

モンテカルロ木探索における シミュレーションの GPU を用いた並列化

明治大学理工学部情報科学科
知的情報処理システム研究室
4 年 15 組 28 番 高野昂平

目次

1	はじめに	2
1.1	研究概要	2
1.2	研究背景	2
2	関連研究	2
3	モンテカルロ木探索について	2
3.1	UCT について	3
3.2	並列化手法	3
4	実験環境	4
4.1	パラメータについて	4
4.2	並列化環境について	4
4.3	実行環境	4
5	実行速度の計測	4
5.1	実験概要	4
5.2	実行時間の比較	5
5.3	並列処理の実行時間の増加について	5
6	強さの比較	6
6.1	実験概要	6
6.2	検定について	6
6.3	勝率の評価	6
7	終わりに	7

1 はじめに

1.1 研究概要

モンテカルロ木探索 (Monte Carlo Tree Search, MCTS) におけるシミュレーション部分を並列化し、高速化を目指す。オセロを対象に実験を行った。

1.2 研究背景

モンテカルロ木探索の処理速度向上を図るうえで、シミュレーション部分が全体の処理速度のボトルネックとなっているのではと考え、シミュレーション部分を並列化することが全体の高速化につながるのではないかと考えた。

高性能な CPU を用いることで高速化を図ることも考えられるが、近年 CPU の性能向上率が以前に比べて伸び悩んでおり、高速化を図るためには CPU に依存した手法では高止まりになってしまう。そこで、別のアプローチで高速化を目指す必要があるが、その手法として GPU を用いた並列化を選択した。近年、GPU を画像処理だけでなく、汎用計算に用いて GPGPU (General Purpose GPU) として活用する動きも見られるため、モンテカルロ木探索に対しても GPGPU として活用できないものかと考えた。

2 関連研究

3 モンテカルロ木探索について

本研究のテーマであるモンテカルロ木探索 (Monte Carlo Tree Search, MCTS) とは、乱数を用いた計算手法であるモンテカルロ法を木探索に応用した手法である。

まず、評価関数を用いないという特徴を持つため、様々なゲームのプレイヤーに対してモンテカルロ木探索を実装することができる。例えば、General Game Playing (GGP) という汎用ゲームに対して人間の介入なしで強いプレイヤーを作るという試みがあるが、そのプレイヤーとしてモンテカルロ木探索を用いることが可能である。

また、不要な探索を行わないという特徴を持つ。合法手が膨大なゲームに対して木探索を行う場合、そのすべてのノードを探索することによって強いプレイヤーの作成するのは計算量が多くなってしまい、現実的でない。例えば、本研究ではオセロのプレイヤーに対して、モンテカルロ木探索を施したが、そのオセロの合法手は 10^{28} 手あると推測されており、すべてのノードを探索することが困難であることがわかる。しかし、モンテカルロ木探索の場合、すべてのノードを探索するのではなく、有効と考えられる手を優先的に深く探索するため、無駄な探索をすることがなく、計算量を削減することができる。

モンテカルロ木探索の流れとしては、選択、シミュレーション、展開、逆伝播の 4 ステップからなっており (図 1)、本研究では 2 ステップ目のシミュレーションを並列化することを目標としている。

3.1 UCT について

選択ステップにおいて、次に訪問するノードを決めるアルゴリズムとして、UCT アルゴリズムを採用した。このアルゴリズムは、以下の式が最大となるようなノードを次の訪問ノードとするアルゴリズムである。

以下の式において、 w はそのノードにおける評価値、 n はそのノードにおける訪問回数、 N はそのノードを親としたときの子ノードの n の値の合計値、 C は任意の値で探索がうまくいくように調整すべきパラメータとなっている。

$$UCT = \frac{w}{n} + C \sqrt{\frac{\ln N}{n}} \quad (1)$$

この式が意味するのは、ノードの評価値だけで訪問先を決めるのではないということ、第 1 項自体は評価値の平均を取っているが、第 2 項において訪問回数が少ないノードを訪問するように UCT 値を調整しているのがわかる。変数 C はこの訪問回数の少ないノードへの訪問をどれくらい重く見るべきか定めるパラメータと言える。

