

Caixeiro Viajante (Multi-Core-GPU)

Introdução

Este projeto consiste no estudo de técnicas multi-core na área de otimização discreta, que se aplica a problemas em que há uma sequência de escolhas, com a característica de que a solução ótima só pode ser encontrada ao enumerar todas as escolhas possíveis. Isto significa que todo algoritmo para sua solução é $O(2^n)$ ou pior.

Neste projeto serão implementados dois métodos de otimização discreta, com processamento multi-core em GPU, estes são:

- Solução Aleatória: Consiste em sortear um número de soluções aleatoriamente e selecionar a que tiver menor custo, esta é uma solução rápida e simplória, com resultados não ideais, porém trata-se de uma solução inicial satisfatória.
- 2-opt (local search): Consiste em aplicar a heurística 2-opt sobre o método de solução aleatória, aplicando operações em cada iteração para melhorar o parâmetro de critério. Seu valor final, em algumas execuções, será apenas próximo do valor ideal, pois esse método é afetado por mínimos locais, logo o processo de inicialização aleatória e melhoria é executado diversas vezes, para reduzir esse efeito colateral.

Esse estudo de otimização já foi realizado em outro projeto porém aplicando técnicas multicore em CPU, A fim de tornar a análise de otimização mais rica, os resultados obtidos no projeto em CPU serão adicionados à seção de análises para serem comparados com os resultados gerados pela implementação em GPU. A fim de concretizar possíveis vantagens e desvantagens de cada aplicação, definindo a que possui melhor desempenho sobre tal tema.

[Projeto em CPU](#)

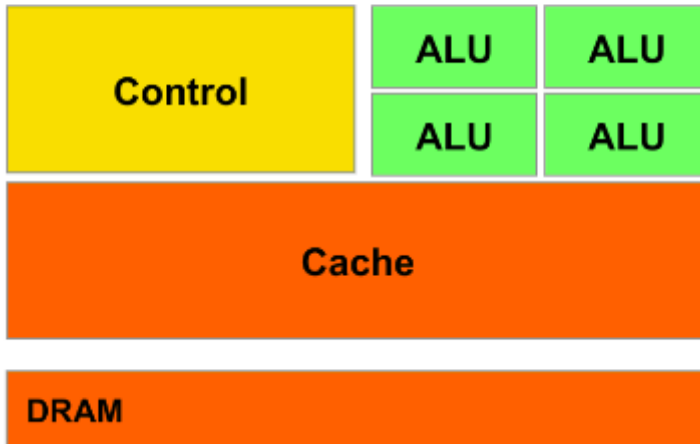
GPU vs CPU

Enquanto a CPU possui um design orientado a minimizar a latência, unidades lógicas aritméticas potentes, caches grandes, controle sofisticado e uma quantidade reduzida de cores. As GPUs possuem um design orientado a minimizar o throughput, com unidades aritméticas simples, caches reduzidos e uma quantidade massiva de cores.

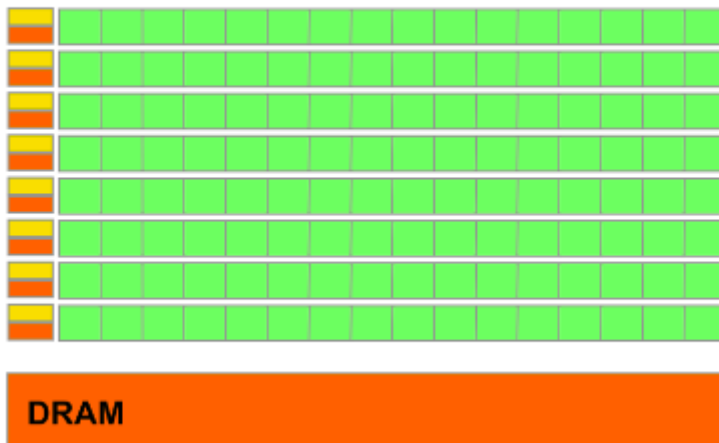
Dada essa diferença, as CPUs podem ser até 10x mais rápidas que as GPUs para código sequencial, enquanto as GPUs podem ser 10x mais rápidas que as CPUs para códigos

paralelos.

