TSUBAME 共同利用 平成 29 年度 学術利用 成果報告書

利用課題名 スマートデータセンター実現に向けたデータ解析基盤の構築 英文:Towards a data analysis framework for smart data centers

利用課題責任者 松岡聡 Satoshi Matsuoka

所属 産総研・東工大 実社会ビッグデータ活用オープンイノベーションラボラトリ AIST – TokyoTech Real World Big-Data Computation Open Innovation Laboratory https://unit.aist.go.jp/rwbc-oil/

データセンターの大規模化・複雑化に伴い、人力でシステムを監視し、運用設定を変更する従来型手法は困難になりつつある。本研究では従来型運用管理手法の問題を解決すべく、データセンターを構成する各種機器が生成するセンサーデータを蓄積し、解析することにより運用のトレンドを把握し、運用改善を行うシステムパラメータを提案し、運用にフィードバックするシステムの実現を目指す。TSUBAME3.0の利用を通じて、センサーデータ、機械学習手段の調査を行い、センサーデータ収集基盤の基本設計を行った。

Operating datacenters by traditional manual operations becomes unfeasible because they become larger and more complicated systems than before. Our focus in this research is to develop a system that collects sensor data generated from components that compose datacenters, analyze the data and propose system parameters that can improve system efficiency. In this work, we surveyed sensor data of a datacenter and evaluated machine learning libraries on TSUBAME3.0 and then propose a basic design of data collector for datacenters.

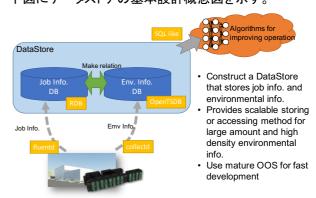
Keywords: smart data center, machine learning, system operation

背景と目的

データセンターの大規模化・複雑化に伴い、人力で システムを監視し、運用設定を変更する従来型手法は 困難になりつつある。その問題を解決すべく、近年では データセンターが生成する各種センサー情報をもとに データ解析を行い、運用改善するパラメータを提案し、 システムに自動適用する研究が行われていたり、 Google など先進的な大規模クラウドデータセンターを 運用する企業ではその技術の一部実用化を進めてい る。しかしながら、それら技術の詳細は公開されていな いこと、また、データセンターのハードウェア構成や、提 供するサービスごとに取得すべきセンサー情報が異な りうるため、他データセンター事業者が模倣することは 困難である。本プロジェクトでは、近年注目されている HPC・AI 向けサービスを提供するデータセンターを事 例に、オープンな技術を用いて、データセンターのデー タ収集・解析・運用へのフィードバックを行うシステムの 実現を目指す。TSUBAME3.0 を用いてセンサーデー タの解析、機械学習手段の調査を行い、センサーデー タ収集基盤の基本設計を行った。

センサーデータ収集基盤基本設計

下図にデータストアの基本設計概念図を示す。



データ取得対象とするデータセンター構成要素として は、以下を想定する。

- 計算ノード
- ジョブスケジューラ(兼、資源管理サービス)
- 付帯設備(電力、冷却システム)

計算ノードについては、各種センサーデータとログデータを時系列データとして取得する。センサーデータについては collectd にて、分単位の粒度で取得することを考える。これは計算ノードで実行されるジョブ 1 つに

つき、ワークロードとして評価されるべきジョブは最低 1 分以上の実行時間は有すべきであること、その時に最 低 1 データポイント以上のデータを取得することが必要 という判断からである。より細粒度で取得することも不 可能ではないが、ネットワーク転送量、蓄積するデータ 量が増える。なお、TSUBAME2.0/2.5 計算ノード相当 のデータを産総研 ABCI 規模(1088 ノード)のシステム で、1 分間隔で取得した際のデータ量は以下の通りとな る。

- 950MB/Node/年
- 1010GB/システム/年

ログデータについては fluentd にて、収集対象のログ出力のたびに取得することを想定する。取得するログ内容については、GPU、SSD などの各種デバイスのヘルスチェック結果、重要なプロセスの死活チェック結果などが考えられるが、運用状況に応じて判断するべき項目である。

ジョブレコードについては、利用者のジョブ完了のたびに取得する。1 つのジョブレコードは複数のデータ型の異なるデータから構成されるため、fluentd にてテキストベースで収集する。付帯設備については、各種センサーデータを時系列データとして取得する。付帯設備は使用する機材によってサポートするデータ取得用のプロトコルが異なるため、サポートされているプロトコルに応じたデータコレクタを採用する。

データストアを構成するデータベースには次の 2 種類を用いる。

- ジョブレコード保存用の Relational データベー ス
- センサーデータ、ログデータ保存用の時系列データベース

ジョブレコードを格納するために Relational データベースを用いる。利用者数・利用状況に応じて、ジョブレコード数は年間数万~100 万規模まで変動すると予想されるが、PostgreSQL や MySQL、MariaDB などの枯れたサーバクライアント型 DB であれば性能・機能的に十分と考えている。Relational データベースではFluentdにて送られてきたデータをレコード形式に変換して DB に登録する。DB 上の 1 レコードサイズを仮に4KB とすると、100 万レコードの場合には4GB のサイ

ズとなる。5年間運用したとしても20GBであり、1ノード 構成で問題ないと考えられる。

ジョブレコード以外のデータは時系列データであり、 その格納には時系列データベースを用いる。現状では、 実装が十分に枯れている OpenTSDB の採用を考えて いるが、近年時系列データベースの開発競争が活発で あることから、システム構築時の状況に応じて変更する ことも考えている。OpenTSDB は collectd、Fluentd の両方に対応したプラグインを有するため、それぞれで 送られてきたデータはそのまま格納できる。ただ、格納 データ形式は数値に限定されるため、ログデータを Fluentd で転送する場合には、意味のある数値データ に変換して送る必要がある。「過去の運用データ解析を 通じてのシステム運用改善」を目的としていることより、 時系列データベースでは格納したデータの時間圧縮は 行わない。これにより、ABCI 規模のシステムを 5 年間 運用した場合、計算ノードのセンサーデータだけでおよ そ 5TB、ログデータや付帯設備のデータも含めると 7TB ほどになると思われる。1 ノードでも蓄積不可能な データ量ではないが、ジョブデータとの関連づけのため には検索性能が重要となるため、分散システムとして 構築した方が良いと考えている。

データ解析用プログラム、運用改善アルゴリズムからのデータストアへのアクセスについては、本データストア用に専用の API を用意するのではなく、様々な側面からデータを解析する需要があると考え、ジョブレコードを保存する Relational データベース、時系列情報を保存する時系列データベースのそれぞれの API をそのまま提供することを考えている。ただし、頻出する検索項目、例えば、特定ジョブが使用した計算ノード群における資源利用データの取得、に関しては、その機能を実現する API を用意する。

まとめ、今後の課題

本研究ではTSUBAME3.0上でのTSUBAME2.0のセンサーデータの解析、機械学習ライブラリの評価を通じて、大規模データセンター用センサーデータ収集基盤の基本設計を行った。今後、性能評価を含め詳細設計を行い、産総研 ABCI システムへの配備を行う。ABCIでの運用を通じて運用データを蓄積し、その解析

技術を確立するとともに、得られた成果は広く公開する。