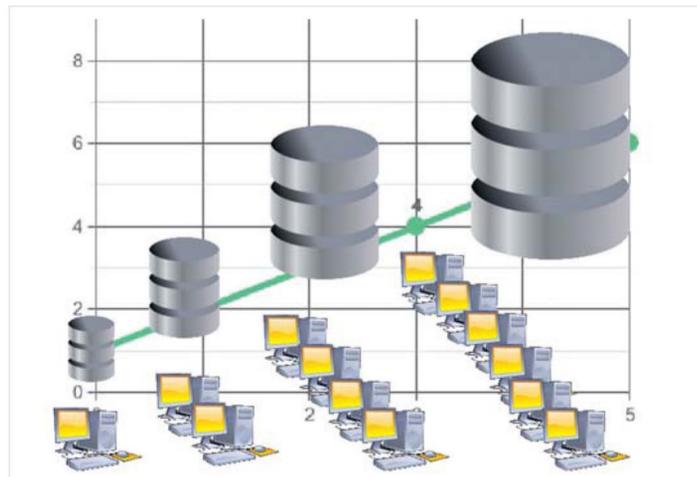


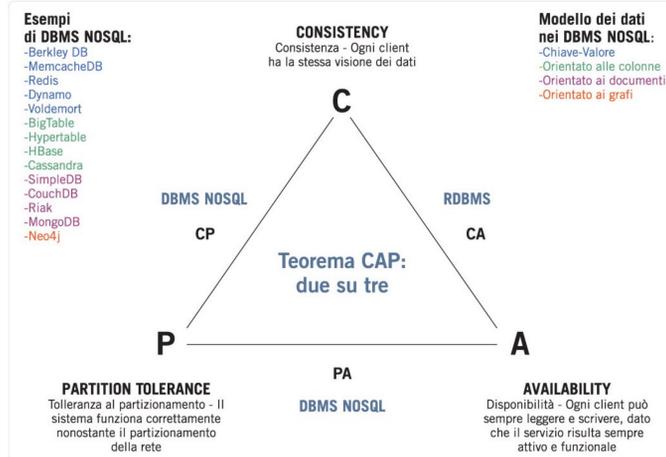
# Not Only SQL

Rodolfo BORASO Diego GUENZI

Internet è sede di continue innovazioni ed è grazie alla sua evoluzione che si sta affermando una generazione di tecnologie in grado di sfruttare al meglio alcune caratteristiche come l'affidabilità e la capacità di banda per offrire un nuovo paradigma legato alla memorizzazione dei dati. Attraverso le reti, è oggi possibile effettuare operazioni massive sui dati andando a suddividere il carico su macchine differenti e, allo stesso tempo, garantendo la sicurezza e l'integrità dei propri dati. In questo ambito si colloca il NOSQL, un movimento che promuove una classe non ben definita di strumenti di archiviazione di dati; questi ultimi si differenziano dalle basi di dati relazionali (RDBMS) in quanto non utilizzano il linguaggio di interrogazione SQL, non richiedono uno schema fisso e puntano a scalare orizzontalmente utilizzando un comune hardware da aggiungere ai propri sistemi. Il nome NOSQL non vuole indicare una contrapposizione all'utilizzo del linguaggio SQL ma intende segnalare la possibilità di utilizzare una vasta gamma di strumenti per compiti particolareggiati, secondo la filosofia del **the right tool for the job**. I database NOSQL, infatti, nascono con l'obiettivo di coprire alcuni settori specifici dove i RDBMS offrono prestazioni ridotte o presentano delle complessità di gestione, affiancandosi al loro utilizzo come strumenti complementari. Esempi classici di problematiche che vanno ad affrontare sono l'utilizzo in ambienti distribuiti geograficamente, la replica dei dati, la fault tolerance e l'high availability, tutti fattori molto sentiti nelle attuali reti di calcolatori.



Negli ultimi anni sono nati molti esempi di database NOSQL (definiti structured storage nel settore accademico/scientifico), buona parte in ambito open source. Il loro modello dei dati differisce dal classico concetto di tabella ed è di diverse tipologie: andiamo dal semplice storage di coppie chiave-valore ad un modello basato su strutture a grafo, dai database orientati alle colonne a quelli orientati ai documenti. Ognuno di questi modelli ha pregi e difetti, oltre ad adattarsi più o meno bene alle diverse esigenze che possono sorgere. Scegliere quale database utilizzare significa comprendere le problematiche che si vogliono risolvere e studiare, all'interno della moltitudine di prodotti disponibili, quale si adatta maggiormente alle proprie necessità.



