

# NVMe-IP リファレンス・デザイン説明書

Rev2.1J 2017/07/31

本ドキュメントは NVMe -IP デモのリファレンス・デザインに関する説明書となります。リファレンス・デザインを実装した NVMe-IP の具体的なデモ手順については以下のドキュメントを参照してください。

文書名: NVMe-IP デモ手順書

ファイル名: dg\_nvmeip\_instruction\_jp.pdf

## 1. NVMe 規格と IP コアについて

NVM Express (NVMe)は PCI Express 接続の SSD(ソリッド・ステート・ドライブ)をアクセスするホスト・コントローラのインターフェイスを定義した規格のひとつです。NVMe はコマンド発行と完了を、コマンド発行/完了サイクルのわずかな遅延で最適処理されています。また、NVMe では単一キューでも最大 65,536 コマンドを同時に並列実行できます。このためシーケンシャル・ランダムのどちらのアクセスもパフォーマンスが改善されています。

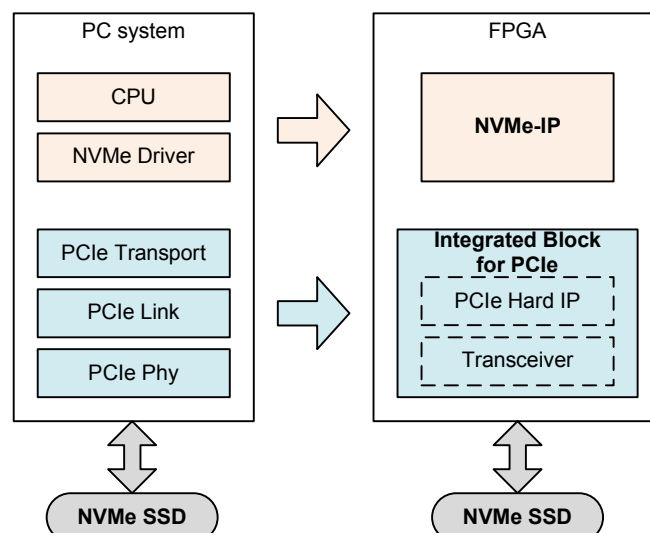
PCIe SSD のマーケットでは2つの規格を選択できます、一つは AHCI でもうひとつが NVMe です。AHCI は SATA ハード・ディスク向けの古い規格であるのに対し、NVMe は SSD のような不揮発メモリ装置に最適化された規格です。両者の規格についてより詳細に比較したドキュメントは以下で参照できます。

[https://sata-io.org/system/files/member-downloads/NVMe%20and%20AHCI\\_%20\\_long\\_.pdf](https://sata-io.org/system/files/member-downloads/NVMe%20and%20AHCI_%20_long_.pdf)

また、NVMe ストレージ・デバイスのリストは以下から参照可能です。

<http://www.nvmexpress.org/products/>

一般的にはユーザは図 1 に示すように NVMe SSD をアクセスするためには NVMe ドライバをインストールする必要があります。NVMe SSD のコネクタ形状としては M.2 コネクタのような PCIe タイプとなります。NVMe-IP コアは NVMe ドライバ機能と CPU で実行するタスク処理機能を純ハードワイヤード・ロジックで実装しています。このため、NVMe-IP コアを FPGA に内蔵することで CPU なしで NVMe SSD にアクセスすることが可能です。



(上図左は一般的な PC でのレイヤ構成、右は FPGA システムによるレイヤ構成)

