

5

スケールアウトの技術

首藤 一幸 東京工業大学

クラウドのソフトウェアには、コンピュータの台数に応じた性能を引き出すこと、つまりスケールアウトが強く求められる。本稿ではスケールアウトの手法を解説し、データストアや、データストアから取得したデータをキャッシュするソフトウェアへの応用を紹介する。

分散データストア

データストアへのアクセス性能や容量を稼ぐためには、次の方策が考えられる。

- スケールアップ / vertical scaling
サーバの性能を強化する。
- スケールアウト / horizontal scaling
サーバの台数を増やす。

スケールアップではすぐに限界に達するし、何よりも、ハードウェアが高つく。2倍の性能・容量を得るために10倍の金額が必要となるといった具合である。そこで、クラウドの構築、またはそこまで行かずとも、1台では不足という場合、スケールアウトを目指すこととなる。

アクセス性能を稼ぐためには、全データを別のサーバに複製するという方法がある。複製元、つまりマスタへの書き込みを、別サーバ上の複製にコピーする。これによって読み出しアクセスを複数台に分散させることが可能となり、読み出しの性能を稼げる。しかし、台数を増やしても容量は稼げず、また、複製を増やすに従って、一貫性の維持、つまり複製の更新がコスト高になっていく。

より多くのサーバを活用したい場合や容量を稼ぎたい場合、データの分割が必要となる。リレーショナルデータベース (RDB) の場合、表ごとに異なるサーバに持たせる方法、表の垂直分割、水平分割という選択肢がある (図-1)。垂直分割とは1つの行を複数の表に分割することであり、通常、表の正規化やセキュリティ等止むを得ない理由で行われる。また、分割後の表の数もたか

が知っている。そのため、以下では水平分割、つまり sharding について考える。

データ (各行) を各サーバにどう分担させるか? データストアへの問合せ内容から担当サーバが決まるのが重要である。さもないと、問合せのたびに全サーバに問い合わせることとなり、アクセス性能を稼ぐことはかなわない。そのためには、データベース設計者が特定の列を分割キー (partition key) として指定し、データベース管理システムはその列の値に従って各行の担当サーバを決める、ということを行う。

これは、リレーショナルデータベースではなく key-value ストアでも同様である。key-value ストアとは、キーにデータ (値) を結びつけて格納し、キーを指定することでデータを検索できるデータストアである。分散 key-value ストアでは、与えられたキーに従って担当サーバを決める。

キーから担当サーバを決める規則は、いろいろ考えられる。たとえば、対応表をあらかじめ用意しておく方法や、キーの値の範囲に応じて担当サーバを決める方法などがある。後者はたとえば、顧客 ID が 0 から 9999 までは1番サーバといった方法である。これらの方法は、やってくるキーやキーごとの頻度があらかじめ分かっている場合はよいが、そうでない場合はうまく働かない。たとえば Web コンテンツの URL がキーとして次々やってくる場合を考えると、対応表や範囲指定では、サーバ群への負荷分散、つまり担当するキーの数を均等に保つことは難しい。

