

# Relační a NoSQL databáze: dvě strany téže mince?

Jaroslav Pokorný

Katedra softwarového inženýrství, MFF UK Praha  
Malostranské nám. 25, 118 00 Praha

`pokorny@ksi.mff.cuni.cz`

**Abstrakt.** Analýza vlastností relačních a NoSQL databází vede k závěru, že tyto systémy pro zpracování dat jsou do jisté míry komplementární. V současných aplikacích pro Big Data, speciálně tam, kde jsou nutné rozsáhlé analýzy, se pak ukazuje, že je netriviální navrhnout infrastrukturu zahrnující software obojího typu. Z hlediska výkonu může být dokonce přínosné transformovat schéma SQL databáze do NoSQL anebo provádět oboustrannou migraci dat mezi relační a NoSQL databází. Cílem článku je diskutovat tyto možnosti a zejména některé nové metody návrhu takových databázových architektur stojící na dědictví tříúrovňové ANSI/SPARC architektury.

**Klíčová slova:** Relační databáze, NoSQL databáze, Big Data, Big Analytics

## 1 Úvod

V poslední době se zdá, že se většina velkých podniků minimálně stará o údržbu podnikových aplikací stávajících systémů. To způsobuje, že se používají "špatná" databázová schémata a obecně dochází k "úpadku databází" [10]. Autoři tvrzení vychází z diskusí s téměř dvaceti správci databází (DBA) u tří velkých podniků. Databáze se mění v závislosti na podmínkách byznysu, běžně jednou za čtvrtletí i více. Heterogenní a dynamické datové prostředí vede k tomu, že často mizí role centrálního správce a objevuje se spíše decentralizovaný přístup s více skupinami DBA zabývajících se databázemi v podniku.

Databáze jsou většinou relační. Od zavedení relačního modelu dat bylo sice zavedeno několik databázových modelů, jako je objektově-orientovaný (OO), objektově-relační (OR), XML či RDF. OO a OR SŘBD reagovaly na objektové přístupy v softwarovém inženýrství z 90. let. Tyto prostředky, však nikdy na trhu nebyly skutečně konkurenceschopné. Důvody by mohly být v nedostatku jejich teoretických základů a omezené výkonnosti.

Dnes situaci v databázovém světě ovlivňují tzv. Big Data. Jejich základní V-charakteristiky jsou objem (Volume), rychlost (Velocity) a různorodost (Variety). Autor práce [9] uvádí dokonce 14 takových V. Ty zásadně ovlivňují infrastrukturu ukládání a zpracování Big Dat. Efektivní využívání systémů zahrnujících zpracování velkých objemů dat vyžaduje v mnoha aplikačních scénářích odpovídající nástroje

*J. Steinberger, M. Zíma, D. Fiala, M. Dostal, M. Nykl (eds.)  
Data a znalosti 2017, Plzeň, 5. - 6. října 2017, pp. 8-14.*

pro jejich ukládání na nízké úrovni a analytické nástroje ve vyšších úrovních. Zdá se, že z pohledu uživatele je nejdůležitějším aspektem zpracování velkých objemů dat na počítači právě jejich analýza, jak se dnes říká - Big Analytics. Bohužel, velké kolekce dat obsahují data v různých formátech, např. relační tabulky, XML data, textová data, multimediální data nebo RDF trojice, což může působit potíže při jejich zpracování algoritmy pro dolování dat (DM). Rovněž zvyšující se objem dat v úložišti a počet jeho uživatelů vyžaduje spolehlivé řešení škálování v těchto dynamických prostředích a pokročilejší prostředky pro zajištění vysokého výkonu, než nabízejí tradiční databázové architektury.

Je zřejmé, že Big Analytics se provádí i nad velkým množstvím transakčních dat rozšířením metod používaných v datových skladech (DW). Technologie DW ale vždy byla zaměřena na strukturovaná data ve srovnání s bohatší variabilitou typů dat, tak jak je dnes aktuální pro Big Data. Analytické zpracování velkých objemů dat proto vyžaduje nejen nové databázové architektury, ale také nové metody pro analýzu dat.

Pro ukládání a zpracování Big Dat lze dnes volit tradiční SŘBD, paralelní DBS, distribuované souborové systémy (např. HDFS), datová úložiště typu klíč-hodnota (NoSQL databáze) a nové databázové architektury (NewSQL databáze). Pro volbu technologie jsou rozhodující aplikace, které mohou být jak transakční tak analytické. Požadují obvykle různé architektury software i hardware, často v jedné infrastruktuře.

Cílem článku je diskutovat vztah SQL a NoSQL databází v tomto polyglotním světě, a to hlavně směrem k Big Analytics. Důležitá je skutečnost, že tyto databáze mají komplementární vlastnosti [4], což i motivuje používat je v jedné infrastruktuře. V sekci 2 stručně popíšeme koncept Big Analytics. Sekci 3 podává stručný přehled technologií NoSQL databází. V sekci 4 ukážeme dualitu mezi SQL databázemi a NoSQL databázemi. Sekce 5 obsahuje závěry a výzvy pro databázovou komunitu.

## 2 Analytické zpracování Big Dat

Big Analytics slouží k proměňování informací ve znalosti pomocí kombinace stávajících a nových přístupů aplikovaných na Big Data. K souvisejícím technologiím patří:

- správa dat (uvažující nejistotu, zpracování dotazu v téměř reálném čase, extrakci informací, explicitní správu časové dimenze),
- nové programovací modely,
- strojové učení (ML) a statistické metody,
- komponentové architektury systémů ukládání a zpracování dat,
- vizualizace informací.