図 1: NVMe のプロトコル層

## 2. 概要

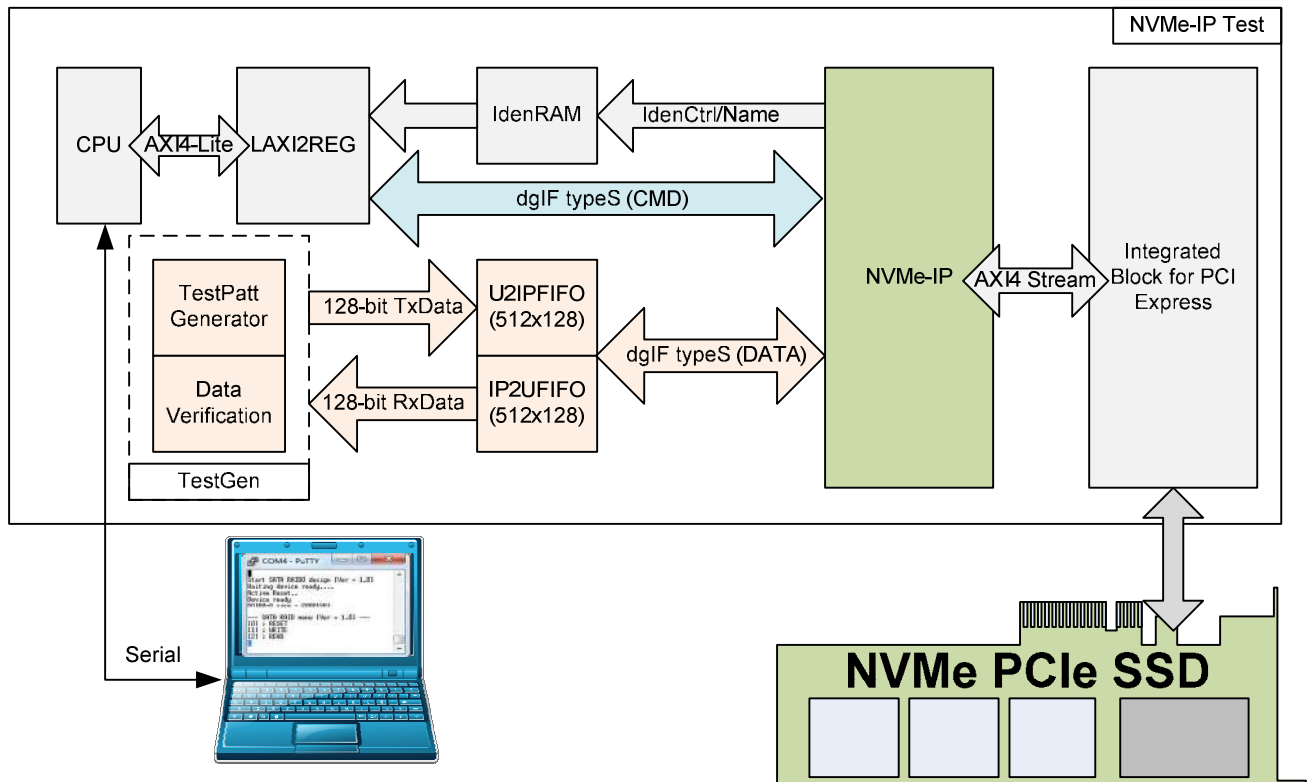


図 2: NVMe-IP デモ(リファレンス・デザイン)のブロック図

このリファレンス・デザインでは NVMe-IP コアにシンプルなロジックを加えて NVMe PCIe SSD に対して高速でライト・リードを実行するシステムを構築しています。デザイン内では CPU が使われていますがそれは主にシリアル・コンソールのユーザ・インターフェイスを実装するためのもので、SSD へのアクセスに CPU は必須ではありません。

このシンプルなテスト・システムにてユーザはリード・ライトのアクセス先セクタ・アドレス、転送長、コマンドをキーボードから入力し、その入力パラメータは NVMe-IP コアのインターフェイス信号に変換されます。指示動作が完了すると、CPU は計測したコマンド所要時間と総転送データ量から SSD のライト/リード・パフォーマンスを計算しコンソール上に表示します。

CPU バスと接続するため LAXI2REG モジュールによりアドレスがデコードされ CPU からデータが NVMe-IP コアの制御/ステータス信号と接続されます。一方 NVMe-IP コアのデータ・ポートは外部 FIFO と接続し、データはコアに内蔵した 256K バイトのデータ・バッファへ転送されます。これら制御およびデータのユーザ・インターフェイスを非常にシンプルで使いやすい仕様で定義したものが dgIF で、dgIF typeS は DesignGateway 社のストレージ系 IP コア間で共通のユーザ・インターフェイスです。