3.2 並列化手法

本研究では、シミュレーション部分を並列化するが、具体的には以下の図 1 のように、4 ステップあるモンテカルロ木探索の 2 ステップ目を複数回のシミュレーションを同時に実行することで実現する。

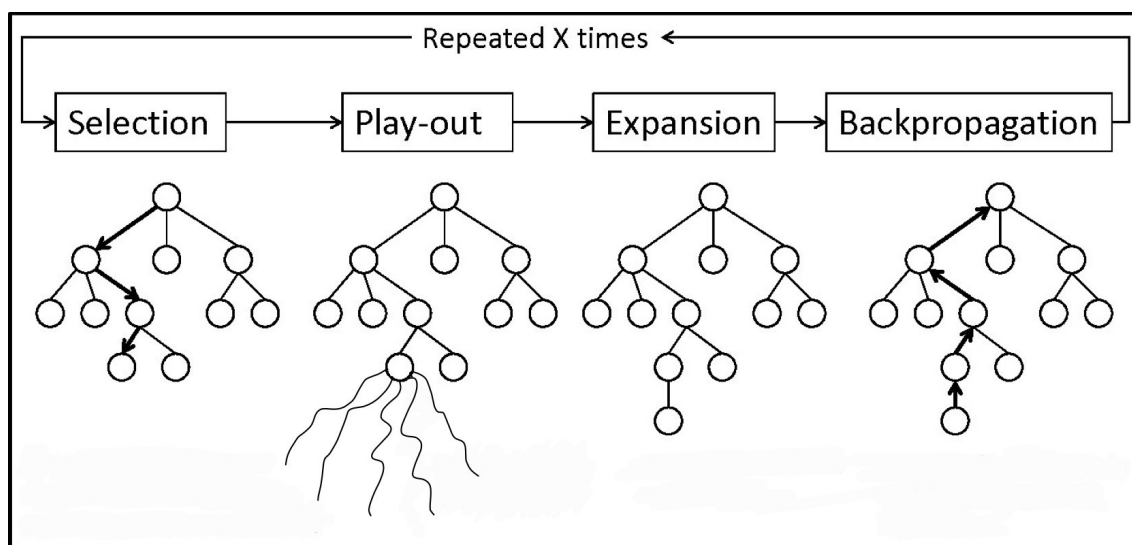


図 1 シミュレーションの並列化

複数回のシミュレーションを行う際、CPU での逐次実行と GPU での並列実行の違いは、まず、逐次実行の場合には図 2 のように 1 回のシミュレーションを for 文によって複数回順番に実行し、それぞれのシミュレーションで得た値を足し合わせることで処理を行っている。

一方、並列実行の場合、図 3 のように 1 回のプレイアウト関数の呼び出しによって、複数回のシミュレーションを終えている。厳密には、`playout_cuda` 関数内でのカーネル呼び出しによって並列にシミュレーショ

ンが行われる。

```
for (int i = 0; i < n_playout; i++){           // n_playoutはシミュレーション回数
    value += Node::playout(tmp, base_player);  // tmpは盤面情報の一時保管場所
}                                              // base_playerは基準となるプレイヤー
```

図2 CPU による逐次実行

```
value += playout_cuda(tmp, base_player);
```

図3 GPU による並列実行

4 実験環境

4.1 パラメータについて

ある盤面ノードに対する評価回数 (図1における X) を 100 回、ノードの展開の基準値を 20 回訪問とした。また、CUDA カーネルを 1 ブロックあたりのスレッド数を 512 として、実験を行った。

4.2 並列化環境について

並列化にあたって、NVIDIA 社の CUDA(Compute Unified Device Architecture) を用いた。CUDA は NVIDIA 社の GPU 上で動作する汎用並列コンピューティングプラットフォームである。なお、言語は C++ とし、CUDA C を用いて実装を行った。

4.3 実行環境

CUDA を用いる関係上、GPU は NVIDIA 社のものとなっている。

OS	Ubuntu 20.04.2 LTS
CPU	Intel(R) Core(TM) i9-10900X CPU @ 3.70GHz
GPU	GeForce RTX 3080 (CUDA コア 8704)
メモリ	32GB

