

FPGA を用いた負荷分散ストレージシステムの実装と評価 A Study of FPGA-based load balancing toward a Parallel Storage System

梅津 直弥[†] 山口 佳樹[‡]
Naoya Umezu[†] Yoshiki Yamaguchi[‡]

1. はじめに

新型コロナウイルス感染症 (COVID-19) の影響は、従来の勤務・学習形態を初め、我々の生活様式を一変させた。これを可能にしたのは、物理的な距離感を情報通信技術の積極利用により隠蔽し、円滑なコミュニケーションを使用に耐えるレベルで提供できたことが一因と言える。そして今後は、この情報通信技術をキーワードに、新しい価値観や新ビジネスが提案されると考えられる。例えば、自動運転・医療分野 (8K 画像などの高解像解析)、防災 (多様かつ膨大なセンサデータ)、企業・自治体 (文書の電子化) などに加えポストコロナを考えると、膨大なデータをリアルタイムに保管する技術 (ストリームデータ保存) が一つのキーテクノロジーになると思われる。

ストリームという観点において、アーキテクチャ上の特徴を活かしやすい FPGA を用いた演算システムは既に数多く提案されている。例えば、マイクロソフト社の Catapult project [1], Brainwave project [2], Azure [3] などは産業界で代表するシステムと言える。また、Amazon EC2 F1 インスタンスなど、クラウド大手各社も FPGA の積極的な活用により性能向上を考えている [4]。また、学界においても、EZH Zurich, UCLA, UIUC, NUS の国際的研究クラスター [12] や我が国でも ACRi [13] による産学連携によるアダプティブコンピューティング研究が進められている。そして、これにストレージを加えたストリームストレージ技術の飛躍的な性能向上は議論の活発化が予想される。例えば、2019年に日立からはVSP5000シリーズが発表されており、このストレージシステムのエンジンには FPGA が利用されている。

そこで本プロジェクトでは、FPGA に大量のストレージを密にかつ並列に接続し、大容量かつ高速・低遅延なストレージシステムを構築することを提案する。加えて、FPGA の特徴を活かし、ロードバランシング、データ整理、信頼性の向上など、また必要に応じてデータ演算についても実時間内に演算するよう組み込む予定である。本論文では、データストレージ部に焦点を絞り、その動作および性能について報告する。

以下、本論文は以下のように構成される。第 2 章では本研究の背景と意義について述べ、第 3 章で提案するシステムについて説明する。続く第 4 章で実験環境を、第 5 章で実験結果について議論し、第 6 章で本論文をまとめる。

[†] 筑波大学大学院システム情報工学研究科,
Graduate School of Systems and Information Engineering,
University of Tsukuba

[‡] 筑波大学システム情報系,
Faculty of Engineering, Information and Systems,
University of Tsukuba

2. 研究背景

近年、クラウド大手各社では FPGA を活用したデータセンタ構築が選択肢の一つとなりつつあり、Microsoft 社の Azure[3], Amazon 社の Amazon EC2 F1 インスタンス [4], IBM 社の SuperVessel[5] などが挙げられる。

高性能計算分野においても、筑波大学計算科学研究センターで運用していたスパコン「HA-PACS (PACS 8)」では、密結合並列演算加速機構(TCA:Tightly Coupled Architecture)として、システム内通信を加速させる目的として FPGA が利用されていた。

CPU の性能向上や FPGA の導入の寄与もあり、演算性能やネットワーク性能が年々向上している一方で、データストレージ部に関しては大きな性能向上が見られず、システム全体の処理性能の向上において大きな課題となっている。

そこで、計算機システムの処理性能のブレイクスルーを実現することを目標として、システムに対する大規模ストレージノードではなく、各演算ノードからメモリ帯域と同等以上のアクセス帯域を保証する中～大規模ストレージシステムの構築を提案する。また、高速かつ低遅延であるストレージシステムの実現にむけ、ストレージ制御のコア部分には FPGA を利用する[15]。