CPU



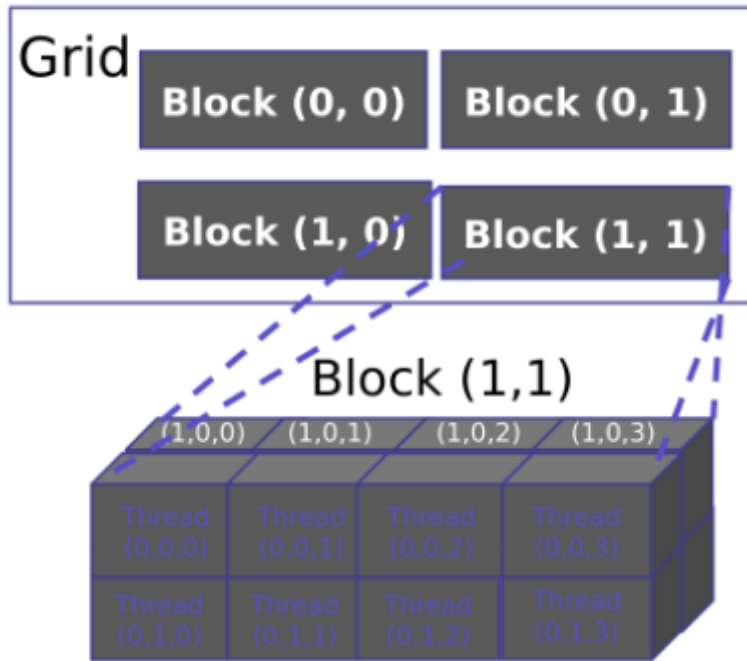
GPU



Para utilizar todo o poder de processamento da GPU utiliza-se o CUDA(Compute Unified Device Architecture), que é uma engine de computação desenvolvida pela Nvidia, que trabalha sobre a estrutura de divisão de memória da GPU em que a thread é a menor unidade de tarefa, essas threads são agrupadas em blocos, onde as threads de um mesmo bloco podem se comunicar e se sincronizar.

Os blocos por sua vez são organizados em grids, onde deverão ter a mesma dimensão, executando a mesma função de kernel, podendo a grid ser bidimensional ou tridimensional.

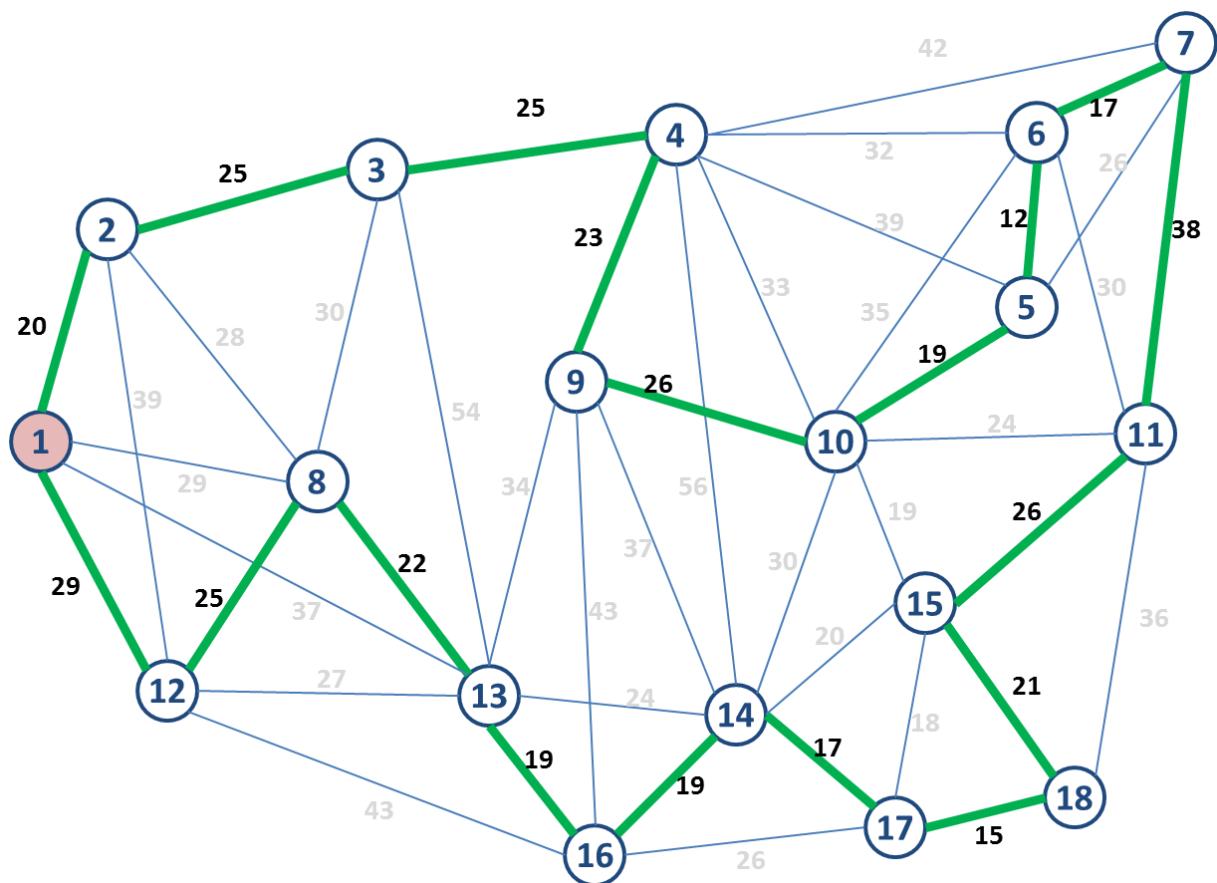
Device (GPU)