ここでよく行われるのが、キーのハッシュ値に基づく担当サーバ決めである (図-2)。サーバが4台なら、キーに対して0から3までの4通りの値を返すハッシュ

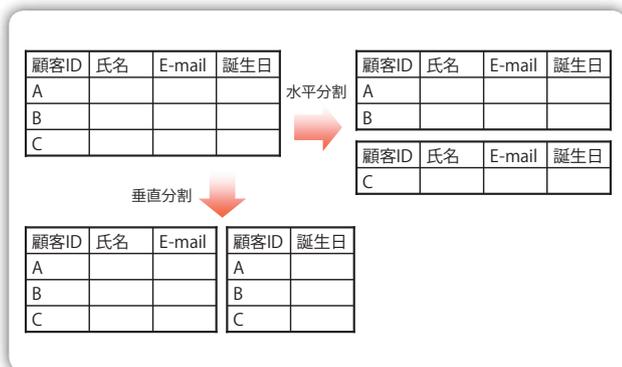


図-1 リレーショナルデータベースでのデータ分割

関数を用意する。キーに対するハッシュ値が1番であれば1番サーバ、というように担当サーバを決める。こういったハッシュ関数の例を挙げる。

$$\text{担当サーバ番号} = H(\text{キー}) \bmod \text{サーバの台数}$$

ここでHとしては暗号的ハッシュ関数、SHA-1やMD5などを使える。ただし、ここで使うハッシュ関数に一方性は不要なので、暗号的ハッシュ関数である必要は特にない。

コンシステントハッシング

ハッシュ値に基づく前述の方法を用いると、未知のキーに対してもサーバごとの担当キー数を均等にできる。しかしこの方法にも、サーバ数の増減には弱いという弱点がある。

性能や容量を増強するためにサーバを増やす、または、減らすという状況を考える。ハッシュ値に基づく方法では、サーバの台数が変わると、ほとんどすべてのキーに対する担当サーバが変わってしまう。たとえば関数Hが返す値が1234だった場合、台数が3から4に増えると、担当サーバは $1234 \bmod 3 = 1$ から $1234 \bmod 4 = 2$ に変わる。サーバがN台からN+1台に増えた場合、全キーの $N/(N+1)$ について担当サーバが変わってしまう。サーバを9台から10台に増やすと全キーの90%について担当サーバが変わる計算である。担当サーバが変わったデータをサーバ間で再配分することは、データが多い場合、ほとんど不可能である。

また、サーバの台数が増えるに従って故障は日常茶飯事となり、増減は日常的に起こるもの、という前提を置かざるを得なくなる。平均故障間隔(MTBF)が15,000時間、つまり1年8~9カ月というサーバがあったとして、その1台が24時間稼働し続ける確率は99.8%である。しかし、10台となると、わずかに44日間で、故障が起き

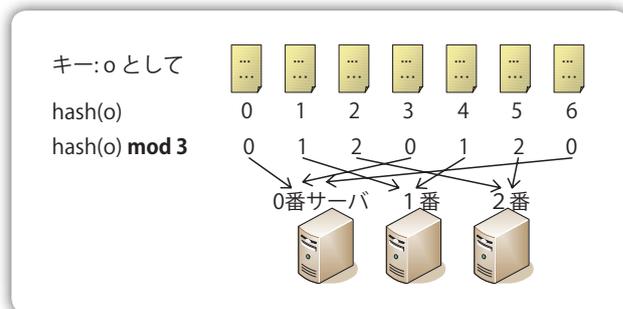


図-2 ハッシュ関数を用いた担当サーバ決め

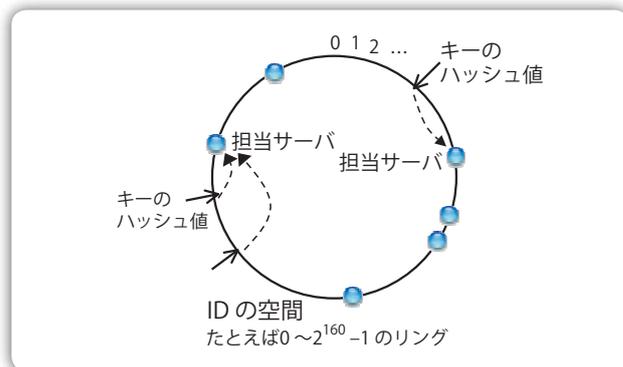


図-3 コンシステントハッシングの条件を満たす担当サーバ決め手法の例(時計回り方式)

ない確率は49.5%と5割を切る。

そこで、サーバの増減に耐える担当サーバ決めの方法が考え出された。ある種の手順で担当サーバを決めると、台数の増減があつた場合でも、台数Nとして、担当サーバが変わるのは全キーの $1/N$ 程度で済む。この性質を満たす担当サーバ決め手法をコンシステントハッシングと言う。これはもともと、複数のサーバでWebのキャッシュサーバを構成するという状況を想定して考え出され、1997年に論文が発表された^{1), 2)}。

この条件を満たす担当サーバ決めの方法は、たとえばこういうものである。 $0 \sim 2^{160} - 1$ の整数値をとるIDの空間を考える。これは暗号的ハッシュ関数SHA-1が返す値の範囲である。図-3の通り、最大値 $2^{160} - 1$ の次の値は0となるリングを想定する。各サーバに $0 \sim 2^{160} - 1$ の範囲からIDを割り当てる。キーにも同様にIDを割り当てる。具体的には、たとえばSHA-1の計算結果をIDとすればよい。図-3の通り、キーのIDからID空間を時計回りにたどり、最初に行き当たったサーバを担当サーバとする。

この手順に従って担当サーバを決めると、サーバの増減時にも、キーの $1/N$ 程度しか担当サーバは変わらない。つまり、コンシステントハッシングの条件を満たす。図-4の通り、3台サーバA, B, CがあつたところにサーバDを追加する状況を考える。すると、サーバBが担当していた範囲の一部がDの担当となる。サーバA