3. 実験方法

3.1 目標システム

FPGA の特徴として、SATA を始めとする高速シリアル通信に使用可能な I/O であるトランシーバーを内蔵していることが挙げられる。ストレージの同時接続可能台数はこのトランシーバーの数によって左右される。また、ストレージのアクセス速度はこのトランシーバーのサポートする帯域によって決まる。

FPGA におけるトランシーバーの性能は年々向上しており、Xilinx 社の大規模 FPGA である VIRTEX Ultrascale+ XCVU13P においては 32.75Gb/s で通信可能な GTY トランシーバーが 128 チャンネル搭載されている。

一方で、シリアル通信規格によっては、デバイス 1 台の接続に対して複数の高速シリアル IO を要求する規格がある。複数のシリアル I/O を用いる規格の方がデバイスあたりの転送性能は高くなる傾向にあるが、その分多く I/O を一台が専有することになり、同時接続可能台数は少なくなる。また、FPGA 内に制御コアを実装することも考えると、期待される性能を引き出すことは、ハードウェア・ソフトウェア両面の制約から難しいことが推測される。

本研究では、大容量の SSD を可能な限り多数同時に FPGA に接続し、補助記憶装置として重要視される要因の一つであるストレージ容量についての高い要求を満たしつつ、比較的シンプルな規格であることから高いカスタマイズ性のもとシステム全体の性能を最大限に引き出すことが

期待できるという観点から、SATA インターフェースによる SSD の接続を選択した。

将来的には、前述の 128 チャンネルの高速シリアル I/O 全てに SATA SSD を同時接続し、大容量かつ高速、低遅延なストレージシステムの構築を目標とする。

3.2 システム実装

入手可能かつユーザが FPGA 部を書き換え可能なシステムを考えたとき、SATA が 8 ポートあるシステムが基本的に最大である。しかし、8 ポートにおける理論性能の合計も高々 6GB/s 程度であり、目標とする性能との乖離が激しい。そこで、著者らは SATA コネクタが 16 ポートまで、SDI コネクタが 4 ポートまで拡張可能なボードを製造した。最終目標に向けてはこのサブ基板 8 枚（合計 128 台の SATA SSD）が 1 枚のメイン基板に搭載された 1 個の FPGA で制御される。本報告では、このサブ基板の評価実験について報告する。



図 1 今回実験に用いたサブ基板

サブ基板に搭載された FPGA は、Xilinx 社製 Kintex UltraScale FPGA の XCKU060-2FFVA1156I である。SSD は、複数種類の確認を行った知見[4]より、性能のバラツキが小さく信頼性の高い SAMSUNG 社製 860EVO を選定した。

SSD のボードへの接続は、図 2 に示される専用の拡張モジュールを介して行われる。1 台の拡張モジュールに対し、SATA SSD が同時に 8 台まで接続できる。



図 2 SSD 拡張モジュールに SAMSUNG EVO 860 SSD を 8 台同時に接続した図

本実験において用いたシステムは、図 3 に示されるように、中央のヒートシンクが付いた FPGA 基板と、上側の SATA デバイス用の I/F 基板（赤い基板）から構成されている。SATA I/F 用基板は、FPGA 用基板と変換コネクタを介して接続されている。

今回は、本ボードの最大接続可能台数である 16 台の SSD を接続し、そのうち制御が安定しているチャンネルを用いて実験を行った。



図 3 SATA サブ基板 1 台に対し、最大接続可能台数である 16 台の SSD を接続した図

また、本ボードは、別途カメラモジュールの拡張ボードとフラットケーブルを介して接続することができ、4 系統の TX と RX および同期用の計 9 ポートと FPGA が直接接続される。

カメラモジュールはそれぞれが 12G-SDI により FPGA と接続される。本構成における SSD に対する総理論転送性能は 12 [GB/s] であり、仕様上は 4K 画像を 4 系統まで扱うことが可能である。本ボードの用途の一つとして、将来的には 4K カメラからの高解像度映像をリアルタイムにストレージに保存することが期待される。