O Problema

Para aplicar os conceitos de otimização discreta foi implementado o problema do Caixeiro Viajante, que consiste em determinar a menor rota para percorrer uma série de cidades(sem repetições) e retornar a cidade de origem.

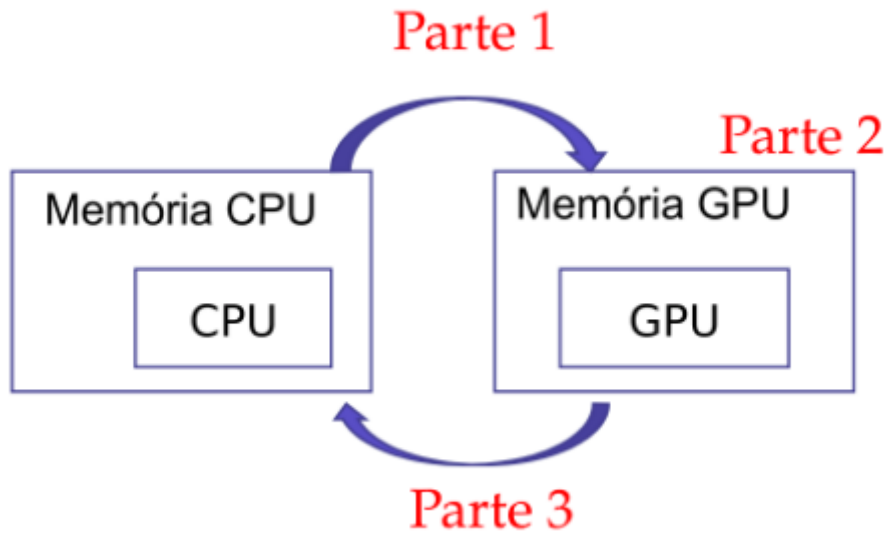
Este problema foi inspirado na necessidade dos vendedores em realizar entregas em diversos locais, percorrendo o menor caminho possível, e consequentemente reduzindo o tempo necessário para a viagem e os possíveis custos com transporte e combustível.



O projeto está estruturado de forma que cada método implementado possui seu respectivo arquivo.cpp, esses arquivos possuem uma função main, responsável por definir o fluxo de dados entre CPU e GPU. Para aplicar o método de otimização discreta de forma paralela sobre o conjunto de dados, há dois kernels, definidos pela flag **global**:

- Também há duas funções auxiliares em GPU, definidas pela flag **device** :

- 4/12



- o Parte1: Copia dados CPU -> GPU
- o Parte2: Processa dados na GPU
- o Parte3: Copia resultados GPU -> CPU

Testes

Os inputs selecionados para testar o desempenho da aplicação de técnicas de otimização discreta foram formulados de forma a tentar levar ao máximo o processamento da GPU. Como o programa exige mais processamento quando há muitos pontos que o caixairo deve percorrer, esse valor foi testado de forma crescente, e com grandes intervalos de iterações máximas, a fim de tentar tornar mais visível o ganho de desempenho ao aplicar essa técnica para diferentes demandas de processamento.

