

MADOCA II データ収集と蓄積システム

MADOCA II DATA ACQUISITION AND STORAGE SYSTEM

山下明広、籠正裕、酒井久伸
 Akihiro Yamashita*, Masahiro Kano, Hisanobu Sakai
 Japan Synchrotron Research Institute/SPring-8

Abstract

In SPring-8, we are constructing MADOCA II, next generation accelerator control framework. It will be installed in the spring of 2014. We describe the part of the data acquisition and the storage system of MADOCA II. MADOCA was built on the bases of ONC-RPC for communication between processes and a relational database for data management. We designed the new framework with the long experience on MADOCA. We employ Zeromq messages packed by Messagepack for communication. NoSQL databases, Redis and Apache Cassandra, store log data. We obtained a high performance, highly reliable, well scalable and flexible data management system. In this paper, we will discuss requirements, design, implementation and the result of the long run test.

1. はじめに

SPring-8では1997年のコミッショニング以来、加速器とビームライン制御フレームワークとしてMADOCA (Message And Database Oriented Control Architecture)^[1]を使用している。

MADOCAはその名前が示す通り、データの管理をデータベースに依存している^[2]。単一のRelational Database Management System (RDBMS)が機器の管理、運転パラメータの管理、およびログデータの保存に使用されてきた。ログデータベースは加速器とビームラインの全ての機器のログデータを常時記録することにより安定運転と高度化に寄与してきた。

そのなかで、RDBMSをログデータの管理に使用することに以下のような問題点が見えてきた。

1つは新しいデータを登録する際に手間がかかることである。MADOCAでは、1つのRDBMSのテーブルの1つのローに多数のデータを入れる形式をとっているため、新しい信号を登録する際にはテーブルの新規作成をおこなわなければならない。1つのローに多数の信号を入れるのはデータ挿入時のSQL命令を減らし、書き込み性能を向上させるためである。

また1つのローに多数のデータを入れるという性格上、それらを同一時刻に取得せざるを得ず、データ取得の時間的柔軟性に欠けると言う問題点もあった。

性能上の問題も見えてきた。RDBMSはデータの一貫性を再重要視するため処理は同一のメモリー空間内で行うという原則がある。このため性能向上はメモリーを共有したマルチプロセッサシステムの性能向上に頼らざるを得なくなり、性能向上のためには、マルチプロセッサのコア数を増やすか、プロセッサ自身の性能向上が必要である。さらに信頼性向上のためには、ディスクを共有化したクラスター構成か、ハードウェア的にコンポーネントを二重化した構成を取る等の方策がとられる。SPring-8では以前は前者の方法を採用していたが、クラスター構成特有の管理の複雑さやクラスター切り替えの際の空白時間の問題があり、後者に切り替えた。それによりクラスター構成の問題点は克服されたも

のハードウェアが高価になり性能向上には多大なコストがかかる問題点がでてきた。

また現MADOCAのデータ収集系はONC-RPCによりサーバーとクライアントが密結合されていたため、相互の依存性が高くまたOS(Unix系に限定)や言語環境(Cに限定)に制限があった。

これらの問題点を解消し次世代のSPring-8の運転に対応するため、MADOCA IIの開発の一環としてデータ収集データ蓄積のシステムも一新すべく開発をおこなった。新システムの設計方針としては以下の通りである。

- データ管理が容易であること、信号の登録に際してほとんど手間がなく行えること。
- データ蓄積単位が1信号で独立していること。柔軟性と管理性を向上させる。
- 時系列データの蓄積、取得に適していること。
- スケーラブルであること、データ容量、性能向上を低コストで実現出来るものであること。
- 現在のMADOCAシステムからの移行が容易であること、MADOCAは信号にヒューマンリーダブルな名前を命名することにより制御を行っているが、新システムでもそれを採用すること。
- 高信頼であること、データベース、データパスとも複数化し1つの機器の故障が致命的システムダウンに至らないこと。
- データ収集系が疎結合であり、相互依存性を極力排除する。オペレーティングシステムや言語環境にも依存も極力少なくする。
- データベース構成はパラメータは現在のRDBMSを使用し、ログデータベースを改変する。

