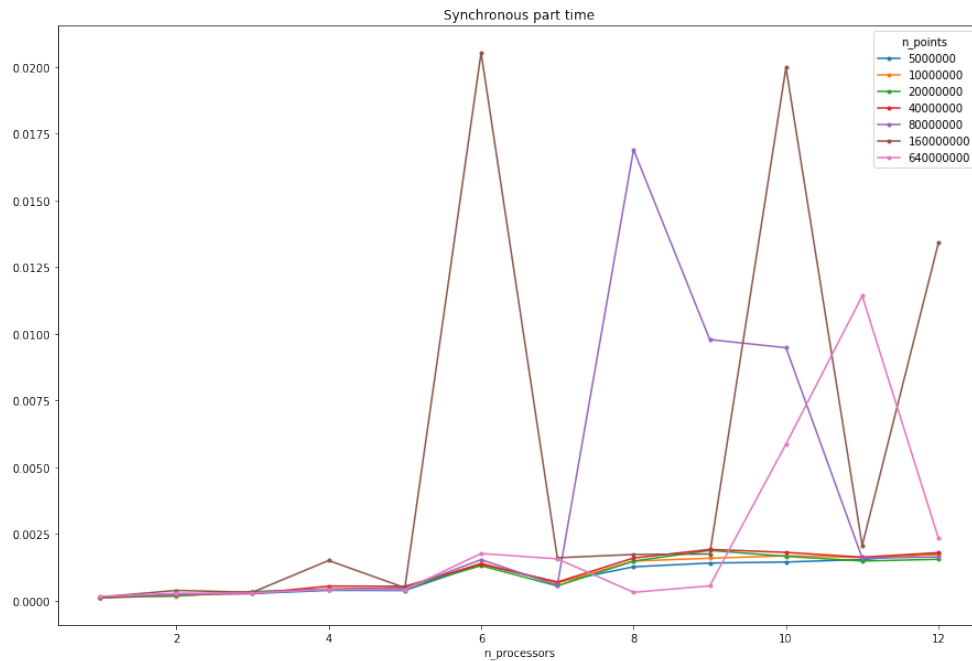
Na podstawie otrzymanych w poprzednim kroku przyspieszeń, obliczyłem wykres efektywności jako stosunek przyspieszenia do ilości procesorów.



Wbrew oczekiwaniom, zaobserwować możemy spadek wydajności nawet dla problemów o znaczącym rozmiarze.

2.4 Pomiar czasu wykonania części synchronicznej

Wykonałem także pomiar części synchronicznej programu - odpowiedzialnej za zbieranie wyników obliczeń. Wyniki przedstawiłem na poniższym wykresie.



Zauważyć możemy, że czas poświęcony na zbieranie wyników w ogólności nie rośnie wraz ze wzrostem rozmiaru problemu obliczeniowego. Wynika to z faktu, że zbieranym wynikiem jest jedynie liczba punktów wylosowanych wewnątrz koła, dlatego ilość danych nie rośnie wraz ze wzrostem problemu.

2.5 Kod programu

```

1000 #!/usr/bin/env python
1001 from mpi4py import MPI
1002 import numpy as np
1003 import os
1004 from sys import argv
1005 import random
1006
1007 N_POINTS = int(argv[1])
1008 N_PROC = int(argv[2])
1009
1010 def get_message():
1011     return os.urandom(MSG_SIZE)
1012
1013 def point():
1014     return (random.random(), random.random())
1015
1016 def point_in_circle(point):
1017     x, y = point
1018     return (x * x + y * y) < 1
1019
1020 def node(comm):
1021     res = 0
1022     node_points = N_POINTS / comm.Get_size()
1023     i = 0
1024     while i < node_points:
1025         if point_in_circle(point()):
1026             res += 1
1027         i += 1
1028     return res
1029
1030 if __name__ == '__main__':
1031     ROOT_NODE = 0
1032
1033     comm = MPI.COMM_WORLD
1034     rank = comm.Get_rank()

```

```

1042 comm.Barrier()
1043 time=MPI.Wtime()
1044
1045 node_res = node(comm)
1046
1047 comm.Barrier()
1048 time_calc=MPI.Wtime()-time
1049 time = MPI.Wtime()
1050
1051 global_sum = 0
1052 global_sum = comm.reduce(node_res, global_sum, op=MPI.SUM, root=ROOTNODE)
1053
1054 time=MPI.Wtime()-time
1055
1056 if rank == ROOTNODE:
1057     print NPROC, ';', global_sum, ';', NPOINTS, ';', 4.0 * float(global_sum)/float(
1058         NPOINTS), ';', time_calc, ';', time
1059
1060 MPI.Finalize()

```

../ex2/pi.py

```

1000 #!/bin/bash -l
1001 #SBATCH --nodes 1
1002 #SBATCH --ntasks 12
1003 #SBATCH --time=01:00:00
1004 #SBATCH --partition=plgrid
1005 #SBATCH --account=plgmr22
1006 #SBATCH --sockets-per-node=2
1007
1008 module add plgrid/tools/openmpi
1009 module spider mpi4py
1010 module add plgrid/libs/python-mpi4py/3.0.0-python-2.7
1011
1012 for points in 1000000 36000000 1280000000 ; do
1013     for n in 1 2 3 4 5 6 7 8 9 10 11 12 ; do
1014         echo n $n >> res.txt
1015         mpiexec -np 12 ./pi.py $points >> res.txt
1016     done
1017 done

```

../ex2/script.sh