5 実行速度の計測

5.1 実験概要

実験の流れとしては、オセロの初期盤面における評価 (図1における 1 サイクル) を 100 回行う。それぞれの評価を行う中で、CPU でのシミュレーションと GPU でのシミュレーションを実行する。各シミュレー

ションで実行時間を計測し、100 回分の実行時間を計測し終わった後、それらの平均を求める。この平均をもとに実行時間の比較を行った。

5.2 実行時間の比較

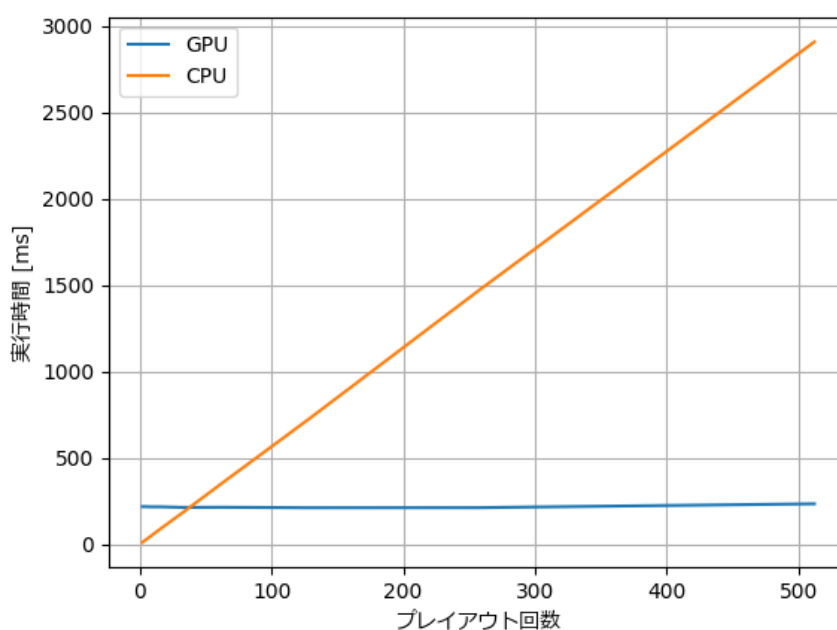


図 4 実行時間の比較

図 4 はプレイアウト数を 1 から順に 2 倍ずつ増やしたときの、実行時間の変化をグラフにしたものである。

図 4 からわかるように、CPU の逐次実行での実行時間が線形的に増加しているのに対して、GPU での実行時間は並列処理であるため、実行時間が増えていない。

しかし、CUDA コア数以上の並列数 (プレイアウト数) で実行するとすべてのプレイアウトを同時に処理できず、並列処理自体を順次実行していくため、実行時間が増加していく。これについては後述する。

5.3 並列処理の実行時間の増加について

図 5 は、さらにプレイアウト数を増やした場合の GPU での実行時間の変化をグラフにしたものである。本実験ではプレイアウト数を 2 倍ずつ増やしているため、その実行時間をグラフにした際に、線形的に増加しているように見える。しかし、プレイアウト数を細かく変化させ、実行時間を計測した場合、天井関数に近い形のグラフになることが予想される。予想される天井関数の概形が図 5 における点線部分にあたる。

並列に実行した場合でも、プレイアウト数を増やしていくと CPU での逐次実行と同様に実行時間が増えていくことがグラフから見て取れるが、しかし、その増加幅は CPU での逐次実行に比べて緩やかなものとなっており、このままプレイアウト数を増やしたとしても、並列処理での実行時間の短縮は実現できるものと考えられる。