3.3 SATA (Serial AT Attachment)

SATA とは、2002 年 2 月に発足した Serial ATA ワーキンググループが発表したインターフェース規格であり、1.5Gbps (Revision 1.x)、3Gbps (Revision 2.x)、6Gbps (Revision 3.x) のデータ転送速度を持つ 3 つの中心規格が存在する。

表 1 : Serial ATA 規格

	SATA Generation		
	Revision 1.x	Revision 2.x	Revision 3.x
Transfer Rate [MB/s]	187.5	375	750
Theoretical Bandwidth [MB/s]	150	300	600

SATA では、8bit のデータを 10bit のデータにエンコードする、8B/10B 変換を行った後に送信される。そのため、実効転送速度は物理転送速度の 80% の値となるが、データの偏りによって、データの変化点がわからなくなったり、電荷のバランスが崩れたりするといった問題を回避できる。また、信号は差動信号として送受信されるほか、信号が Hi-Z となる OOB(Out of Band)信号もサポートする必要がある。

これらの信号を送受信するためには高速シリアル通信をサポートするトランシーバ (HSSIO : High Speed Serial IO) が必要である。本稿で用いた FPGA (XCKU060) は GTH トランシーバが搭載されており、これを利用することで SATA SSD の通信を実現できる。本稿のシステムはカスタムで開発が必要な部分が多く、市販の IP コアを利用することが難しい。そこで、SATA SSD 制御用回路はスクラッチに近い形で書かれており、GTH トランシーバ周囲の高速シリアル IO に関連する部分については Xilinx 社による IP コアを用いて行った。

3.4 FPGA と演算回路部

FPGA は、任意の回路に再構成可能な回路であり、HDL (Hardware Description Language) を用いて記述する。特定の計算に特化した回路を記述することで、高速並列処理を得意とすることから、特に近年では高性能計算分野で活躍している。

本研究でも、低遅延高速ストレージアクセスの要として、各種データ通信の制御に用いており、ハードウェアとして高速シリアル通信を可能にするギガビットトランシーバを搭載している。

3.5 回路実装

本実験におけるシステムのブロック図は図 4 に示す。

最大 16 の SATA CORE を管理するモジュールが存在し、速度計測や、その後の負荷分散アクセスについての制御をスタンドアロンで行う。

各 SATA CORE は、対応する SATA デバイスとの通信を制御し、物理層の GTH トランシーバを介して SSD にアクセスする。SSD のリセットなどの初期化処理も、SATA 制御コアが行う。

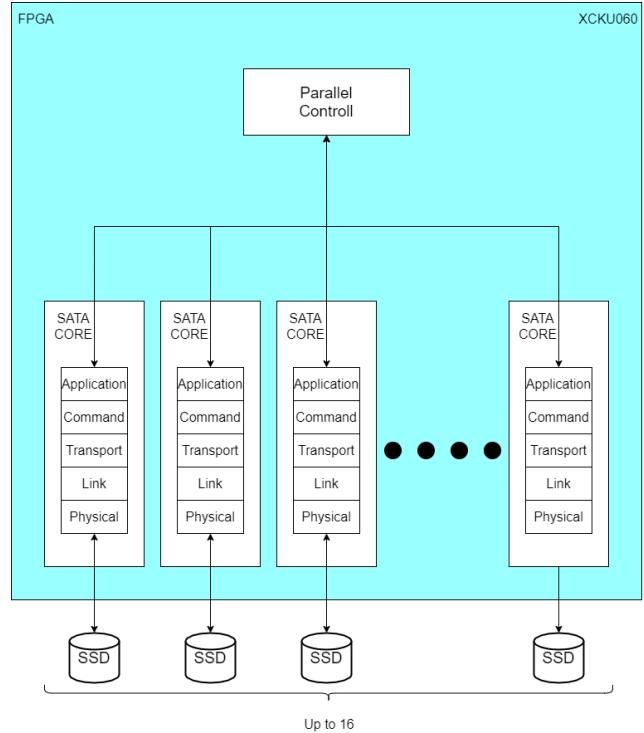


図 4 SATA の各層と機能の全体図

3.6 負荷分散アクセス

本システムは、データアクセスに先立って、各デバイスへの転送速度の計測を行い、その速度の比に基づき、各コアに配分するデータ量を調節し、全体の転送速度を SATA の規格の理論値に近づける。図 5 に、全体の処理のフローチャートを示す。

