



StreamBase boort enorme markt voor real-time stream processing aan

Sensorrevolutie stelt nieuwe eisen aan datamanagement

Peter Verkooijen

De relationele database loopt tegen zijn beperkingen op. Michael Stonebraker, uitvinder van PostgreSQL, werkt sinds de jaren negentig aan andere technologieën die de database moeten aanvullen of vervangen. StreamBase voert query's op datastromen uit. Opslag is optioneel.

Computerwetenschapper Michael Stonebraker was in de jaren zeventig en tachtig het brein achter de Ingres en Postgres databases. Hij bracht zijn technologie op de markt met de bedrijven Ingres, Ilustra, Cohera, Informix en Required Technology. Met zijn nieuwe onderneming laat Stonebraker de database achter zich. "Er zijn allerlei redenen waarom databasesystemen tekort schieten", zegt de oprichter en CTO van StreamBase. "Een reden is dat je data eerst moet invoeren in het systeem."

Met name voor elektronische handel voldoet dat niet meer. Acht jaar geleden verwerkten financiële markten data met een snelheid van 900 berichten per seconde, de snelheidslimiet van relationele databases. Alternatieve oplossingen hebben de verwerkings-snelheid sindsdien elk jaar bijna verdubbeld. "Je moet 50.000 berichten per seconde kunnen verwerken met minder dan tien milliseconden vertraging", vat Stonebraker de eisen aan real-time analytische software samen. "Tien milliseconden is een eeuwigheid in elektronische handel. De huidige database-systemen kunnen die lawines van real-time data niet goed aan."

Real-time posities

StreamBase is een van de marktleiders in een nieuwe generatie software die analyse op datastromen toepast nog voor de gegevens de database bereiken. Producten voor event stream processing (ESP) die nu op de markt zijn, kunnen 25.000 tot 150.000 berichten per seconde verwerken. Een sleutelement in stream processing is een event processing language (EPL) of een event query language (EQL) die query's op de binnenkomende data kan uitvoeren. Afhankelijk van de geformuleerde regels kan het systeem actie op binnenkomende data ondernemen.

Relationele databases zijn ontwikkeld voor analyse van statische, opgeslagen informatie. Querytaal SQL kan data alleen op waarden verbinden en de output is meer data. Stonebraker en collega's hebben de taal StreamSQL ontwikkeld die aan SQL de factor tijd

toevoegt en acties als output heeft. "Onze software leest berichten van de TCP/IP socket", zegt Stonebraker. "StreamSQL is zeer flexibel in wat je het kunt laten berekenen. Je kunt bijvoorbeeld zeggen, geef een TCP output bericht als het momentum van IBM over de laatste vijf tikken meer dan twintig procent groter is dan het momentum van HP."

Stonebraker was als professor aan het Massachusetts Institute of Technology (MIT) een van de theoretische grondleggers van stream processing. De focus van het onderzoek waaruit StreamSQL voortkwam was niet Wall Street, maar sensornetwerken bij het Amerikaanse leger. Stonebraker werkte met professor Stan Zdonik van Brown University en een team van dertig faculteitsleden en studenten samen in het project Aurora. Het leger gaf de wetenschappers van 2001 tot 2003 de gelegenheid prototypen te ontwikkelen en in de praktijk te testen.

Het is onacceptabel voor een real-time engine om te blokkeren en te wachten

"De mensen van het leger gaven ons een hele serie taken", vertelt Stonebraker. "Het is een heel serieuze zaak voor het Amerikaanse leger om zicht te houden op alle levenstekens van hun mensen en materiaal. Ze plakken een draadloze sensor op iedere soldaat en een verzameling sensors in ieder voertuig. Een van onze taken was de real-time posities berekenen van een peloton van twaalf soldaten, zodat een bewegende stip op de kaart in het hoofdkwartier van de generaal te zien is. We ontdekten dat we voor allerlei subtiele problemen oplossingen moesten vinden. Wat doe je als iemand over de heuvel buiten het radiocontact loopt?"



Michael Stonebraker: paradigkawijziging in datamanagement is onvermijdelijk.

Je kunt het systeem niet laten wachten tot de twaalfde man rapporteert. Het is onacceptabel voor een real-time engine om te blokkeren en te wachten. In de database-wereld moet je een operatie gelegenheid geven time-out te gaan. In de real-time streaming datawereld moet je iets doen met ontbrekende data. In de opgeslagen SQL-wereld bestaat dat niet."

