



**(要件2) 固定サイズ、小サイズメモリ**

データストリームを解析対象にした場合、1つ1つのデータは小さいがその数が膨大であるために大きな計算機リソースが必要となることが多い。また、データストリームは、データ発生が動的に変化するため、この解析には、最大流量に対応できる計算機リソースを予測し、準備しておく必要がある。また、長期間安定して動き続けるためには必要なリソースが一定である必要がある。

これらの要件を満たすために、ストリーム処理アルゴリズムの多くの研究では以下の戦略をとる。「厳密解を計算することはあきらめて、近似解、確率的な解を求める」

一般に、解析精度と、解析に必要な計算機リソースはトレードオフ関係にあり、データストリーム分析アルゴリズムの研究では、解析精度より必要計算機リソースを優先する。これは、この技術の典型的な適用先では、正確な分析結果であるが分析時間がかかることよりも、分析結果の精度は落ちるが分析時間がかからないことを重要視しているためである。

この戦略をふまえ、データストリーム分析アルゴリズムの基本技法を述べる。表-1に示した研究の多くのアルゴリズムでは、その工夫点に注目すると以下の3種類の技法に分類できる。

- (1) **確率統計的計算**：データをサンプリングして確率統計をベースに解析結果を算出する。たとえば、整数カウンタを用いずに確率統計的にかぞえることで効率の良い計算が得られる<sup>1)</sup>。
- (2) **データ粗視化**：詳細な情報を得るためにすべてのデータをメモリに保持するのではなく、データをグループ分けしておおまかな情報、たとえば、グループの代表元や統計量などを保持することでデータを圧縮してメモリを節約する方式である。
- (3) **適応的計算**：あらゆる解の可能性を最初から用意すると計算リソースが大きくなることが多い。そこで、最初は軽い計算により粗い解を求め、計算が進んで情報が増えてきたら、より詳細まで計算する戦略である。

処理内容	例	技法
異なり数 <sup>1)</sup>	購買された商品の種類数	(1)
ホットリスト <sup>2)</sup>	売り上げ最上位K個の商品	(1)
滑り総和 <sup>3)</sup>	一定期間内の売り上げ総和	(3)
組合せ <sup>4)</sup>	同時に買われる商品を検出	(3)
分類規則 <sup>5)</sup>	購買傾向予測ルールを検出	(3)
クラスタリング <sup>6)</sup>	売り上げ傾向が似た商品を検出	(2)

表-1 ストリーム分析アルゴリズム

**分析アルゴリズムの例**

車両を動くセンサと見立て時々刻々と変化する車両情報（位置、速度、ワイパーの利用などであり、以下「プローブ情報」と呼ぶ）をサーバに集約し分析することにより渋滞情報などの交通情報をリアルタイムに提供する「プローブ情報システム」を例にデータストリーム分析アルゴリズムの実例を説明する。プローブ情報システムでは大量のクルマからのプローブ情報をいかに軽量に地図上の道路にマッチングさせるか（マップマッチング）が課題の1つである。これに対し、前記のデータストリーム分析の戦略を応用した(1)モザイクマッチング方式と(2)計算コスト可変近似方式の2つの方式で解決する。

**(1) モザイクマッチング方式**

本方式は、データ粗視化の技法を使った方式である。クルマと道路のマッチングは、交差点付近では細かくマッチングを行う必要があるが、それ以外ではおおざっぱなマッチングで十分であることが多い点を利用して粗視化する。

道路は交差点を意味するノードと、交差点から交差点までの道路セグメントを意味するリンクのネットワークで記述される。マップマッチングとはクルマの位置とリンクを対応させることである。

図-1のaからjの実線がリンクである。本方式のポイントはグリッドを数メートル程度の非常に小さなサイズにし、緯度方向と経度方向に平行にグリッド分けをすることである（図の格子）。このようにすることにより、ほとんどのグリッドは1つのリンクと重なり合っているために、プローブ情報の緯度、経度からグリッドを特定するだけでリンクが特定できる。たとえば、図-1の例では、クルマがいるグリッドはリンクbとのみ重なっているた



監視し、一度に解析するデータが揃えば、解析関数に対象データを渡し解析を依頼する。解析関数は渡されたデータを解析し、後段の処理ノードの待ち行列キューへ解析結果を書き込む。これを繰り返し、全体として解析が行われる。