TestGen モジュールはテスト・パターンの発生器を含み、生成したテスト・パターンは ライト・テストでは書き込みデータとして SSD へと転送され、リード・テストではベリファイの期待値データとして使われます。IdenCtrl/IdenName データは初期化時に IdenRAM へ転送されるので CPU は SSD の型番などの Identify 情報にアクセスできます。

このリファレンス・デザインでの NVMe-IP クロック周波数は、PCIe Gen3 の場合 275MHz で PCIe Gen2 の場合 200MHz です。NVMe-IP クロック周波数は、PCIe 統合ブロックの周波数と同じかそれ以上とする必要があります。(PCIe Gen2 の場合 125MHz で PCIe Gen3 の場合 250MHz)

NVMe-IP コアのデータシートのダウンロードや評価用の bit ファイルのリクエストは以下の弊社 Web サイトから可能です。また、本リファレンス・デザインの実機評価でのパフォーマンスは NVMe PCIe SSD に依存します。

[NVMe-IP コア(Xilinx 版)紹介ページ]

[http://www.dgway.com/NVMe-IP\\_X.html](http://www.dgway.com/NVMe-IP_X.html)

### 3. CPU

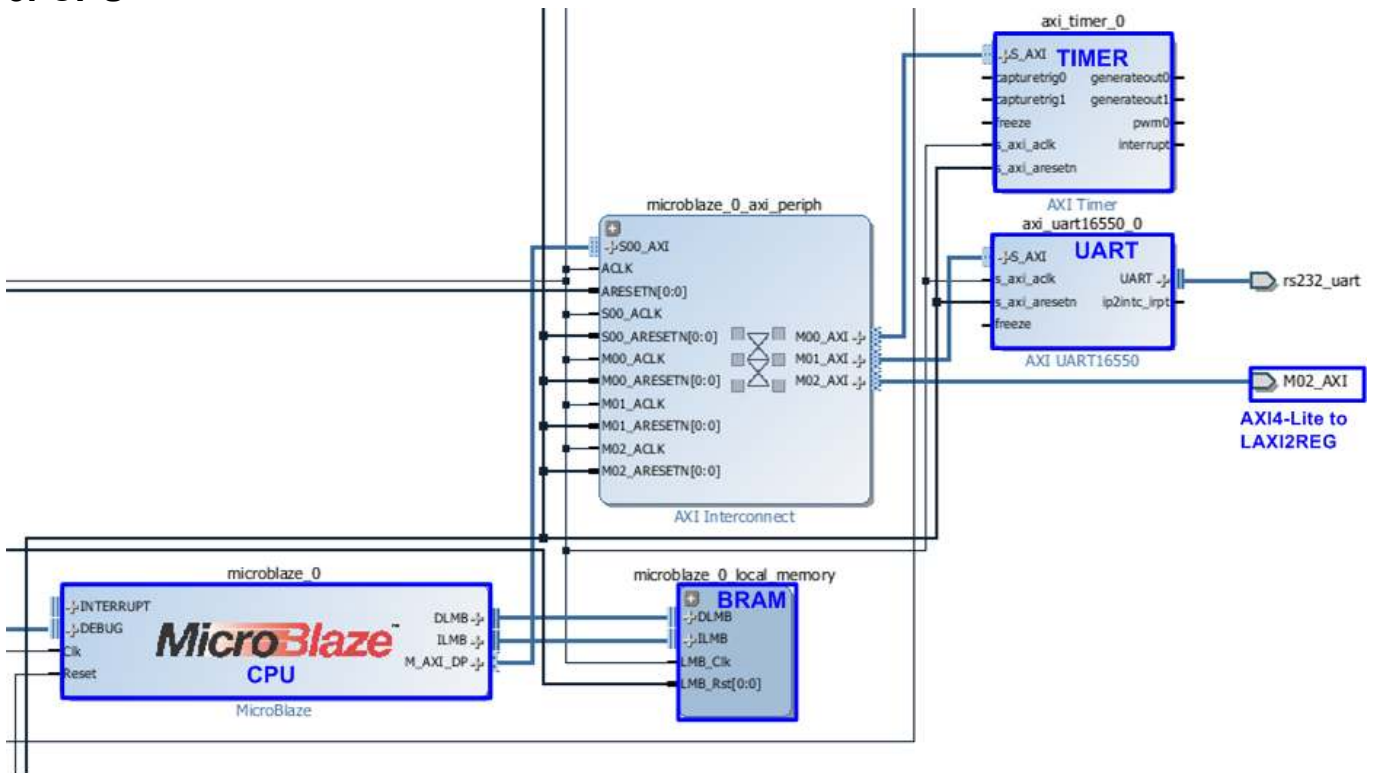


図 3: リファレンス・デザインの CPU システム

リファレンス・デザインにおいて CPU の周辺回路にはユーザ・インターフェイスとして UART、パフォーマンス計測用のタイマー、CPU ファームウェア格納用の BRAM があります。AXI インターコネクトは MicroBlaze や周辺回路の接続として使われますが、NVMe-IP コアを制御/監視するための 32 ビット AXI4-Lite バスの外部接続にも使われます。AXI4-Lite の詳細なメモリ・マップを下表 1 に示します。