2. 設計

2.1 データ収集系

データ収集系の概略図をFig. 1に示す。SPring-8の機器制御系はいわゆる3層標準モデルに従っていて、制御対象機器側に置かれた組み込みコンピュータとネット接続されたサーバーコンピュータによって構成される。組み込みコンピュータからZeromq^[3]メッセージの形式

* aki@spring8.or.jp

によりデータが複数の中継サーバーに組み込みコンピュータのタイミングでプッシュされる。中継サーバーはメッセージを複数の書き込み用プロセスに転送する。

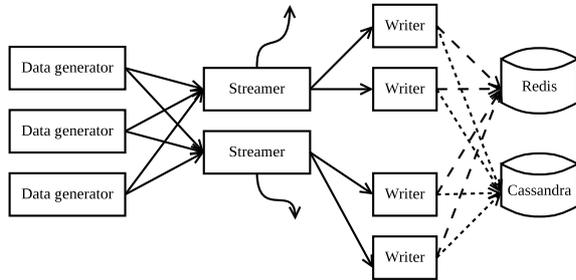


Figure 1: データ収集系 実線は Zeromq Push/Pull、曲線は Zeromq Pub/Sub、破線は Redis API、点線は Thrift RPC。

Zeromq は様々なメッセージングパターンを選択できる多言語多 OS 対応の非同期メッセージングライブラリーである。ここでは組み込みコンピュータから中継サーバーと、中継サーバーと書き込みプロセス間に Push/Pull パターンを使用している。中継サーバーはメッセージを Publish(出版)/Subscribe(購読) パターンでメッセージを送る機能も備えており、購読予約をしたプロセスはデータをリアルタイムで読み込める。Push/Pull パターンは 1 個の送り手から複数の受信者にメッセージをラウンドロビンで送るパターンで、ロード バランシングが容易に実現できる機能を持つ。さらに送り手は受信者のトラブルを検知することができ、トラブルがあった受信者を避けてメッセージを送ることができる。Zeromq のこのパターンを使用し、複数の中継サーバーを使用すれば耐障害性を持ったロード バランシングが実現できる。

中継サーバーは複数の書き込みプロセスにメッセージを送る。後述のデータベースへの書き込み API は、同期型なので、書き込みの間メッセージの受信ができなくなり、メッセージの待ち行列に滞留が起きる可能性がある。これを防止するため、複数の書き込みプロセスに作業を分散させることによりこれを防止している。

また、データ収集プロセスから直接書き込みサーバーに接続せず中継サーバーを使用するのは、Zeromq のソケット管理上の要請と組み込みコンピュータからの接続数を減らし管理を容易にするためである。

Zeromq の Pull ソケットに対してはいつでも複数のクライアントから接続が可能である。したがって任意の時間でデータ収集プロセスを起動し、データ収集を開始または停止することができる。

2.2 データベース系

データベースは永久保存用データベースと現在値専用のキャッシュサーバーで構成される。2つに分離した理由は永久保存用データベースと現在値用のデータベースは要求される機能が異なり、1つのデータベースシステムではそれを実現できないためである。

永久保存用データベースに要求されることは、まず書き込み性能である。現 MADOCA とは異なり、信号 1 つにつき 1 回の書き込み命令を実行するため、現

在の毎秒 7800 信号の書き込みを十分越える性能を要求した。時系列でデータを扱いやすい形式に対応していること、また挿入されるデータが非定型であることスケールアウトにより性能向上が可能なこと、分散された複数のサーバーによって構成されること。データの多重化ができることまた耐障害性のため、各ノードが平等でマスターを作らずの構成されることである。

Google の BigTable^[4] 発表以来大規模な NoSQL (Not only SQL) データベースが web サービスの分野で開発が盛んにおこなわれるようになった。これらは高性能で、安価なサーバーにデータを分散しスケールアウトが可能でデータにも冗長性が持たせられる実装が多い。