### データストリーム処理基盤 DSPP

データフロー計算モデルをコンピュータ上で実行させる DSPP (Data Stream Processing Platform) を紹介する<sup>8)</sup> (図-4)。DSPP は、アプリケーションに依存しない基本機能を提供する。この基本機能を使ってアプリケーションに依存する部分、すなわち、解析ウィンドウ、解析関数と、処理ノードネットワークを開発することでシステムを構築する。

DSPP の基本機能は以下のとおりである。

- **処理ノードを実現する基本機能** : DSPP は処理ノードに共通の基本機能を提供し、アプリケーションを構築する際にはこの基本機能に独自の解析関数をプラグインすることにより各処理ノードを実現する。さらに、他の処理ノードへリンクする API (Application Program Interface) を利用して処理ノードネットワークを構築する。
- **処理ノードの実行プロセス、スレッド管理機能** : 各処理ノードを非同期に独立に動作させる。
- **処理ノード間のデータ転送、データ共有機能** : 小サイズのデータの逐次処理に特化した独自のメモリプールを利用して、高速にデータの受け渡しおよび、処理ノードのキューを管理する。
- **処理ノード間のフロー制御機能** : ノードマネージャと呼ばれるモジュールが、処理ノードネットワーク全体の処理の実行状態を管理する。ノードマネージャは、各処理ノードの待ち行列キューが溢れないように、あるいは系全体のパフォーマンスを向上させるために、各処理ノードへ流れるデータ量、処理ノードの処理優先度の制御を行う。

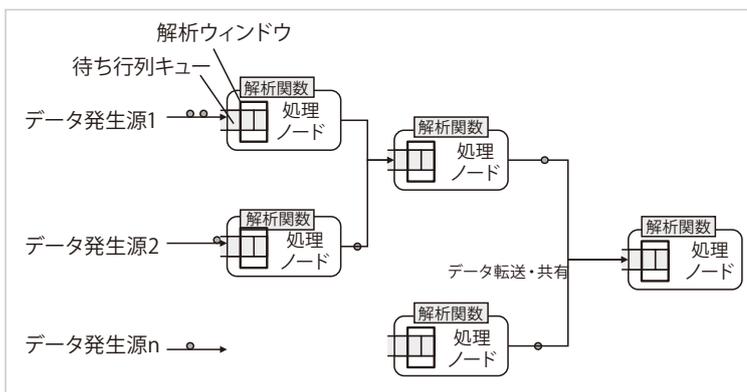


図-3 データフロー計算モデル

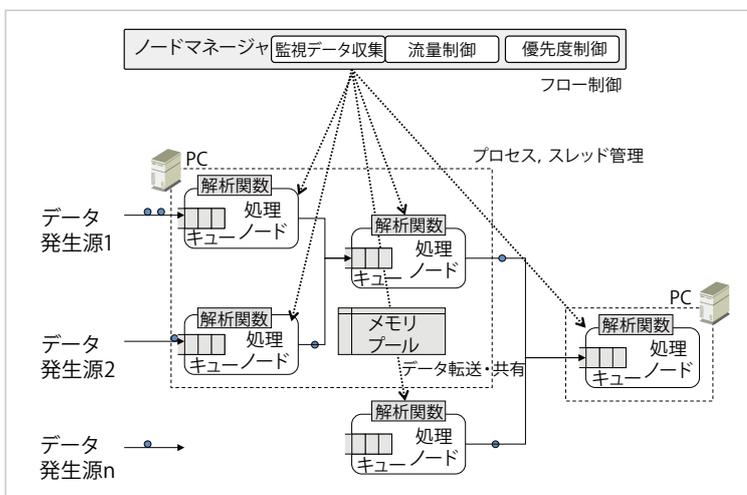


図-4 DSPP アーキテクチャ

## 交通情報システムへの応用

### 実証実験概要

本章では、データストリーム分析を応用して実現した渋滞状況提供実証実験を紹介する。実証実験は、名古屋市にて実施されたプローブカーの実験であり、業務車両（タクシー）から、位置、速度、方向、ワイパーの動作などのデータを収集し、渋滞情報、旅行時間予測、降雨情報を提供するサービスを試行した。実験においては、プローブカーの台数は平均 1,700 台であり、ピーク時は最大 4,000 台が走行し、1 分間隔でサーバにアップロードし、5 分間隔で結果を更新した (表-2)。