Os seguintes inputs foram testados:

- o Input1: N = 200
- o Input2: N = 300
- o Input3: N = 400
- o Input4: N = 500

Em que N representa o número de pontos no sistema. Como o método implementado constrói diversas soluções, numericamente definido como 10000 no projeto, o número de elementos utilizados é $N * 10000$.

Definido como variável de desempenho o tempo de execução do programa, foi utilizada a estrutura `cudaEvent_t` para mensurar esta informação, estrutura esta que permite obter este parâmetro no contexto da GPU.

Como deseja-se medir apenas o desempenho do algoritmo do caixeiro viajante a medição não leva em conta o tempo gasto na leitura do arquivo de input e outros processos similares.

Resultados

GPU (Tempo)

```
##import dependences
import matplotlib.pyplot as plt
import pandas as pd
import subprocess
import sys
from IPython.display import display

##files
dir_fil = "../build/"
files    = ["random-sol", "2-opt-sol"]

dir_inp = "../inputs/"
inputs  = ["input1", "input2", "input3", "input4"]

n_sizes = []
dic      = {}

## Run files and storage output
for input in inputs:
    dic[input] = {}

    for f in files:
        command = dir_fil + f + " < " + dir_inp + input
        stout    = subprocess.check_output(command, stderr=subprocess.STDOUT,
        dic[input][f] = (float(stout.split("\n")[-2].split(":")[1]), float(s

    with open(dir_inp + input, 'r') as f:
        nrect = f.read().split("\n")[0]

    n_sizes.append(int(nrect))

## Generate Dataframe
s0 = [dic["input1"]["random-sol"][0], dic["input2"]["random-sol"][0], dic["ir

s1 = [dic["input1"]["2-opt-sol"][0], dic["input2"]["2-opt-sol"][0], dic["inpu

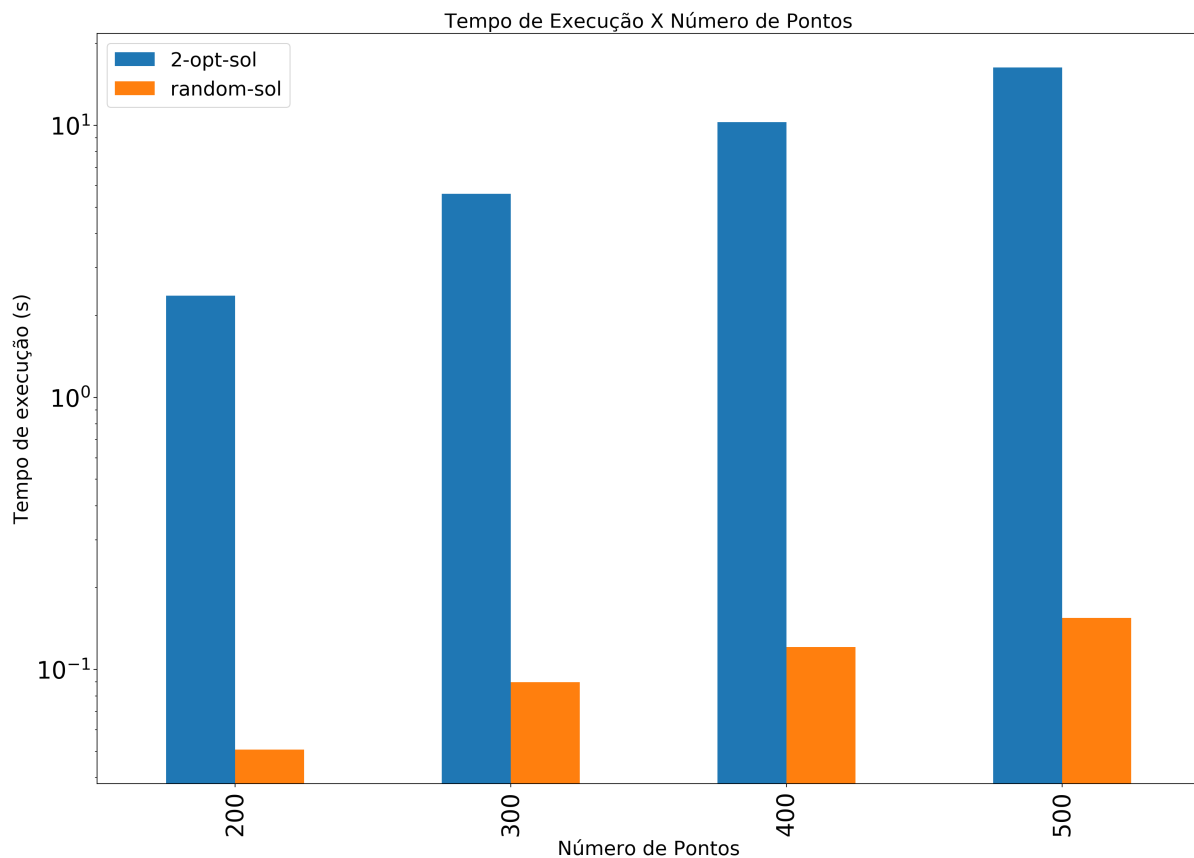
df = pd.DataFrame({"random-sol":s0, "2-opt-sol":s1}, index = n_sizes)

display(df)
```

```
df.plot.bar(figsize = (20,14), fontsize=24, logy = True)
plt.legend(prop={'size': 20})
plt.title("Tempo de Execução X Número de Pontos", fontsize = 20)
plt.xlabel("Número de Pontos", fontsize = 20)
plt.ylabel("Tempo de execução (s)", fontsize = 20)
```