表 1: レジスタ・マップ

アドレス	レジスタ名	説明
Rd/Wr	(“nvmeiptest.c”内のラベル名)	
BA+0x00	ユーザ・アドレス(下位)レジスタ (USRADRL_REG)	[31:0]: NVMe-IP のアクセス先開始アドレス下位 32 ビット UserAddr[31:0]
BA+0x04	ユーザ・アドレス(上位)レジスタ (USRADRH_REG)	[15:0]: NVMe-IP のアクセス先開始アドレス上位 16 ビット UserAddr[47:32]
BA+0x08	ユーザ転送長(下位)レジスタ (USRLENL_REG)	[31:0]: NVMe-IP の転送セクタ数下位 32 ビット UserLen[31:0]
BA+0x0C	ユーザ転送長(上位)レジスタ (USRLENH_REG)	[15:0]: NVMe-IP の転送セクタ数上位 16 ビット UserLen[47:32]
BA+0x10	ユーザ・コマンド・レジスタ (USRCMD_REG)	[1:0]: NVMe-IP のユーザ・コマンド UserCmd “00”-Identify device, “10”-Write SSD, “11”-Read SSD 本レジスタが書き込まれると NVMe-IP に対して新たなコマンド実行の要求を発生します。
BA+0x14	テスト・パターン・レジスタ (PATTSEL_REG)	[2:0]: テスト・パターン選択 “000”-インクリメンタル, “001”-デクリメンタル, “010”-オール 0, “011”-オール 1, “100”-32 ビット LFSR