Overheidsmarkt

Voor financiële instellingen die datastromen van verschillende leveranciers krijgen, bleek het vaak nodig duplicaten uit te filteren en de beste data bij elkaar te vegen. "Je moet bepalen of een tik die binnenkomt de eerste is of een herhaling", zegt Stonebraker. "Je slaat ze dus op in een tabel zodat je kunt vergelijken en doublures kunt weggooien. StreamSQL kan dat in een paar regels code doen. In conventioneel SQL heb je SELECT FROM en het ding waaruit je selecteert is de tabel. StreamSQL breidt SQL uit zodat de dingen in de FROM clause zowel berichtenstromen als opgeslagen data kunnen zijn."

Project Aurora leidde in 2003 tot het bedrijf StreamBase na een investering van Bessemer Venture Partners en Highland Capital Partners. In februari 2005 lanceerde de startup de eerste versie

van zijn product, de StreamBase Stream Processing Engine.

"Je krijgt van ons de programmeernotatie StreamSQL en we compileren StreamSQL in wat je een applicatie zou kunnen noemen. Die applicatie draait op een of meer servers, leest de berichten, voert de noodzakelijke berekeningen uit en produceert de output-berichten. We zijn over het algemeen een alternatief voor custom C++ code op de hardware."

Ook StreamBase is nog grotendeels maatwerk, maar StreamSQL en een *drag and drop* Java-interface maken het mogelijk in een paar dagen een toepassing te ontwikkelen. StreamBase biedt potentiële klanten meestal een pilot-programma om te bewijzen hoe snel de oplossing is in te zetten. Een ander verkoopargument dat Stonebraker benadrukt is de snelheid van zijn product ten opzichte van de concurrentie. StreamBase claimt in specifieke toepassingen op een doorsnee PC 140.000 berichten per seconde te kunnen bewerken.

De eerste klanten van StreamBase waren behalve Wall Street, het leger en inlichtingendiensten zoals de NSA en CIA. "Als Tom Ridge, het voormalige hoofd van Homeland Security, zei dat de 'chatter' was toegenomen, dan waren dat de drie letter-instansies die analyse toepasten op real-time spraakverkeer", zegt Stonebraker. "We hebben veel succes in de overheidsmarkt." Sindsdien heeft Streambase zijn klantenbestand uitgebreid in onder andere de netwerkbeveiliging- en telecom-sector. "En raar maar waar, een recente klant gebruikt onze software om een internet-game met duizenden spelers te managen."

Prestatie

StreamBase is niet de enige ESP-leverancier. Concurrenten zoals Apama, Celequest en Actimize passen volgens Stonebraker een andere benadering toe gebaseerd op *rule engines* uit de wereld van kunstmatige intelligentie. Coral8 en Aleri Labs nemen net als StreamBase SQL als uitgangspunt. "Het verbaasde mij dat Coral8 ervoor heeft gekozen hun engine als een database-applicatie op het primaire geheugen database-systeem te bouwen", zegt Stonebraker. "Dat gaat gewoon niet erg snel. Onze prestaties zijn aanmerkelijk beter dan die van Coral8 en ook Aleri. Deze markt draait volledig om prestatie."

Weegt de hardware niet veel zwaarder in die prestaties dan de software? Komt de echte concurrentie voor StreamBase van de nieuwe generatie multi-core en multi-thread chips en parallel processing? "Dat is een van mijn favoriete onderwerpen", zegt Stonebraker. "Multi-core chips zijn de toekomst. Om te kunnen concurreren moet je systeem multi-threaded zijn en alle beschikbare processor-kernen kunnen inzetten. Onze engine voldoet aan die eisen en werkt ook op verschillende besturings-systemen zodat je parallel processing kunt inzetten."

De gedistribueerde, multi-processor capaciteiten zijn op MIT ontwikkeld in het Borealis-project dat volgde op Aurora. "Een CPU-board met vier kernen kost tegenwoordig 700 dollar", zegt Stonebraker. "Ons standpunt is dat je gewoon meer systemen op elkaar moet kunnen stapelen wanneer je meer rekenkracht nodig hebt. Het probleem met een hardware-oplossing is dat je een

kaart van 700 dollar vervangt door een kaart van 10.000 dollar. Dat moet dan wel een factor 13 of zo aan prestatievoordeel opleveren. In ons soort toepassing zie ik niet waar dat vandaan moet komen."