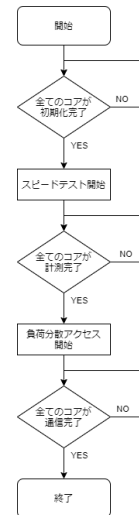


図 5 負荷分散計算フローチャート

4. 実装結果

4.1 実験及び実験結果

本実験では、対象となる SATA デバイスに対し、65536 セクタの DMA Write/Read を単位として、デバイスの先頭アドレスからシーケンシャルにアクセスを行う。データア

アクセス量は、この DMA Write/Read を繰り返し行うことで調節する。その際、繰り返しの 2 回目以降の DMA Write/Read では、ループ回数に応じ、全体を通してシーケンシャルなアクセスになるように、デバイスにアクセスを始めるアドレスを変更する。

本実験におけるシステムの状態の観測は、Xilinx の IP である VIO を用いて、USB-JTAG 経由で開発環境である Vivado 上に表示する。

その際、一部内部波形は Xilinx の IP である ILA を用いて、USB-JTAG 経由で開発環境である Vivado 上に表示する。

4.2 リソース使用量

今回の実験で実装した回路が使用した FPGA のリソースの割合を表 2 に示す。

表 2 FPGA リソース使用量

Resource	Utilization	Available	Utilization %
LUT	61846	331680	18.65
LUTRAM	144	146880	0.10
FF	46742	663360	7.05
BRAM (in RAMB18 equivalent)	104	1080	9.63
IO	3	520	0.58
GT	16	28	57.14
BUFG	57	624	9.13
MMCM	1	12	8.33

16 コア分のデータ量の配分の計算について、16 ビット整数の剰余を用いたため、LUT の使用量が多くなっている。

また、内部状態の観測のために、レジスタを多く用いて特定の状況での内部信号の保存を行ったため、FF の使用量も多くなっている。

5. 性能評価

5.1 SATA CORE の性能評価

負荷分散システムの評価に先立って、SAMUSUNG 860EVO 500GB 1 台に対する DMA Write によるスピードテストを行い、実装した SATA CORE の性能を検証した結果を図 6 に示す。ここでは標準偏差のエラーバーを用いた。以降の結果も同様である。

テストデータが 8GiB の時に最も DMA Write のスピードが速く、506.0[MB/s]であった。これは、SATA Rev 3.x における転送速度の 83.7%である。