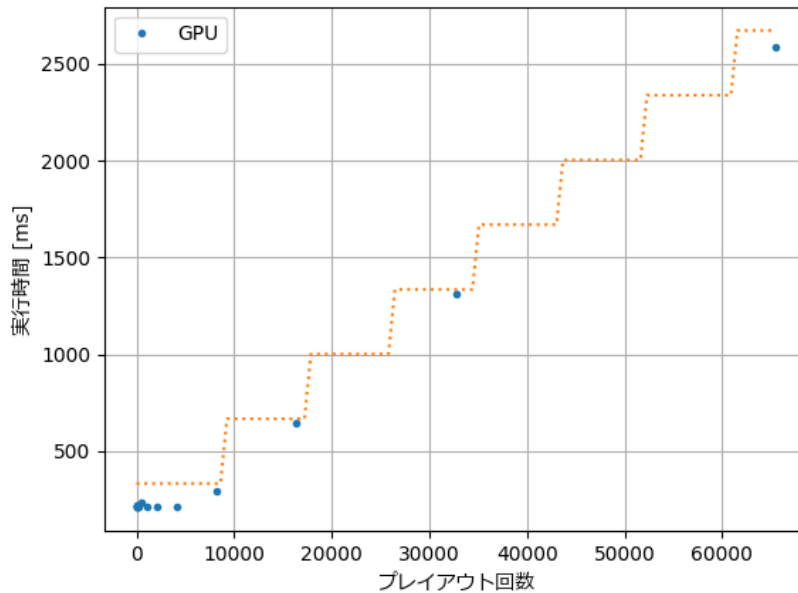


図5 GPUでの実行時間

6 強さの比較

6.1 実験概要

プレイアウト数を増やすことでプレイヤーがどれくらい強くなるのか実験を行った。

後手白の強さを固定し、先手黒のプレイアウト数を増やしていったときにどのように勝率が変化するかを見ていくことで強さの評価を行った。後手白の強さはプレイアウト回数を 2048 回とした。

6.2 検定について

本実験を行う上で t 検定を行い、勝率を評価した。帰無仮説を”先手黒と後手白の強さは同じ”とし、また、対立仮説を”白が黒より強い”とした。対立仮説をこのようにしたのは、今回の実験において後手白のプレイアウト回数が先手黒のプレイアウト回数より多いことが保証されており、白の強さが黒を下回ることは予想されないためである。この対立仮説を示すために片側検定を行った。よって、有意水準を 0.05 としたとき、z 値が -1.64 を下回る勝率に関しては帰無仮説を棄却し、対立仮説を採用するため”白が黒より強い”ということが有意に言える。

6.3 勝率の評価

表 1 より、プレイアウト回数を 1 回のとき、先手黒が後手白より圧倒的に弱いことがわかる。徐々にプレイアウト回数を増やしていくと後手白の強さに近づいていくことがわかる。特に、プレイアウト回数 128 回及び

プレイアウト数	勝ち数	負け数	勝率 [%]	z 値
1	4	146	2.67	-11.59
2	18	131	12.08	-9.26
4	27	120	18.37	-7.67
8	36	113	24.16	-6.31
16	55	94	36.91	-3.2
32	65	84	43.62	-1.56
64	64	82	43.84	-1.49
128	81	69	54.00	0.98
256	74	72	50.68	0.17
512	67	81	45.27	-1.15

表 1 勝率

256 回では、先手黒のプレイアウト回数が後手白よりも少ないにもかかわらず勝ち越しているため、ある程度のプレイアウト回数があれば、それ以上増やしても強さに大きな影響はないといえる。

7 終わりに

モンテカルロ木探索におけるシミュレーション部分を並列化することができたため、逐次実行ではプレイアウト数を増やしたとき、実行時間が膨大となっていたところを大幅に短縮することができた。しかし、プレイアウト数を増やしたとしても、ある程度のプレイアウト数で強さに限界を迎えてしまう。よって、今後の課題としては、モンテカルロ木探索の別箇所の並列化を行い、より強くより速いプレイヤーを作る必要がある。モンテカルロ木探索において他に並列化できそうな箇所としては評価を行う関数を並列化することが考えられる。つまり、今回の実験においてある盤面ノードに対して 100 回の評価を行っていたのを、並列に実行し、高速化を図れるのではないかということである。しかし、評価関数の複数回実行を並列化するためには、それぞれのカーネルで評価を行って生成したゲーム木を 1 つに統一する必要がある。