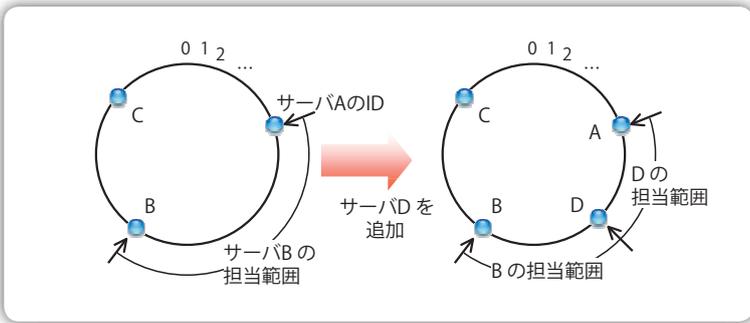


図-4 サーバが増えた場合の担当範囲の変化(時計回り方式の場合)

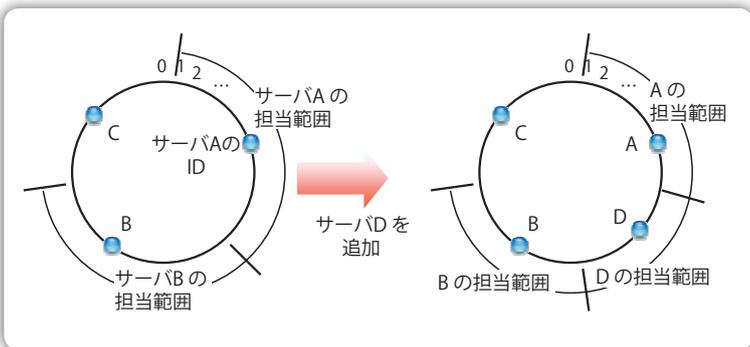


図-5 サーバが増えた場合の担当範囲の変化(最近値方式の場合)

とCに影響はない。

コンシステントハッシングの条件を満たす方法はほかにもある。図-5は、キーのIDから数値的に最も近いIDを持つサーバを担当サーバとする手法を表している。

分散ハッシュ表 (DHT)・構造化オーバーレイネットワークの方式で言うと、Chordが前者の時計回り方式、Pastryが後者の最近値方式を採用している。コンシステントハッシングの条件を満たす担当サーバ決め手法はほかにもあり、リング状のID空間が必須というわけでもない。DHTの方式で言うと、Kademlia、CANなどは他の距離空間を採用している。

サーバ側 vs. peer-to-peer

数台という規模を大きく超えたスケールアウトを狙うデータストアでは、前述の通り、コンシステントハッシングが非常に有用である。コンシステントハッシングを応用したデータストアには、次の2種類が見られる。

- サーバ側で動作するデータストア
- peer-to-peer のデータストア：
分散ハッシュ表(DHT)・構造化オーバーレイ

前者にはデータベース管理システム (DBMS) 自体や、トランザクション機構を持たない、たとえば DBMS が

ら取得したデータなどをキャッシュするデータストアが含まれる。規模としては、数台から多い場合で数百台での動作を狙う。

後者、DHTとは、peer-to-peer、つまり、多数のコンピュータ群が対等な関係で連携して達成するハッシュ表である。正確には、データ構造がハッシュ表である必要はないので、ハッシュ表というよりはkey-valueストアである。規模としては数百万台やそれ以上を狙う。

両者は別種のソフトウェアと見なされているが、実のところ、1つの枠組みの中に位置づけることが可能であり、有効である。以降、具体例を紹介しつつ、両者の関係を論じていく。

■サーバ側データストア

memcachedは、ネットワーク越しにアクセスするkey-valueストアのサーバソフトウェアである。名前の通り、主に、DBMSなどから取得したデータをメモリ上にキャッシュする目的に使われる。Webでの著名なサービス、たとえばFacebook、YouTube、Twitter、Digg、はてな、mixiなどで使われてきている。mixiでは100台以上での動作実績があり、Facebookでは1,000台近くで合計28テラバイトのメモリ容量という実績がある。

memcachedへはクライアントライブラリを通してアクセスすることが一般的である(図-6(a))。クライアント、とあるが、これはたいてい、Webサーバ側で動作するWebアプリケーションである。memcached自体は複数台で連携する機能を持たない。キーに対する担当サーバはライブラリが決める。ライブラリには、ハッシュ値に基づいて決めるライブラリ、コンシステントハッシングで決めるライブラリなど、さまざまなものがある。

memcachedへはクライアントライブラリを通してアクセスすることが一般的である(図-6(a))。クライアント、とあるが、これはたいてい、Webサーバ側で動作するWebアプリケーションである。memcached自体は複数台で連携する機能を持たない。キーに対する担当サーバはライブラリが決める。ライブラリには、ハッシュ値に基づいて決めるライブラリ、コンシステントハッシングで決めるライブラリなど、さまざまなものがある。