	2-opt-sol	random-sol
200	2.36230	0.050745
300	5.60886	0.089715
400	10.27450	0.120887
500	16.30620	0.154622

```
Text(0, 0.5, 'Tempo de execução (s)')
```



GPU (Distância)

```
s0 = [dic["input1"]["random-sol"][1], dic["input2"]["random-sol"][1], dic["ir
s1 = [dic["input1"]["2-opt-sol"][1], dic["input2"]["2-opt-sol"][1], dic["inpu
```

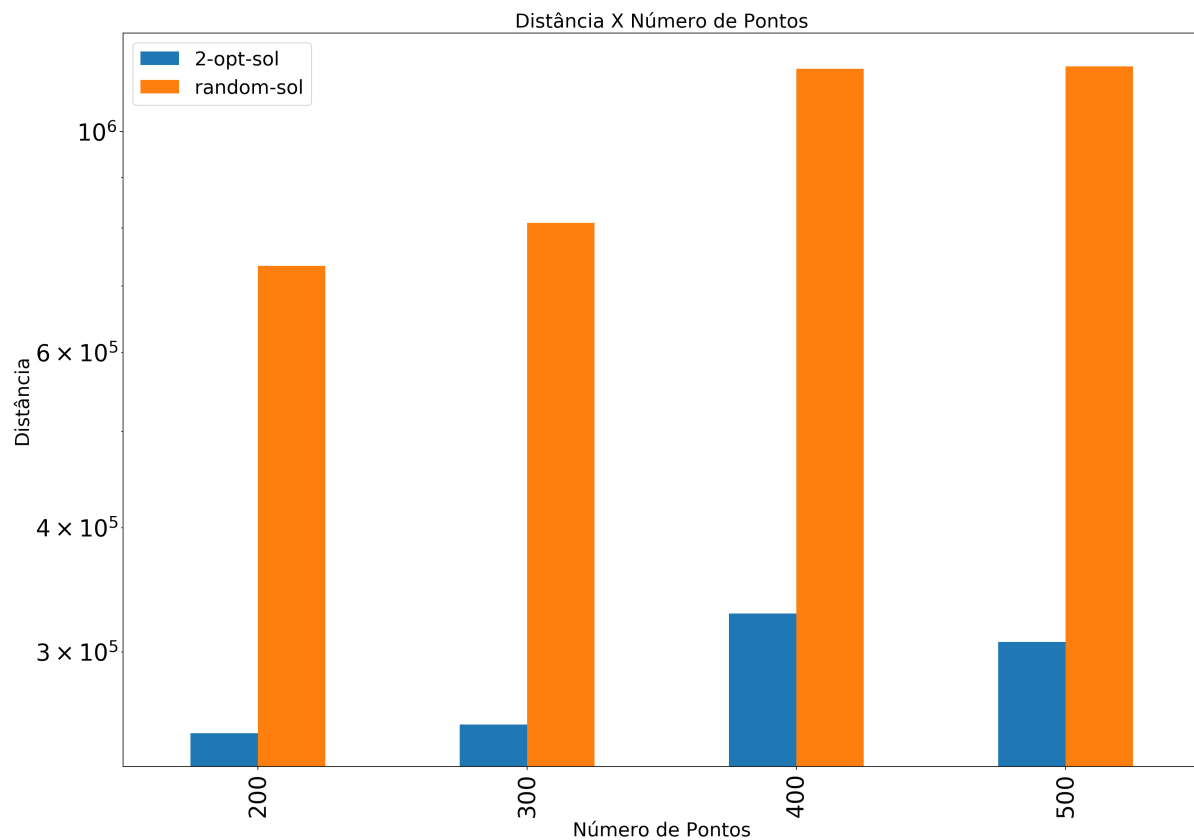
```
df = pd.DataFrame({"random-sol":s0, "2-opt-sol":s1,}, index = n_sizes)

display(df)

df.plot.bar(figsize = (20,14), fontsize=24, logy = True)
plt.legend(prop={'size': 20})
plt.title("Distância X Número de Pontos", fontsize = 20)
plt.xlabel("Número de Pontos", fontsize = 20)
plt.ylabel("Distância", fontsize = 20)
```