われわれが選択したのは Apache Cassandra^[5] でカラム指向データベースとよばれる NoSQL データベースである。これはわれわれが調査した限りでは上記の条件を満たす代表的なデータベースである。これをテストの上実用に十分耐えられるものと判断した。

Apache Cassandra はデータを複数のノードに分散させる。そのため CAP 定理^[6]のうち可用性と耐分断性は満たすものの、一貫性は保証されない。Cassandra では最終的には一致するもののデータ挿入の初期段階においてはノード間でデータが一致しないということがある。Cassandra ではこれを結果整合性と呼ぶ。われわれの測定では 6 ノードでデータを 3 重化した状態でデータを取得した場合、矛盾が起きるのは挿入後 1 秒程度以内である。

この対策としてキャッシュサーバーを用意し、Cassandra とキャッシュサーバーにデータを並行して書き込み最新値はキャッシュサーバーから得るという方式で Cassandra の結果整合性を補償する。

キャッシュサーバーとして Redis^[7] を採用した。キャッシュサーバーに一要求されるのは書き込み、読み込みとも高速で行えるものであることである。書き込みは Cassandra 並みであることが必要であるが、Spring-8 制御では最新データを頻繁に要求する大量のクライアントからのクエリーを処理する必要がある。また 1 つのクエリーで多数の信号の値を取得できる機能も必要である。われわれはインメモリ型のキーバリュデータベースである Redis はこの要求を満たすことがテストの結果わかった。キャッシュサーバーとして以前から人気のある memcached^[8] と比較しても値に 1MB の壁が無い、ストアドプロシージャのような組み込みのアプリケーションが使用できるなどの利点がありこれをテストした結果採用することになった。また mget 命令により一度に大量のデータを高速に取得できる。

3. 実装

3.1 メッセージ構造

組み込みプロセスから中継サーバーに送られるメッセージは 3 つの Zeromq のマルチパートメッセージで構成される。第 1 のメッセージはキーバリュメッセージの 2 文字の接頭辞+信号名+終端文字(:) 第 2 のメッセージはメタ情報を Messagepack^[9] でパックされたマップデータ。第 3 のメッセージも Messagepack でパックされたデータである。図 2 にメッセージの例を示す。第 1 のキーにより中継サーバーから信号名で subscribe

が可能になる。また最初の 2 文字のキーワードはそのメッセージのタイプを示すため、ログデータの場合は LG としているが、他の用途には別のキーワードを使用する。例えばアラーム発生メッセージなどが収集したい場合には AL というような接頭辞をつければ同じ信号名でも使い分けが可能である。

メタデータとデータを別のメッセージにしたのは、メタデータのみで判断分岐が可能となるためである。現在のところはデータ取得時刻のみがメタデータとして登録されているが、マップ形式なので将来の項目追加も容易である。

Messagepack は自己記述型のバイナリ形式のオブジェクトシリアライザである。これは JSON 形式をバイナリ化と言ってよく、JSON で表現できるオブジェクトはこれで高速かつでシリアライズ/デシリアライズすることが可能である。また出力の文字列も JSON でエンコードされたものより短い。

3.2 コンフィギュレーション設定ファイル

このシステムでは多くのホストコンピュータ中で多数の関係するプロセスが複数の通信ポートを使用する。それらを統一的に管理するため、JSON 形式による設定ファイルで全系での統一をおこなっている。

3.3 中継サーバー

中継サーバーは Python でプロトタイプを制作し、動作性能を確認後、C++ で実装した。機能はごく単純で Zeromq の Pull ソケットから受け取ったメッセージを Push ソケットと Publish ソケットに送っているだけである。トラブル時の再接続などの機能は Zeromq ライブラリーに任せてある。1 つの中継サーバーは複数の書き込みプロセスに Push/Pull パターンを使用して接続されている。