アドレス Rd/Wr	レジスタ名 (“nvmeipstest.c”内のラベル名)	Description
BA+0x100 Rd	ユーザ・ステータス・レジスタ (USRSTS_REG)	[0]: NVMe-IP のビジー・フラグ ('0': アイドル, '1': ビジー) [1]: NVMe-IP からのエラー出力 ('0': 通常, '1': エラー) [2]: データ・ベリファイ・エラー ('0': 通常, '1': ベリファイ・エラー発生) [4:3] NVMe-IP からの PCIe リンク速度 (“00”: 未リンク状態, “01”: PCIe Gen1, “10”: PCIe Gen2, “11”: PCIe Gen3)
BA+0x104 Rd	総ドライブ容量(下位)レジスタ (LBASIZEL_REG)	[31:0]: NVMe-IP で報告される総ドライブ容量(単位:セクタ)下位 32ビット LBASize[31:0]
BA+0x108 Rd	総ドライブ容量(上位)レジスタ (LBASIZEH_REG)	[15:0]: NVMe-IP で報告される総ドライブ容量(単位:セクタ)上位 16ビット LBASize[47:32]
BA+0x10C Rd	ユーザ・エラー・タイプ・レジスタ (USRERRTYPE_REG)	[31:0]: NVMe-IP で報告されるユーザ・エラー・ステータス UserErrorType[31:0]
BA+0x110 Rd	PCIe 割込みステータス・レジスタ (PCIeINTSTS_REG)	[31:0]: NVMe-IP で報告される PCIe 割込みステータス PCIeIntStatus[31:0]
BA+0x114 Rd	完了ステータス・レジスタ (COMPSTS_REG)	[15:0]: NVMe-IP コア内 Admin 完了 (AdmCompStatus[15:0])ステータス [31:16]: NVMe-IP コア内 IO 完了 (IOCompStatus[15:0])ステータス
BA+0x118 Rd	CAP レジスタ (NVMCAP_REG)	[21:0]: NVMe-IP コアからの NVMeCAPReg[31:0]出力
BA+0x11C Rd	NVMe-IP テスト・ピン・レジスタ (NVMTESTPIN_REG)	[31:0]: NVMe-IP コアからの TestPin[31:0]出力
BA+0x120 Rd	比較エラー・アドレス(下位)レジスタ (RDFAILNOL_REG)	[31:0]: リード(ベリファイ)にてデータ比較エラーが発生したアドレス(バイト単位)の下位 32ビット[31:0]
BA+0x124 Rd	比較エラー・アドレス(上位)レジスタ (RDFAILNOH_REG)	[23:0]: リード(ベリファイ)にてデータ比較エラーが発生したアドレス(バイト単位)の上位 24ビット[56:32]
BA+0x130 Rd	期待値ワード 0 レジスタ (EXPPATW0_REG)	[31:0]: リード(ベリファイ)での期待値データ・ワード 0 [31:0]
BA+0x134 Rd	期待値ワード 1 レジスタ (EXPPATW1_REG)	[31:0]: リード(ベリファイ)での期待値データ・ワード 1 [63:32]
BA+0x138 Rd	期待値ワード 2 レジスタ (EXPPATW2_REG)	[31:0]:リード(ベリファイ)での期待値データ・ワード 2 [95:64]
BA+0x13C Rd	期待値ワード 3 レジスタ (EXPPATW3_REG)	[31:0]:リード(ベリファイ)での期待値データ・ワード 3 [127:96]
BA+0x140 Rd	実リード値ワード 0 レジスタ (RDPATW0_REG)	[31:0]: リード(ベリファイ)での実リード値データ・ワード 0 [31:0]
BA+0x144 Rd	実リード値ワード 1 レジスタ (RDPATW1_REG)	[31:0]: リード(ベリファイ)での実リード値データ・ワード 1 [63:32]
BA+0x148 Rd	実リード値ワード 2 レジスタ (RDPATW2_REG)	[31:0]:リード(ベリファイ)での実リード値データ・ワード 2 [95:64]
BA+0x14C Rd	実リード値ワード 3 レジスタ (RDPATW3_REG)	[31:0]:リード(ベリファイ)での実リード値データ・ワード 3 [127:96]
BA+0x150 Rd	現在テスト・バイト(下位)レジスタ (CURTESTSIZEL_REG)	[31:0]: TESTGEN モジュール内の現在テスト・データ・サイズをバイト単位で表示 (bit[31:0])
BA+0x154 Rd	現在テスト・バイト(上位)レジスタ (CURTESTSIZEH_REG)	[31:0]: TESTGEN モジュール内の現在テスト・データ・サイズをバイト単位で表示 (bit[56:32])
BA+0x2000 ~ 0x2FFF	Identify コントローラ・レジスタ (IDENCTRL_REG)	4K バイトの Identify コントローラ・データ空間
BA+0x3000 ~ 0x3FFF	Identify ネームスペース・レジスタ (IDENNAME_REG)	4K バイトの Identify ネームスペース・データ空間