同様に、DMA Read によるスピードテストも行った。その結果を図 7 に示す。

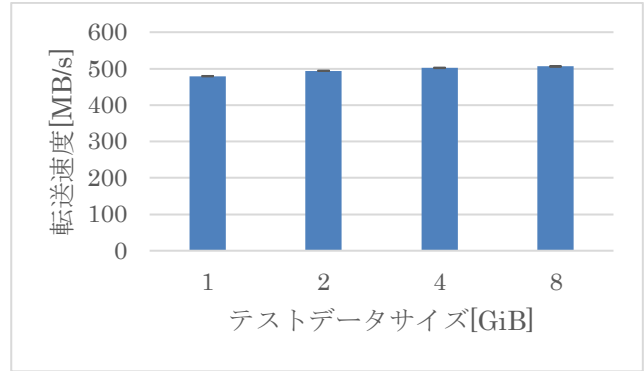


図 6 SAMSUNG 製 SSD 860 EVO 250GB に対する書き込みスピードテスト結果

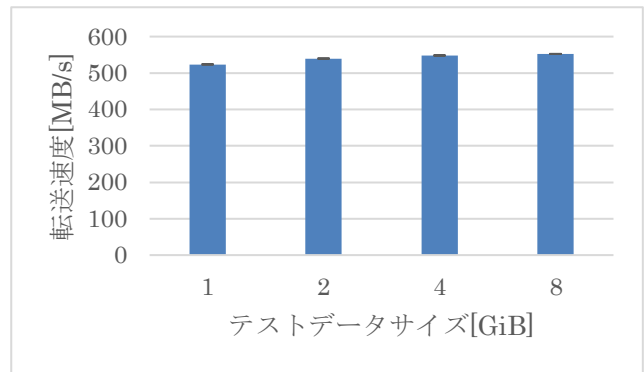


図 7 SAMSUNG 製 SSD 860 EVO 250GB に対する読み出しスピードテスト結果

テストデータが 8GiB の時に最も DMA Read のスピードが速く、552.4[MB/s]であった。これは、SATA Rev 3.x における転送速度の 92.0%である。

SAMSUNG 860 EVO のシーケンシャルアクセス性能の公称値は読み出しが 550MB/s、書き込みが 520MB/s であることから、十分な性能を発揮できていると言える。

5.2 負荷分散システムの評価

高速なストレージである SSD に対し、低速なストレージとして、東芝製 HDD MQ01ABF050 を用いた。

システムの動作として、スピードテストにより低速なストレージを検知した後、本番のデータアクセスに際して、スピードテストにより得られた転送速度の比に基づいてデータ配分を行い、全体の転送速度を各デバイス単独でのアクセス時の転送速度の総和に近づけることが期待される。

スピードテスト時には、各デバイスに同じずつの量のデータアクセスを試みる。スピードテストが完了した後に、全体としてのデータアクセス量は保ちつつ、計測結果の転送速度に基づきデバイスごとにデータアクセス量を割り振る。

以下、スピードテスト時のテストデータサイズを 0.5, 1, 2GiB と変化させて実験を行った。なお、本番のデータアクセス量は

スピードテスト時のテストデータサイズ * デバイス台数となる。

5.2.1 HDD の性能評価

今回の実験に用いた HDD の読み出し速度の測定結果を図 8 に示す。

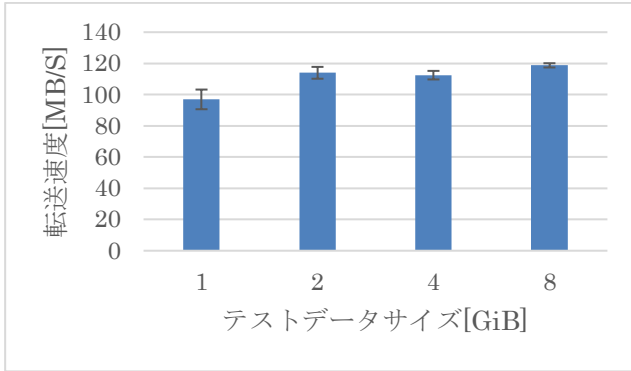


図 8 東芝製 HDD MQ01ABF050 500GB に対する読み出しスピードテスト結果

テストデータが 8GiB の時に最も DMA Read のスピードが速く、117.8 [MB/s]であった。これは、SATA Rev 3.xにおける転送速度の 19.6%である。

読み出し速度の測定結果にばらつきが見られたが、スピードテスト開始時のディスクの回転の状態によっては、読み出しが実際に始まるまでのタイムラグが生じる場合があることに起因すると考えられる。また、連続アクセス時間が長いほど、その影響が相対的に小さくなっていると考えられる。

5.2.2 負荷分散システムの実効転送速度計測

以下、実効転送速度に差のあるストレージが複数導入されている環境を想定し、複数の SSD に対して、低速な HDD を 1 台混ぜて負荷分散アクセスを実行し、実効転送速度を計測した。