3.4 書き込みプロセス

書き込みプロセスも Python でプロトタイプを制作し、動作性能を確認後、C++ で実装した。機能として受信したメッセージを Redis と Cassandra に並行して書き込むこのため内部で構成は受信スレッドと 2 つの書き込みスレッドで構成される。受信スレッドは Zeromq の pull ソケットから受け取ったメッセージを Zeromq のプロセス内通信で publish する (Fig. 2)。2 つの書き込みプロセスはそのメッセージを subscribe し、メッセージを Redis, Cassandra 用に加工し、書き込む。Cassandra は Thrift RPC を使用する。受信スレッドと書き込みスレッドの独立性が高いため、例えば Redis, Cassandra 以外のデータベースに書き込む場合でも、専用のスレッドはメッセージを subscribe するだけなので、受信スレッドには何の変更が必要無い。

3.5 キャッシュサーバー

Redis はその特長としてリスト、ハッシュや集合など豊富な型システムが選べるが本実装では単純な文字列 1 つを値として挿入する。

複数のデータストリームを持つ本データ収集系では必ずしも組み込みのサーバーから送られたメッセージが時間順序通りキャッシュサーバーに到着することは保証できない。対策をしなければ、最新値の後に古い値

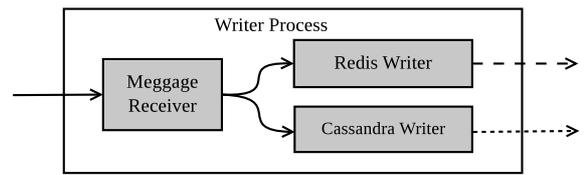


Figure 2: 書き込みプロセスの内部構造。灰色箱はスレッド。メッセージを受け取ったスレッドは Publish/Subscribe パターンを使用してデータベース書き込みスレッドにメッセージを中継する。

が上書きされる怖れがある。この対策のために Redis の機能である組み込みスクリプティング機能を使用した。書き込みスレッドは書き込みの際に Redis の set 命令ではなく Lua 言語で書かれたスクリプトを実行する。このスクリプトは Messagepack でエンコードされた文字列をデコードして、時刻を現在ストアしているものと比較し、それが新しければ上書きする。

キーはメッセージのキー部分と同じ。バリューは時刻と信号値をまとめて Messagepack エンコードされた文字列である。

1 台の Redis の性能が飽和した場合にそなえて、複数の Redis サーバーで並列処理をおこなった。各キーのハッシュ値に対応したサーバーにデータを振り分ける。

3.6 永久保存サーバー

カラム指向型データベースである Cassandra で効率良く時系列データを扱うため 1 つの信号について 1 日分のデータには 1 つの行キーを与える (Fig. 3)。1 つのイベントデータの追加には、時刻をカラム名としたカラムを追加する。という形でデータを蓄積する [10]。この方式は読み込み、書き込み性能は良いものの、1 つのイベントデータが 1 つの文字列に制限される。そこで、複雑なデータ構造は Messagepack 化した文字列で挿入する。この方法では信号名+日データを新たに定義することなく書き込み、信号登録にかかる手間を大幅に低減することができる。

| Column label | | | |
|---------------|------------------------|------------------------|------------------------|
| sig1:20130415 | 2013/04/15 12:00:00 | 2013/04/15 12:00:01 | 2013/04/15 12:00:02 |
| | 0.123 | 0.276 | 6.984 |

Key label value → Columns

Figure 3: Cassandra のデータ構造。ある信号に対する時間データはカラムを増やして収納する。信号はキーを増やすことで登録する。

4. テスト

4.1 書き込みテスト

書き込みの安定性をテストするため、Table 1 のような構成で書き込みの長期間テストを行った。このデータは実際に SACL A で使用されている全信号とそれらのデータ型を模擬したものである。ただしデータ書き込み周期は現行より高速である。この構成で 3 ヶ月間の長期書き込み試験を行った。