StreamBase bewerkt data wanneer ze binnenkomen. Data opslaan is optioneel. "Sommige mensen willen een historisch overzicht houden", zegt Stonebraker. Aangescherpte boekhoudregels in de VS en Europa dwingen dat de afgelopen jaren af. Stream processing gaat tegen die trend in. Voor Stonebraker is een paradigma-wijziging in datamanagement onvermijdelijk. "Er is een sensornetwerkrevolutie gaande waarvan RFID maar een klein onderdeel is. Er zal een enorme markt zijn voor real-time stream processing. De vraag komt van Wall Street, het leger, netwerk-beheer en een hele reeks sensornetwerk-toepassingen."

Verwerking van snelle datastromen is niet het enige terrein waarop de relationele database tekortschiet. "De wetenschappelijke gemeenschap heeft array data en is erg ongelukkig met database-systemen", zegt Stonebraker. "Google gebruikt geen database voor tekst-processing. De relationele database-leveranciers verkopen architecturen die zijn bedacht voor verwerking van bedrijfsdata een kwart eeuw geleden. Applicaties die nu verschijnen, zoals de sensornetwerken, passen niet op die bestaande architecturen. De leveranciers kunnen proberen voor elk van die markten een aparte engine te schrijven en zich het hoofd breken hoe ze die weer aan elkaar kunnen lijmen. Of ze kunnen markten die niet groot genoeg zijn gewoon negeren. Hoe dan ook, ze staan voor een interessante uitdaging."

Federated

De bedenker van Postgres heeft de relationele database nog niet helemaal opgegeven. "Relationele databases werken prima voor verwerking van bedrijfsdata", zegt Stonebraker. "Maar andere technologieën zullen de database aanvullen. De boel gaat anders ontspreiden. Je hoort veel over web-diensten, maar web-diensten over grenzen van bedrijven heen gebruiken zal toch een vorm van semantische federatie van databases vereisen. Dat is een ander

Progress Apama marktleider?

Michael Stonebraker doet concurrent Apama af als een 'rule engine' die onmogelijk de verwerkingsnelheid van systemen op SQL-basis zoals StreamBase kan benaderen. Een nieuw onafhankelijk rapport van Bloor Research roept Apama echter uit tot de onbetwiste marktleider in event stream processing (ESP).

Het Britse bedrijf Apama is vorig jaar overgenomen door Progress Software. Volgens het rapport heeft deze gevestigde leverancier, met een wereldwijd klantenbestand in de financiële wereld, Apama's software geïntegreerd in een breder platform. Progress Apama zou 'het enige complete event processing platform' zijn. Het kan niet lang duren voordat een andere zakelijke software-leverancier StreamBase overneemt.

MIT-project waar ik aan werk." Stonebraker heeft eerder aan een federated database gewerkt in het Mariposa-project waaruit het bedrijf Cohera voortkwam.

XML is geen panacee voor gegevensuitwisseling tussen databases. "De mensen die web-diensten verkopen zetten XML op een voetstuk om de semantische problemen met web-diensten over bedrijfsgrenzen heen te maskeren", zegt Stonebraker. StreamBase kan XML-berichten verwerken, maar het gaat niet van harte. "Om in de zakelijke markt mee te kunnen spelen moet je alle bestandsformaten kunnen lezen. We hebben een adapter die XML meteen in binaire code vertaalt. Data uit XML halen is het eerste dat je moet doen als je snelheid wilt. Je kunt onmogelijk 100.000 berichten per seconde verwerken in XML. Als je tegen Wall Street zegt dat je een engine hebt die onvertaald XML verwerkt, lachen ze je uit."

Peter Verkooijen is freelance journalist.

BI-matrix

De BI-matrix geeft u een overzicht van leveranciers op het gebied van Business Intelligence op de Nederlandse markt en werd voor u samengesteld door Paul van der Linden, senior consultant DWH en BI bij Atos Origin.

In DB/M 5 drukten wij een verkleinde vorm van de BI-matrix af. Vanaf nu bieden wij u de uitgebreide vorm van de BI-matrix aan op onze website www.dbm.nl. De BI-matrix on-line is een doorzoekbare database, waar u zelf kunt selecteren welke gegevens u van welke BI-leverancier naast elkaar wilt zien. Handig als u een shortlist van bij u passende leveranciers wilt samenstellen.



De BI-matrix is een initiatief van