なお、今回の実験において、計測は読み出しにて行った。

5.2.3 SSD 1 台 + HDD 1 台

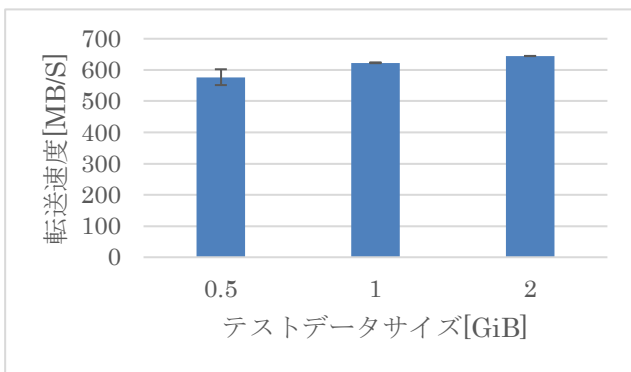


図 9 SSD 1 台 + HDD 1 台に対する負荷分散読み出しスピードテストの結果

テストデータが 2GiB の時に最も負荷分散読み出しのスピードが速く、644.7[MB/s]であった。これは、SATA Rev 3.xにおける転送速度 2 倍の 1200[MB/s]の 53.7%である。

各デバイスの 2GiB ずつの読み出しが全て終わるまでを実効転送時間とすると、その時の実効転送速度は 212.0[MB/s]であったことから、データを当分割してアクセスする方式に対して速度が高くなっている。

また、スピードテスト時に算出した実効転送速度の理論値は 642.3[MB/s]であることから、負荷分散は理想的に行われていると言える。

5.2.4 SSD 3 台 + HDD 1 台

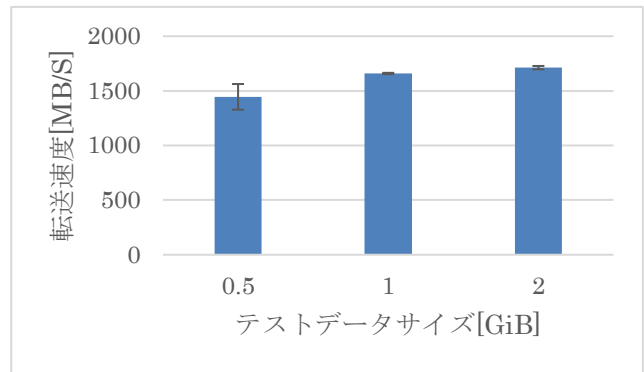


図 10 SSD 3 台 + HDD 1 台に対する負荷分散読み出しスピードテストの結果

テストデータが 2GiB の時に最も負荷分散読み出しのスピードが速く、1712.3 [MB/s]であった。これは、SATA Rev 3.xにおける転送速度 4 倍の 2400[MB/s]の 71.3%である。

各デバイスの 2GiB ずつの読み出しが全て終わるまでを実効転送時間とすると、その時の実効転送速度は 430.2 [MB/s]であったことから、データを当分割してアクセスする方式に対して速度が高くなっている。

また、スピードテスト時に算出した実効転送速度の理論値は 1722.5 [MB/s]であることから、負荷分散は理想的に行われていると言える。

5.2.5 SSD 7 台 + HDD 1 台

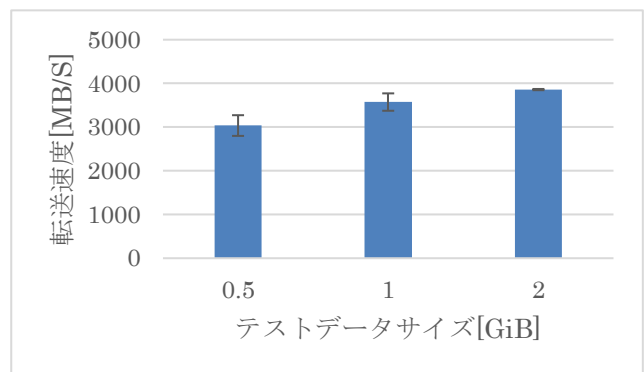


図 11 SSD 7 台 + HDD 1 台に対する負荷分散読み出しスピードテストの結果

テストデータが 2GiB の時に最も負荷分散読み出しのスピードが速く、3859.0 [MB/s]であった。これは、SATA Rev 3.xにおける転送速度 8 倍の 4800[MB/s]の 80.3%である。