本デザインでの CPU ファームウェア動作は以下のシーケンスとなります

- シリアル・コンソールからのユーザ指示に従い、IDENTIFY DEVICE か WRITE か READ を実行します

#### IDENTIFY DEVICE コマンドの場合

- 1) USRCMD\_REG="00"をセットします。するとテスト・ロジックはコマンドを生成し NVMe-IP に対してコマンドを指示します。ビジー・フラグ(USRSTS\_REG[0])は'0'から'1'に変化します。
- 2) CPU は USRSTS\_REG の値をモニタしコマンドが完了するか又は何らかのエラーが発生するかを確認します。 コマンドが完了した場合 BIT[0]は'0'にネゲートし、何らかのエラーが検出された場合 BIT[1]が'1'にアサートします。エラー発生を検出した場合、エラー・メッセージを表示します。
- 3) 本コマンド実行結果として LBASIZEL/H\_REG レジスタに格納された SSD 全容量情報をシリアル・コンソールへ表示します。

#### WRITE/READ コマンドの場合

- 1) アクセス先開始アドレス、転送セクタ数、テスト・パターンをシリアル・コンソールから受信します。ここで無効なパラメータが入力された場合、コマンド動作はキャンセルされます。
- 2) 入力された各パラメータを USRADRL/H\_REG, USRLENL/H\_REG, USRCMD\_REG (USRCMD\_REG のセット値はライトの場合"10"でリードの場合 "11")にセットします。
- 3) IDENTIFY DEVICE コマンドのステップ 2) と同様にコマンド完了かエラー発生を確認します、ただしリード・コマンドの場合は USRSTS\_REG[2]もチェックしベリファイ・エラーが発生していないかを確認します。
- 4) コマンド実行中は転送済みのデータ数情報を毎秒表示しコマンドが進行中であることを示します。コマンドが完了すると計算したパフォーマンス結果を表示します。

## 4. LAXI2REG

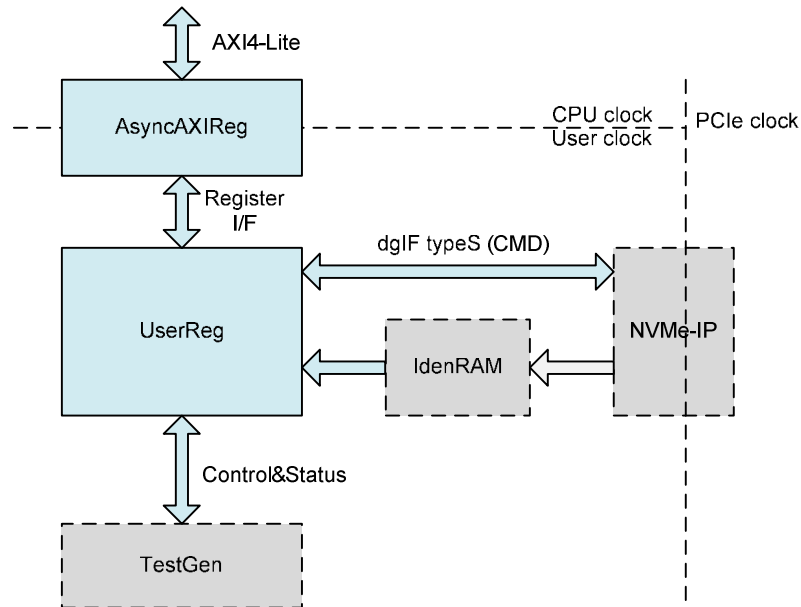


図 4: LAXI2REG モジュール

このモジュールは AsyncAXIReg と UserReg の 2 サブ・モジュールで構成されます。AsyncAXIReg は AXI4-Lite バスをレジスタ・インターフェイスに変換しさらに CPU クロックとユーザ・クロック間のクロック・ドメイン間の変換を行います。UserReg モジュールはライト/リード・アドレスをデコードするロジックを内蔵し、アクセス先のレジスタを選択します。アドレスは表 1 に示すマッピングでデコードされます。ユーザから指定された転送方向・サイズ・アドレスなどの転送パラメータは dgIF typeS の制御インターフェイスに変換され、TestGen モジュールへの制御信号に変換されます。転送中 CPU はレジスタをリードし NVMe-IP コアや TestGen 結果や Identify Device データなどのステータスを監視します。