Dynamo³⁾はAmazon社が開発したkey-valueストアである。同社のWeb上サービスで使われているという。当初の目標は数百台での動作であり、ショッピングカートといったきわめて高い可用性を求められる应用を見据えている。非常に弱い一貫性モデルである結果整合性(eventual consistency)を採用し、不整合の解決はvector clocksを参照して利用側が行う点、また、クライアントとの間でSLA(service level agreement)について合意する点が特徴である。

各キーを担当するサーバはコンシステントハッシングで決まる。読み書きの担当サーバへの誘導は、クライアントライブラリが行う(図-6(a))こともサーバが行う

(図-6 (b)) こともできる。

Microsoft Azure Platform は、Microsoft 社が 2008 年 10 月に開発者会議 PDC 2008 で発表したクラウド基盤である。Azure 向けのソフトウェアは、Microsoft 自身が運営するデータセンターだけでなく、従来通り、各組織・個人が自身で運用するサーバ上でも動作させられる。これによってクラウド以前のユーザや開発者の Azure への移行を支援する。

Azure も、データストアにコンシステントハッシングを活用している。それどころか、PDC 2008 での当初発表によれば、次に述べる peer-to-peer 由来の技術である DHT・構造化オーバーレイを各所で活用する。たとえば、Blob, Table, Queue といったデータ構造を扱う Windows Azure Storage^{☆1} の基礎部分は DHT そのものである。

■ peer-to-peer のデータストア—さらなるスケラビリティ・耐故障性のために

ここまで説明したサーバ側データストアでは、どこかに全サーバの一覧表があることを前提としている。図-6 (a) では全クライアントが、(b) では全サーバが全サーバの一覧表を持つ必要がある。(a) に描かれているクライアントは 1 台だが、実際の運用では複数のクライアント(アプリケーションサーバ)が動作する。

サーバの台数が増えるに従って、散在する一覧表を最新の状態に保つことが困難になってくる。サーバを落とす際、サーバ側から一覧表を更新できればよいが、故障などで意図せず落ちる場合を想定すると、表を管理する側からサーバに対しての生存確認は必須となる。実際は存在しないサーバが表に載っていた場合、存在しないことはタイムアウト待ちでしか検知できず、この待ち時間がアクセス性能を大きく損ねるからである。

ここでたとえば、サーバ一覧表を持つコンピュータのうち 2 台が代表してサーバ 100 台の生存確認を行う状況を考える。生存確認のたびに問合せと応答合わせて 400 のメッセージが飛ぶので、5 秒に一度生存確認を行うとすると、毎秒平均 80 のメッセージが飛ぶことになる。80 は許容できたとしても、1,000 台で 800, 10,000 台で 8,000 メッセージ/秒を許容できるだろうか。どう工夫しても、秒あたりのメッセージ数は台数に比例せざるを得ない。また、サーバ追加の際は、一貫性を保ちつつすべてのサーバ一覧表を更新しなければならない。完璧を期するためには分散トランザクションが必要となり、一覧表を持つコンピュータの台数に応じてコスト高(可用性が犠牲)、そしてほとんど不可能になっていく^{☆2}。

分散したサーバ一覧表すべての間で一貫性を保とうと

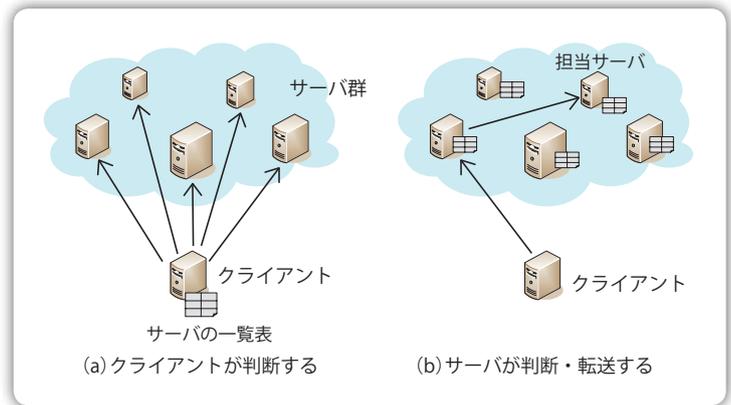


図-6 サーバ側データストア：担当サーバへの到達方法

するから大変になる。ここで、一貫性の条件を緩める方法がある。データストアに対する問合せのサーバ間での転送を許すのである(図-7)。つまり、データの格納、読み出し時に、問合せを受けたサーバが担当サーバではなかった場合、ほかの、より適切と思われるサーバに問合せを回す。サーバ自身が問合せを転送してもよい、問合せ元に対してより適切なサーバを知らせてもよい。前者の方式を再帰探索(recursive lookup)、後者を反復探索(iterative lookup)と呼ぶ。探索のほか、フォワーディング、ルーティングと言うこともある。