各デバイスの 2GiB ずつの読み出しが全て終わるまでを実効転送時間とすると、その時の実効転送速度は 871.3[MB/s]であったことから、データを当分割してアクセスする方式に対して速度が高くなっている。

また、スピードテスト時に算出した実効転送速度の理論値は 3913.3[MB/s]であることから、負荷分散は理想的に行われていると言える。

5.2.6 評価結果まとめ

負荷分散アクセスによる実効転送速度の評価結果を表 3 に示す。表中の%表記は、実験に用いた SATA デバイスの台数分の SATA 規格理論上の転送速度の総和に対する、実験結果の転送速度の割合を示す。

表 3 評価結果

構成	最大転送性能[MB/s]	平均転送性能[MB/s]	最小転送性能[MB/s]
HDD*1 SSD*1	645.0 [53.8%]	644.7 [53.7%]	644.3 [53.7%]
HDD*1 SSD*3	1720.6 [71.7%]	1712.3 [71.3%]	1681.3 [70.1%]
HDD*1 SSD*7	3867.8 [80.6%]	3859.0 [80.4%]	3831.2 [79.8%]

アクセス速度の遅いデバイスの割合が小さいほど、実効転送速度が理論値に近づいている一方で、いずれの場合もデータを当分割してアクセスした場合に対して実効転送速度の向上が見られる。

6. 結論

本実験では、転送速度に差のある複数のデバイスについて、スピードテスト結果に基づいて FPGA 各々のデータアクセス量を最適に配分することで、全体の転送速度を各デバイスの理論値の総和に近づけ、ひいては SATA 規格の理論値に近づけることができた。

テストデータサイズが小さいケースでは、テストデータサイズが大きいケースに対して性能の低下が見られる。これは、テストデータサイズが大きいほど SSD の実効転送速度が高くなるという計測結果の他に、FPGA 内部で整数による剰余の演算が行われる結果、途中計算結果の値の丸め込みなどによりアクセス速度の遅い HDD へのデータ配分が相対的に多くなってしまいうケースの一部由来すると考えられる。

今回は 16 台の SATA デバイスを搭載可能なシステムを用いて実験を行い、16 台の SAMSUNG EVO 860 SSD での同時読み出しにも成功したが、多数のデバイスを同時に接続したときの物理層の動作に不安定性が見られるため、今後改善を図る。

本研究では、FPGA を用いて SATA 経由で SSD を多数接続し、更に、FPGA が事前にスピードテストによって計測した速度を元に、データアクセス量を適切に分配することで、デバイスごとの実効転送速度の差がシステム全体の転送速度に悪影響を及ぼすことを回避することができる事が分かった。

一方で、ストレージシステムとしての観点からは、RAID システムに代表されるような冗長性や、各種ファイルシステムの実装が求められる。

また、将来的には最大 128 台の SATA SSD を同時に制御することを検討している。今回使用した Xilinx Ultrascale XCKU060 の持つトランシーバーは 28 本が使用可能であったが、より大規模 FPGA を用いて同時接続数を増やすことが考えられる。

今後は、システム全体の安定性や同時接続台数の向上について図りつつ、冗長性やファイルシステムについて検討していく。

また、今回目標とした大容量かつ高速、低遅延のストレージシステムは、近年需要増大傾向にある、運転支援技術における危険な状況のシミュレーションのためのドライビングシミュレータにおいても強く求められている。遠くに映る小さな物体も精細に描画することのできる、4K、8K 画質のシミュレータが、今後主流なシミュレーションシス

