

算術強度によるオフロード先振分けの検討

Study of appropriate offload destination based on software arithmetic intensity

山登庸次
Yoji Yamato

日本電信電話（株） ネットワークサービスシステム研究所
Network Service Systems Laboratories, NTT Corporation

1. はじめに

近年，IoT[1]-[4]等新たな領域で，GPU 等のヘテロデバイスの利用が増えているが，それらの活用には壁が高い。私は，記述コードを，配置先デバイスに合わせ，自動変換等し高性能に動作させる，環境適応ソフトウェアを提案しており，GPU での自動オフロード等実現してきた[5]。

現在，クラウド[6][7]等で GPU や FPGA，メニーコア CPU が混在の環境に対しては，全てのオフロード先性能を検証し最高速の移行先を選択しており，長大な時間がかかっている。混在環境にて，適切なオフロード先を事前にある程度振分できれば，検証時間削減が可能である。

本稿では，アプリケーションループ文の算術強度 (Arithmetic Intensity) を用いて，GPU 向き，メニーコア CPU 向きの推測が出来ないか検討，検証する。

2. 算術強度を用いた振分検討

私は，GPU 等への自動オフロードを提案してきたが，オフロードパターンを作成し，検証環境で性能測定し，進化計算等を用いて新たなパターンを作成し，徐々に高速化する事は共通である。これは，性能は，コードだけでなく，処理データ，デバイススペック等により変わるため予測困難であるためである。移行先が3つ混在の場合は，3つ全てに対してオフロードを試み，部分 A は GPU が最高速な際は，A は GPU にオフロードする方式である。

しかし，この方式は，3つに移行を試行するため長大な時間がかかる問題があった。特に，FPGA では実機で動作するまで配線等を含めたコンパイルに数時間以上かかる現状があり，数回試行でも大きく時間がかかってしまう。

もし，実測をしなくても，このアプリケーションは GPU 向き等が事前に分かれれば，GPU 試行を先に行い，その時点で十分な性能であれば，FPGA 検証は行わない等の対応ができる。そこで，本稿では，アプリケーションループ文の算術強度を用いて，まず，GPU とメニーコア CPU の2者に対して向き不向きを振分を検討する。

算術強度とは，計算程度を示す指標で，単位は FLOP/Byte で，ループ中データ当たりの浮動小数点計算回数を示しており，ROSE フレームワーク等の分析ツールによって，各 for 文毎の算術強度を取得することが出来る。

これまで混在環境に対する検証で，単純行列計算の Polybench 3mm は GPU の方が高速化でき，ブロック対角ソルバ計算の NAS.BT はメニーコア CPU の方が高速化できることが分かっている。NAS.BT はメニーコア CPU の方が速い理由として，データ量が多く，GPU-CPU のメモリ転送がネックになっていることが原因と推定している。

そこで，算術強度を分析し，例えば，最大算術強度が5以上は GPU 向き，1-5 はメニーコア CPU 向き，1以下は

通常 CPU 向き等ある程度の振分ができないか検証する。

3. GPU オフロード性能と算術強度相関

算術強度が高い程，GPU に向いており，GPU オフロードで大きく高速化出来るのではないかの仮説の検証のため，複数のアプリケーションの GPU オフロードを[5]の手法で行い，ループの最大算術強度との関係を表にまとめた。検証では，NVIDIA Quadro K5200 GPU を備えた，Dell new vostro 15 5000 マシンでオフロードを行った。図1は，各アプリケーションの GPU オフロード時の高速化率と for 文数，for 文中の最大算術強度を記載している。

検証の結果，例えば，NASA 開発の NAS.BT，NAS.FT，NAS.EP を見ても，算術強度が高い程，GPU で高速化できているというわけではなく，仮説は棄却された。算術強度は1ループに対しての値だが，GPU へのオフロードはループ回数も影響するので，gcov 等を用いたループ回数含めた検証や，GPU 転送量と処理データ量は厳密には異なるため転送量を正確に把握した検証等が必要と考える。

アプリ	NAS.BT	NAS.FT	NAS.EP	Correlation	Covariance	Laplace Eq	Himeno bench	DFT
for文数	120	81	12	12	10	5	13	10
最大算術強度	20.8	0.275	2.28	0.25	0.25	0.313	1.21	1.5
高速化率	1	5.36	19.3	25	26.4	13.6	4.82	5.1

図1: GPU オフロード高速化率と算術強度

4. まとめ

ループ文算術強度を用いて，GPU，メニーコア CPU 適合性推測の妥当性を検証した。算術強度だけでは振分が困難であり，別指標と組み合わせた振分を今後検討する。

参考文献

- [1] Y. Yamato, et al., "Security Camera Movie and ERP Data Matching System to Prevent Theft," IEEE CCNC 2017, pp.1021-1022, Jan. 2017.
- [2] Y. Yamato, "Experiments of posture estimation on vehicles using wearable acceleration sensors," IEEE BigDataSecurity 2017, pp.14-17, May 2017.
- [3] Y. Yamato, et al., "Analyzing Machine Noise for Real Time Maintenance," ICGIP 2016, Oct. 2016.
- [4] Y. Yamato, "Proposal of Vital Data Analysis Platform using Wearable Sensor," ICIAE 2017, Mar. 2017.
- [5] Y. Yamato, "Study of parallel processing area extraction and data transfer number reduction for automatic GPU offloading of IoT applications," Journal of Intelligent Information Systems, Springer, Vol.54, No.3, pp.567-584, 2020.
- [6] Y. Yamato, "Automatic system test technology of virtual machine software patch on IaaS cloud," IEEEJ Transactions on Electrical and Electronic Engineering, Vol.10, pp.165-167, 2015.
- [7] Y. Yamato, "Proposal of Optimum Application Deployment Technology for Heterogeneous IaaS Cloud," WCSE 2016, pp.34-37, June 2016.