## 5. TestGen

このモジュールの動作は2種類あります、ひとつは WrFf ポートへ出力するテスト・データの生成でユーザがライト・コマンドを指定した場合に動作します、そしてもうひとつは RdFf ポートからの受信データをペリファイしますがそれはユーザがリード・コマンドを指定した場合に動作します。このモジュールの内部ロジック・デザイン詳細を図 5 に示します。

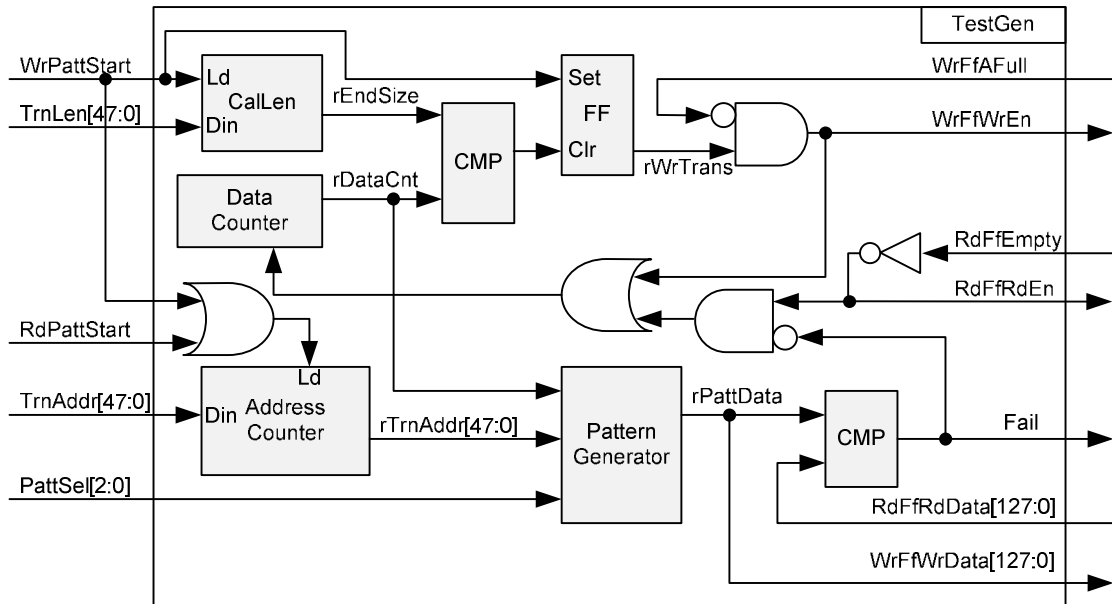


図 5: TestGen モジュール内のロジック・デザイン

ライト転送の開始にて、有効な値がセットされた TrnLen および TrnAddr とともに WrPattStart が '1' にアサートされます。TrnLen はライト転送の終了位置を計算するために使われます。rWrTrans は WrFf ポートへのライト・イネーブル信号として使われますが WrPattStart および終了位置信号(rEndSize)により制御されます。WrFfAFull はライト転送のフロー制御に使います。この信号が '1' アサートされると WrFfWrEn は '0' ネゲートされデータ発生が抑制されます。

このモジュールには 2 つのカウンタが内蔵されます、すなわちひとつはデータ・カウンタ(Sata Counter)で転送の終了タイミングを検出するために総転送サイズをカウントします。もう一つはセクタ (512 バイト) 単位で現在のセクタ・アドレスをカウントするアドレス・カウンタ(Address Counter)です。このアドレス・カウンタは TrnAddr 信号から開始アドレス値をロードします。そして 512 バイトの転送終了毎にインクリメントします。パターン発生器(Pattern Generator)は rTrnAddr からの現在アドレスを各セクタの 64 ビットのヘッダ値として読み込みます。そしてこの値は 32 ビット・インクリメント、32 ビット・デクリメント、32 ビット LFSR のテスト・パターン開始値として使われます。rPattData はライト転送時の WrFfWrData として、あるいはリード転送時の期待値として使われます。