テムに導入されており、高度な安全運転技術の進歩に寄与することが期待されている。一方で、高解像度の画像を高フレームレートで取得するという事は、画像を保存するためのストレージシステムに求められる要件も厳しいものとなるため、本研究の目標とするストレージシステムの特性を要求する分野の一つと言える。

本研究におけるストレージシステムの応用先の一つとして、高解像度映像のリアルタイムでの読み書きも視野に入れて、今後の改善を図る。

謝辞

本研究の一部は、科学研究費補助金(JP17H01707, JP18H03246, JP19H00806)および TIA 連携プログラム探索事業「かけはし」(2019 年度)の助成を受けたものである。また、Xilinx 社より「Xilinx University Program」を通じて開発ソフトウェアの支援を受けており、ここに謝意を表す。

参考文献

- [1] Adrian M. Caulfield, Eric S. Chung, Andrew Putnam, et al., "A CloudScale Acceleration Architecture", Annual IEEE/ACM Symposium on Microarchitecture, pp.1-13, Oct. 2016.
- [2] Eric S. Chung and Jeremy Fowers, "Accelerating Persistent Neural Networks at Detacenter Scale", Hot Chips: A Symposium on High Performance Chips, August 2017.
- [3] Mark Russinovich, "Inside the Microsoft FPGA-based configurable cloud", Build2017, May 2017.
- [4] Amazon Web Services, Inc., "Amazon EC2 F1 インスタンス", (オンライン: <https://aws.amazon.com/jp/ec2/instance-types/f1/>) (引用日:2020年6月19日)
- [5] IBM, "Xilinx and IBM to Enable FPGA-Based Acceleration within SuperVessel OpenPOWER Development Cloud", <https://www.xilinx.com/news/press/2016/xilinx-and-ibm-to-enablefpga-based-acceleration-within-supervessel-openpowerdevelopment-cloud.html>, Press Release, April 2016.
- [6] Steve Crowe, "Waymo autonomous vehicles leave Apple in the dust", ROBOTICS NEWS, ROBOTREPORT, 15 Feb., 2019.
- [7] 菅原 博英, シリアル ATA の基礎と FPGA への実装, 東京, CQ 出版, 2010, 272p
- [8] XILINX, UltraScale アーキテクチャ GTH トランシーバーユーザガイド(), 2018, 492p,
- [9] Design Gateway Co., Ltd : SATA IP Transport & Link Layer Core 2016, 11p
- [10] 紀野國祐太, 山口佳樹: FPGA を用いたストレージコントローラの実装と評価, 情報処理学会第 81 回全国大会, 2L-09, 2019年3月. 2p.
- [11] 段然, 梅津直弥, 山口佳樹: 実車映像コンテンツ作成のための FPGA システム, 自動車技術会 2019 年秋季大会, 2019年10月. 6p.
- [12] Xilinx, Xilinx Teams with Leading Universities Around the World to Establish Adaptive Compute Research Clusters, May 05, 2020. [Online] <https://www.xilinx.com/news/press/2020/xilinx-teams-with-leading-universities-around-the-world-to-establish-adaptive-compute-research-clusters.html>
- [13] ACRI, 「アダプティブコンピューティング研究推進体-ACRI」を設立—日本初, 産学連携で FPGA 検証環境と学習機会を無償で提供—, プレスリリース, <https://www.acri.c.titech.ac.jp/wp/>, 2020年4月3日.
- [14] Hitachi Vantara, "Hitachi Accelerated Fabric", White Paper, WA-592-A, September 2019.
- [15] Kodama, Yuetsu and Hanawa, Toshihiro and Boku, Taisuke and Sato, Mitsuhsa, "PEACH2: An FPGA-Based PCIe Network Device for Tightly Coupled Accelerators", SIGARCH Comput. Archit. News, Vol.42, No.4 (2014)