Questi database tendono ad allontanarsi dal modello ACID (Atomicity, Consistency, Isolation, Durability), tipico delle basi di dati relazionali, per abbracciare il modello BASE (Basically Available, Soft state, Eventually consistent) dove alcuni vincoli vengono rilassati. Come esplicitato dal **teorema CAP**, non si può avere Consistency, Availability e Partition tolerance al tempo stesso ma solamente due di queste caratteristiche alla volta. Se si sceglie di non avere Partition tolerance (la filosofia dei database relazionali, dove i dati sono solitamente presenti su di un'unica macchina, per evitare partizionamenti) avremo dei problemi a scalare in quanto possiamo solamente utilizzare la modalità verticale, sicuramente più costosa. Se scegliamo di rinunciare all'Availability, al verificarsi di un partizionamento dobbiamo attendere che esso venga risolto prima di soddisfare alcune richieste, a discapito, quindi, delle prestazioni globali. Se rinunciamo alla Consistency avremo la possibilità, in alcuni casi e per un determinato periodo, che un partizionamento generi un disallineamento nelle repliche dei dati. Questi ultimi due casi rientrano nel modello BASE e sono quelli che ci permettono di superare il più grosso limite dei RDBMS attuali: la scalabilità. In molti settori, la consistenza e/o la disponibilità offerte dai database relazionali non sono strettamente essenziali; infatti, recentemente, sempre più aziende stanno spostando i propri sistemi informativi dalle tradizionali basi di dati relazionali ai DBMS NOSQL.

Per citare alcune delle principali società che hanno migrato i propri dati su sistemi NOSQL, possiamo prendere in esame il caso di Facebook che ha deciso di sviluppare internamente il proprio database scalabile (ispirato in parte a BigTable di Google e in parte a Dynamo di Amazon Web Services). Questo software, dopo poco tempo, è divenuto open source e uno dei progetti principali di Apache Software Foundation, conosciuto con il nome di Cassandra. In seguito a questo episodio, molti altri social network (Twitter e Digg, per segnalare i principali) hanno deciso di migrare verso Cassandra per la gestione dei loro dati, sebbene esso non sia il solo **database specializzato** nato per uno scopo ben preciso. Oltre ai social network, molte altre società hanno avuto necessità analoghe e hanno fatto ricorso a prodotti della famiglia NOSQL, spesso abbandonando i RDBMS general purpose che adottavano in ambienti di produzione da anni; alcuni esempi possono essere Sourceforge che ha scelto MongoDB come back-end per la gestione dei progetti, Yahoo che ha sviluppato PNUTS per risolvere i problemi di distribuzione geografica dei dati all'interno dei suoi data center, Adobe che dal 2008 ha scelto HBase per la memorizzazione dei dati per uso interno o LinkedIn con l'adozione di Voldemort per lo storage dei dati (a cui è seguito il rilascio dei sorgenti verso le comunità open source, esattamente come è accaduto per Cassandra).



Questi database sembrano sposarsi perfettamente con la crescente metodologia di **computazione cloud-based** che richiede l'utilizzo di ambienti fortemente distribuiti, scalabili e dinamici, all'interno dei quali il numero di macchine (virtuali o meno) è in continua variazione. Inoltre, questi modelli di memorizzazione dei dati si stanno imponendo sul mercato evidenziando i limiti dell'attuale rete globale e delle architetture applicative utilizzate: il connubio fra i DBMS NOSQL e una maggior resistenza e capacità degli apparati di connettività permetterebbe di risolvere il problema dell'accesso e del trasferimento di grosse moli di dati da una parte all'altra del pianeta. Su questi temi, il CSP ha avviato uno studio sullo stato dell'arte delle principali soluzioni di basi di dati NOSQL, al fine di analizzare le problematiche tipiche della memorizzazione di informazioni in ambienti distribuiti geograficamente mantenendo l'accesso ai dati anche in presenza di guasti. Per ciò che riguarda il **futuro di Internet** in relazione a questi sistemi di memorizzazione, la visione che il CSP ha è quella di una rete sempre più flessibile e reattiva, pronta a rispondere ad esigenze differenti e variabili nel tempo; oltre a questo, si prevede un aumento nell'utilizzo di reti dedicate (fisiche o virtuali) e di disponibilità di banda: quest'ultima in particolare, viene dettata dalle esigenze dei servizi di storage distribuito (strutturato e non) ma anche da tutto ciò che, più in generale, utilizza la rete come strumento di interscambio efficiente di informazioni.