この実証実験では、データストリーム方式の有効性を検証するために、メモリデータベースを用いた従来方式と、データストリーム処理の方式の2つの方式で構築し性能比較を行った。



**現状**



幹線道路のみ、数分ごとに更新

**将来 (開発技術を適用)**



細街路を含め、数十秒ごとに更新

図-7 従来方式と本方式で提供する渋滞情報サービスの比較 (地図データはインクリメントP (株) の製品 MapDK を用いて作成されたものです)

ら 30 分前の情報である。渋滞を解消するためには、渋滞の起き始めの対処が重要であることが知られており、可能な限り渋滞情報の遅延を少なくすることが重要となる。これに対し、データストリーム方式の更新頻度は 5 分間隔である。このように鮮度の高い渋滞情報は、信号の制御や、カーナビの経路案内などにより交通をコントロールし渋滞を解消することができる。

**展望**

データストリーム処理技術に関して、分析アルゴリズムと実行基盤に関して解説した。また、交通情報システムへの応用例を紹介した。

今後の ICT (Information and Communication Technology) 産業は PC や携帯電話といった ICT そのものの普及を目的としたものではなく、交通やエ

ネルギー等の社会基盤産業と同期することでさまざまな社会問題の解決に貢献することが期待されている。この実現のためには、実世界の「今」何が起きているのかをリアルタイムに ICT で理解する必要があり、この実現としてデータストリーム分析技術は非常に重要な技術と考えられる。データストリーム分析技術が、人と地球にやさしい社会の実現に貢献できることを期待している。

参考文献

- 1) Flajolet, P. and Martin, G. N. : Probabilistic Counting, in Proc. FOCS'83, pp.76-82 (1983).
- 2) Gibbons, P. B. and Matias, Y. : Synopsis Data Structures for Massive Data Sets, in J. Abello, editor, External Memory Algorithms, Vol.50 of DIMACS Series in Discrete Mathematics, pp.39-70, DIMACS (1999).
- 3) Babcock, B., Babu, S., Datar, M., Motwani, R. and Widom, J. : Models and Issues in Data Stream Systems, in Proc ACM PODS'02, pp.1-16, ACM (2002).
- 4) Manku, G. S. and Motwani, R. : Approximate Frequency Counts over Data Streams, in Proc.VLDB'02, pp.346-357 (2002).
- 5) Domingos, P. and Hulten, G. : Mining High-speed Data Streams, in Proc. ACM SIGKDD'00, pp.71-80 (2000).
- 6) Zhang, T., Ramakrishnan, R. and Livny, M. : Birch : A New Data Clustering Algorithm and its Applications, Data Min. Knowl. Discov., pp.141-182 (1997).
- 7) 喜田弘司, 藤山健一郎, 三津橋晃丈, 中村暢達: 次世代プローブ情報システム (2) ~大規模高速マップマッチングアルゴリズムの提案~, 情報処理学会, マルチメディア分散協調とモバイルシンポジウム論文集 (2007).
- 8) 喜田弘司, 藤山健一郎, 今井照之, 中村暢達: データストリーム処理による大規模プローブカーシステムの開発と評価, 情報処理学会, 高度交通システム研究会, Vol.2008, No.83 (2008).
- 9) IBM InfoSphere Streams, <http://www-06.ibm.com/software/jp/data/infosphere/streams/>
- 10) 日立製作所: uCosminexus Stream Data Platform, <http://www.hitachi.co.jp/Prod/comp/soft1/cosminexus/sdp/>
- 11) Yahoo S4, <http://incubator.apache.org/s4/>

(2012 年 6 月 6 日受付)

謝辞 本研究開発成果の一部は、2008 ~ 2010 年度総務省委託研究「ユビキタスサービスプラットフォーム技術の研究開発」による。

喜田弘司 (正会員) kida@da.jp.nec.com

日本電気 (株) 主任研究員。エージェント, 検索エンジン, システムセキュリティの研究に従事。現在, M2M の研究に従事。博士 (工学)。

藤山健一郎 (正会員) k-fujiyama@cw.jp.nec.com

日本電気 (株) 主任。データストリーム処理の研究に従事。現在, M2M の研究に従事。

磯山和彦 k-iso@bc.jp.nec.com

日本電気 (株) 主任。現在, 複合イベント処理 (CEP) の研究に従事。