リード転送において RdPattStart はアドレス・カウンタのロード信号としてのみ使われます。RdFfRdEn は RdFfEmpty で制御されます。このデザインでは簡単に RdFfRdEn は RdFfEmpty を反転しています。RdFf ポートからの RdFfRdData がテスト・パターンと不一致であった場合に Fail フラグは '1' アサートされます。

## 6. テスト結果例

本デザインを 512GB Samsung 960PRO で評価した結果を下図 6 に示します。

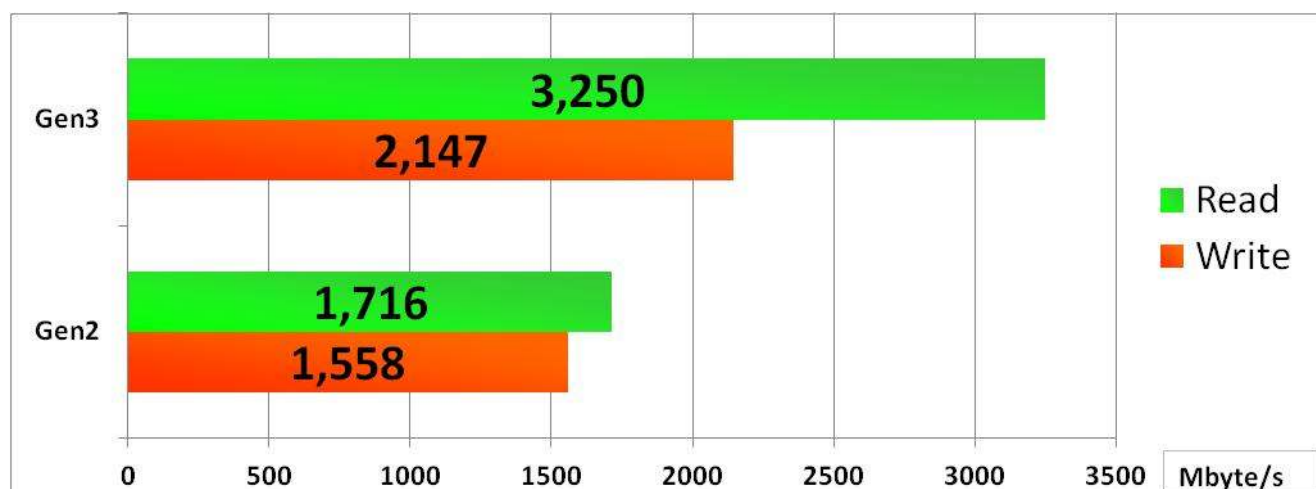


図 6: NVMe-IP デモ・デザインを Samsung 960Pro SSD で評価したパフォーマンス結果

PCIe Gen3 の KCU105 ボードにて Samsung 960Pro SSD を使って評価した結果、ライト・パフォーマンスは約 2100Mbyte/sec でリード・パフォーマンスは約 3200Mbyte/sec でした。一方 PCIe Gen2 の VC707 ボードでは Gen3 よりパフォーマンスは低下し、同じ SSD を使った実機評価にてライトで約 1500Mbyte/sec リードで約 1700Mbyte/sec です。

## 7. 更新履歴

リビジョン	日付	履歴
1.0	1-Jun-16	Initial Release
1.0J	2016/6/7	日本語初期版作成
1.1J	2016/09/06	CURTESTSIZE レジスタを追加
1.2J	2016/12/19	NVMe-IP コアのバッファ・メモリ内蔵版改良によるアップデート
2.0J	2017/06/09	コア改良とデータ・バッファを 256K バイトに固定化した変更に伴う修正
2.01J	2017/07/25	概要の説明文にて誤字を修正
2.1J	2017/07/31	32 ビット LFSR パターンを追加

Copyright: 2016 Design Gateway Co.,Ltd.