転送を許すことで、すべてのサーバ一覧表を常に完璧にしておく必要はなくなる。これによって、サーバ一覧表の維持管理がずいぶん楽になる。メッセージ数を減らせるというだけではない。転送の結果として担当サーバにたどり着ければよいから、全サーバ一覧表の一貫性を保持しなければならない、という制約を緩めることができる。また、一覧表に全サーバを載せる必要もなくなる。表のサイズは、多くの方式で、サーバ数 N として $O(\log N)$ となる。

ただし、転送の結果として担当サーバにたどり着けることは保証されていないと困る。それを保証する方式が peer-to-peer 由来の分散ハッシュ表(DHT)である。ルーティング、つまり転送の経路決め方式という観点では構造化オーバーレイネットワークと呼ばれる。両者の関係は、キーを担当するサーバへのルーティング機能を構造化オーバーレイが提供し、それを使って DHT を実装できる、というものである。

転送を繰り返して担当サーバに到達するということは、つまりはルーティングなのである。ここまでの用語をル

☆1 リレーショナルデータベースである SQL Azure とは異なる。

☆2 複製の作り方によっては、完璧な一貫性までは必要ないようにすることもできる。この種のシステムに複製は必須だが本稿では論じない。

転送なし(図-6)	転送あり(図-7)
<ul style="list-style-type: none"> 担当サーバに直接到達 サーバ側 低遅延 一覧表に全ノードが載る 全一覧表の一貫性維持 ～1,000台? 	<ul style="list-style-type: none"> 担当サーバまで転送数回 peer-to-peer 由来 小さな経路表$O(\log N)$ 経路表間の一貫性は緩くてよい → サーバの頻繁な出入りに耐える ～数百万台以上?

図-9 「転送なし」「転送あり」それぞれの利点

経路表サイズ	経路長	
$O(1)$	$O(\log N)$	constant-degree DHT
$O(\log N)$	$O(\log N / \log \log N)$	
$O(\log N)$	$O(\log N)$	一般的なDHT, 構造化オーバーレイ
$O(N^{1/2})$	$O(1)$	
...	...	
$O(N)$	1	転送なし

図-10 「転送なし」の構造化オーバーレイとしての位置付け

Kai (以上, コンシステントハッシング採用), Flare (グリー(株)), repcached (KLab(株))などを国内の技術者が中心となって開発している。また, 2009年6月に開催された Interop クラウドコンピューティングコンペティションでは参加10チーム中5チームが分散データストアを披露した⁵⁾。開発者のうち最年少は高校1年生という若さであった。

とはいえ, そういったスケールアウト指向データストアも, クラウドを構成する要素技術の一部に過ぎない。他の分野も依然おもしろい。

- (サーバやネットワーク等リソースの) ユーティリティ化
 - 仮想化
 - リソース管理
- 分散システム
 - スケールアウト(本稿)
 - ...
- ...

彼らに続いて, 作れる, 作ろう, という気持ちを持っていただくことが本稿の隠された目的である。敷衍はそう高くない。

参考文献

- 1) Karger, D. et al. : Consistent Hashing and Random Trees : Distributed Caching Protocols for Relieving Hot Spots on the World Wide Web, In Proc. 29th ACM Symposium on Theory of Computing (May 1997).
- 2) Karger, D. et al. : Web Caching with Consistent Hashing, In Proc. WWW8 (May 1999).
- 3) DeCandia, G. et al. : Dynamo : Amazon's Highly Available Key-value Store, In Proc. SOSP 2007 (Oct. 2007).
- 4) Shudo, K. : Overlay Weaver : An Overlay Construction Toolkit, <http://overlayweaver.sf.net/>
- 5) 中田 敦 : memcached を超える成果も, Interop で若手技術者がクラウドを支える技術を競う, 日経 IT Pro, <http://bit.ly/ojRrG> (June 2009). (平成 21 年 9 月 23 日受付)

首藤 一幸 (正会員)

shudo@is.titech.ac.jp

2001年早稲田大学大学院理工学研究科情報科学専攻博士後期課程修了。同年産業技術総合研究所研究員。2006年ウタゴエ(株)取締役最高技術責任者。2008年より東京工業大学大学院情報理工学研究科数理・計算科学専攻准教授。