Table 1: テスト パラメーター

| | |
|-------------------|-----------------------------------|
| CPU | Intel Xeon X3470 2.93GHz 4Core |
| OS | CentOS 6.2 64bit |
| Redis version | 2.6.10 |
| Cassandra version | 1.15 |
| Java | OracleJavaVM 1.6.0 |
| データー生成プロセス数 | 240 |
| 書き込み信号数 | 47397 |
| 書き込み速度 | 1Hz |
| 平均キー長 | 38.7 bytes |
| メタデーター長 | 13 bytes |
| 平均データー長 | 9 bytes |
| 中継サーバー数 | 3 |
| 書き込みプロセス数/中継サーバー | 8 |
| Redis プロセス数 | 4 |
| Cassandra サーバー数 | 6 |
| Cassandra データー冗長度 | 3 |

データーを 3 重化し信頼性を向上させた。この機能を確認するため Cassandra サーバーを 1 台シャットダウンさせるなどの試験をおこなったがシャットダウン時やその後のリカバリー時でも書き込みに影響は出なかった。

非同期型書き込みという性格上書き込み遅延時間測定はおこなわず、データーの欠損を見たが、この間のデーターの欠損はなかった。

4.2 読み込みテスト

読み込みテストは、実運用時の状態に近づけるため書き込みテストと並行して行った。Table 2 に Redis の読み込み速度の測定結果を示す。この他に高負荷時の場合

Table 2: Redis の読み取り速度 (1 クライアント)

| | |
|------|-----------|
| 試行回数 | 1,000,000 |
| 平均 | 0.26ms |
| 標準偏差 | 0.14ms |

を想定した測定もおこなった。200 クライアントから、毎秒計 20,000 の問合せで平均 1ms、最悪値で 6ms の応答速度であった。

Table 3 に Cassandra の読み込み速度の測定結果を示す。Cassandra の 1 日分のデーター読み込み速度は約 1 秒であるがこれを多数のプロセスで分散読み込みをさせることで 0.2 秒程度にまで短縮されることがわかった。

5. 結論

MADOCA のシステムの設計当時から最近までデーター管理システムでは RDBMS が主流であった。しかし近來の NoSQL 技術の進展により、分散化並列化された高性能、高信頼のシステムが低コストで得られるようになった。また分散コンピューティングのためのミドル

Table 3: Cassandra の読み取り速度 (1 クライアント)

| | |
|---------|---------------|
| 試行回数 | 10,000 |
| 対象データー数 | 86,400 (1 日分) |
| 平均 | 1.01s |
| 標準偏差 | 0.18s |

ウェアについても複雑化(例えば CORBA)をたどっていたがより単純で高速なメッセージングが利用できるようになった。

今回これらの最新の技術を使用することによっていままでの MADOCA ではできなかった柔軟で、管理が容易かつ高性能で拡張性のあるシステムを実現できた。このシステムは 2013 年度のさらなる試験を経て 2014 年度に本格導入される予定である。これらを利用することにより加速器機器の管理がより容易にかつ精密におこなわれることが期待される。

参考文献

- [1] R.Tanaka, et al., "The first operation of control system at the SPring-8 storage ring", Proceedings of ICALEPCS 1997, Beijing, China, (1997) p.1.
- [2] A.Yamahita, et al., "The database system for the SPring-8 storage ring control", Proceedings of ICALEPCS 1997, Beijing, China, (1997) p.427.
- [3] Pieter Hintjens, "Zeromq Messaging for Many Applications", O'Reilly Media, (2013).
- [4] C.Fay, et al., "Bigtable: A Distributed Storage System for Structured Data", Research, Google (2006).
- [5] <http://cassandra.apache.org/>.
- [6] Nancy Lynch and Seth Gilbert, "Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services", ACM SIGACT News, Volume 33 Issue 2 (2002), pg. 51-59.
- [7] <http://redis.io>.
- [8] <http://memcached.org>.
- [9] 古橋貞之 "分散システムのためのメッセージ表現手法に関する研究", 筑波大学大学院博士課程システム情報工学研究科修士論文, (2012).
- [10] <http://www.datastax.com/dev/blog/advanced-time-series-with-Cassandra>.