	2-opt-sol	random-sol
200	248553.23263	7.329119e+05
300	253712.12790	8.101004e+05
400	327883.65939	1.156667e+06
500	306977.37416	1.163051e+06

```
Text(0, 0.5, 'Distância')
```



CPU vs GPU


```
s0 = [0.594745, 1.23814, 2.48137, 3.70665] ## Valores obtidos com o projeto e
s1 = [dic["input1"]["2-opt-sol"][0], dic["input2"]["2-opt-sol"][0], dic["input3"]["2-opt-sol"][0]]

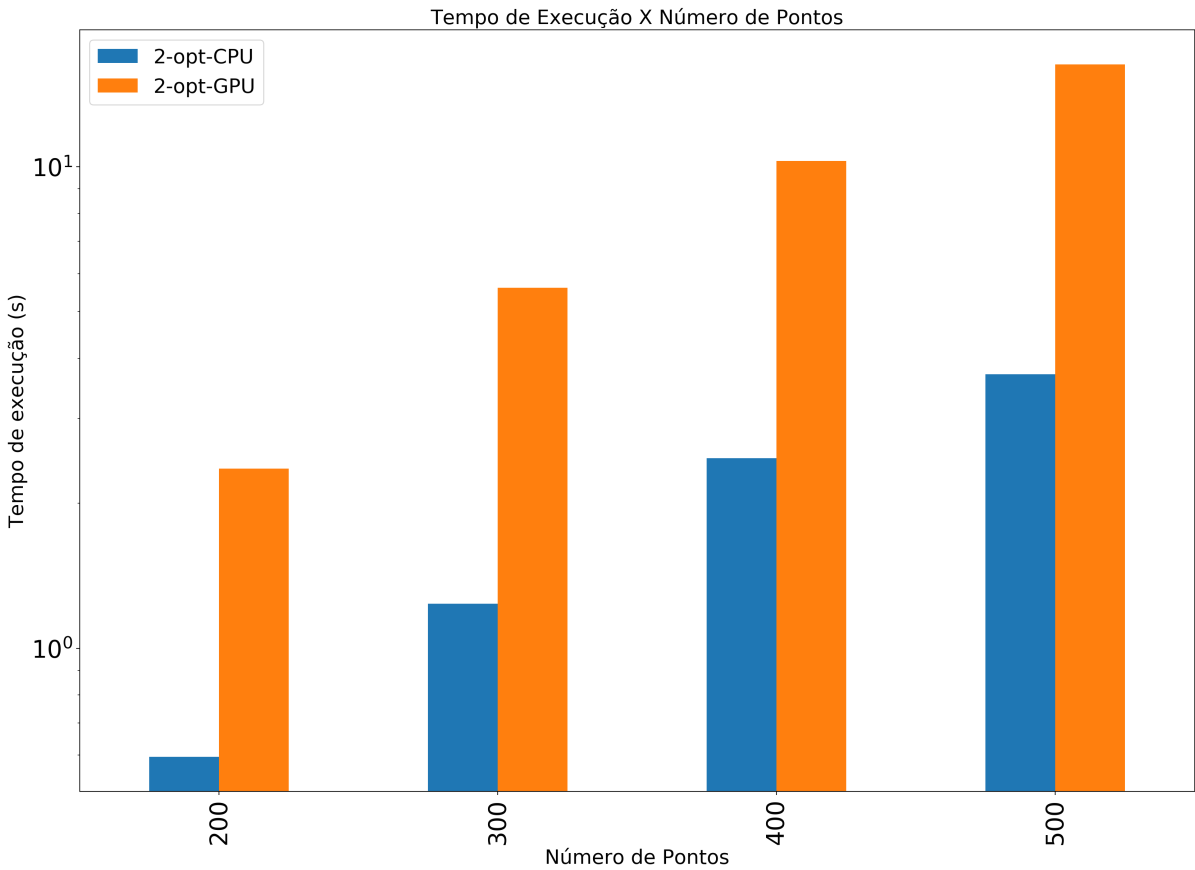
df = pd.DataFrame({"2-opt-CPU":s0, "2-opt-GPU":s1,}, index = n_sizes)

display(df)

df.plot.bar(figsize = (20,14), fontsize=24, logy = True)
plt.legend(prop={'size': 20})
plt.title("Tempo de Execução X Número de Pontos", fontsize = 20)
plt.xlabel("Número de Pontos", fontsize = 20)
plt.ylabel("Tempo de execução (s)", fontsize = 20)
```

	2-opt-CPU	2-opt-GPU
200	0.594745	2.36230
300	1.238140	5.60886
400	2.481370	10.27450
500	3.706650	16.30620

Text(0, 0.5, 'Tempo de execução (s)')



Conclusão

Como pode-se observar nos resultados obtidos ao comparar os métodos de GPU para o tempo de execução, tem-se que a solução 2-opt possui um tempo de execução maior, o que já era esperado, pois este método consiste em acrescentar diversas operações ao método de solução aleatória, porém conforme o obtido pelo gráfico de GPU para distância a solução 2-opt sempre obtêm valores melhores de custo, aproximando-se do valor ideal, devido às operações de melhoria do parâmetro de critério.

Ao comparar o método 2-opt, considerado como um local search, implementado em GPU com o mesmo método implementado em CPU esperava-se que o código em GPU gera-se um desempenho muito melhor, visto que é possível atingir um nível muito maior de paralelização com esta estrutura, e também foi observado que o tempo de transferência dos dados entre a CPU e a GPU não foram expressivos para interferir em seu desempenho.

Essa afirmação foi obtida com o uso do nvprof, pelo qual detectou-se que na totalidade o tempo de execução do programa está contido no processos de cálculo da distância na GPU, representando 99,8% do tempo total de execução, para todas as entradas apresentadas, com uma variação de apenas 0.05%.

Porém conforme o obtido, a implementação em CPU possui melhor desempenho, isso pode ser justificado pela forma como o código está implementado, cada thread da GPU está acessando a memória global do Device (acesso lento), para corrigir este problema pode-se substituir essa leitura para a memória compartilhada de cada bloco (acesso rápido), e balancear a carga que cada thread realiza, esta é uma próxima etapa para o projeto.

Especificações da Máquina utilizada

CPU

Arquitetura: x86_64

Modo(s) operacional da CPU: 32-bit, 64-bit

Ordem dos bytes: Little Endian

CPU(s): 4

Lista de CPU(s) on-line: 0-3

Thread(s) per núcleo: 2

Núcleo(s) por soquete: 2

Soquete(s): 1

Nó(s) de NUMA: 1
ID de fornecedor: GenuineIntel
Família da CPU: 6
Modelo: 61
Nome do modelo: Intel(R) Core(TM) i7-5500U CPU @ 2.40GHz
Step: 4
CPU MHz: 2430.779
CPU MHz máx.: 3000,0000
CPU MHz mín.: 500,0000
BogoMIPS: 4788.89
Virtualização: VT-x
cache de L1d: 32K
cache de L1i: 32K
cache de L2: 256K
cache de L3: 4096K
CPU(s) de nó NUMA: 0-3

GPU

Descrição : VGA compatible controller
Produto : HD Graphics 5500
Fabricante : Intel Corporation
ID físico : 2
Informações do barramento : pci@0000:00:02.0
Versão : 09
Largura : 64 bits
Clock : 33MHz
Capacidades : msi pm vga_controller bus_master cap_list rom
Configuração : driver=i915 latency=0

Descrição : 3D controller
Produto : GK208BM [GeForce 920M]
Fabricante : NVIDIA Corporation
ID físico : 0
Informações do barramento : pci@0000:03:00.0
Versão : a1
Largura : 64 bits
Clock : 33MHz
Capacidades : pm msi pciexpress bus_master cap_list rom
Configuração : driver=nvidia latency=0