Obvyklé se rozlišují dva typy zpracování: reálnodobé zpracování „dat v pohybu“ (*data-in-motion*) a dávkové zpracování dat získaných z různých zdrojů do např. jedné databáze (*data-at-rest*). Dávková analýza pak může být: *malá* (Small Analytics), tj. OLAP a DW, a *velká* (Big Analytics), tj. DM, ML, e-science.

Problémy, které se v této souvislosti vyskytují, vycházejí z faktu, že požadavky na Big Data jsou často dynamičtější než klasické zpracování dat v DW. Dalším problémem je, jak analyzovat Big Data pocházející z relačních DB. Velký objem je nejen

problémem pro ukládání dat, ale ovlivňuje také Big Analytics. S nárůstem složitosti dat je rovněž složitější i jejich analýza. Chceme-li využívat Big Data, musíme škálovat jak infrastrukturu, tak i standardní techniky jejich zpracování. Rychlost může být také problémem, protože hodnota analýzy (a často i dat) se snižuje s časem. Pokud je potřeba více průchodů proudů dat, musí být údaje vloženy do DW, kde lze provést další analýzy. Data mohou být uložena a zpracována pomocí např. NoSQL databáze.

Big Data jsou často zmiňována pouze v souvislosti s BI, nicméně nejen vývojáři BI, ale také vědci v e-science analyzují velké kolekce dat. Výzvou pro počítačové odborníky nebo datové vědce je poskytnout těmto lidem nástroje, které mohou efektivně provádět složitou analytiku s přihlédnutím ke zvláštní povaze zpracování velkých objemů dat. Big Analytics také nezahrnuje pouze fáze analýzy a modelování. Roli hraje zkreslený kontext, heterogenita dat a interpretace výsledků. Tyto aspekty ovlivňují škálovatelné strategie a algoritmy, proto je zapotřebí účinné předzpracování dat (filtrování a integrace) a pokročilé paralelní výpočetní prostředí. Variabilita dat se dnes stává součástí celkového návrhu systému, nicméně výkon je stále požadavkem první kategorie.

Kromě těchto spíše klasických témat DM velkých objemů dat se v posledních letech objevily další zajímavé požadavky, jako rozpoznávání pojmenovaných entit, analýza názorů a mínění (např. pozitivní, negativní, neutrální) a jejich dolování (sentiment analysis). Jejich řešení využívají zejména metody vyhledávání informací a analýzy webových dat. Porovnávání vzorů grafů se běžně používá při analýze sociálních sítí, kde grafy např. zahrnují miliardu uživatelů a stovky miliard odkazů. Technické problémy současných technik DM používaných pro Big Data pak v každém případě pocházejí z jejich nedostatečné škálovatelnosti a paralelizace.

### 3 NoSQL databáze

Pro ukládání a zpracování velkých kolekcí dat jsou často používány NoSQL databáze. NoSQL znamená "ne pouze SQL" nebo "žádné SQL vůbec", což dělá tuto kategorii databází velmi různorodou a ne příliš jasně specifikovatelnou. NoSQL databáze, jejichž vývoj začíná od konce 90. let, poskytují v porovnání s relačními databázemi jednodušší škálovatelnost a vyšší výkon. Popíšeme stručně jejich vlastnosti a klasifikaci včetně použitelnosti pro zpracování Big Dat. Detailnější diskusi těchto témat jsou věnovány v české literatuře např. články [6], [8], či kniha [3], škálovatelnost je diskutována v [7].

To, co je hlavní v klasických přístupech k databázím - (logický) datový model - je v NoSQL databázích popsáno spíše intuitivně, bez jakýchkoliv formálních základů. Terminologie NoSQL je také velmi rozmanitá a rozdíl mezi konceptuálním a databázovým pohledem na data je většinou rozmazaný. Nejznámější NoSQL databáze mohou být podle použitého datového modelu klasifikovány jako:

- úložiště typu klíč-hodnota, např. Redis<sup>1</sup>,
- sloupcově-orientované, např. CASSANDRA<sup>2</sup>,

---

<sup>1</sup> <https://redis.io/>

- dokumentově-orientované, např. MongoDB<sup>3</sup>.

Zdá se, že všechny uvedené datové modely jsou v podstatě typu klíč-hodnota. Odlišují se především v možnostech agregace dvojic (klíč, hodnota) a zpřístupňování těchto hodnot. Obecněji se mezi NoSQL databáze řadí i grafové databáze, XML databáze, RDF databáze a další. Pro naše úvahy vystačíme s třemi výše uvedenými typy.

Tím, že jsou NoSQL určeny hlavně pro ukládání Big Dat, musí být tyto databáze škálovatelné. Významnou roli pak hraje jejich indexace. V NoSQL databázích se používají speciální datové struktury, např. LSM-strom (Log Structured Merge Tree) [5], požívaný např. v Cassandra a MongoDB. LSM-strom je tvořen kaskádou B-stromů. LSM-strom je vhodný speciálně pro data, kde převládá operace INSERT.

NoSQL bývají částí cloudových, datově intenzivních aplikací (hlavně webových). Patří sem zábavné aplikace, obsluha stránek webových míst s vysokým provozem, doručování médií proudovým způsobem, či data vyskytující se v sociálních sítích. Google využívá sloupcově-orientovanou BigTable ve více než 60 aplikacích.

Zkušenosti s NoSQL databázemi ukazují, že je lze použít i na „malá“ data a zejména na aplikace nepožadující transakční sémantiku, např. pro adresáře, blogy nebo systémy zpracování obsahu, rovněž pro Big Analytics i dat v reálném čase (např. proudy kliknutí na webové místo). V prostředí mobilního zpracování dat jsou navíc transakce ve větším rozsahu technicky nemožné. NoSQL systémy jsou tedy vhodné spíše pro prostředí s interaktivními datovými službami. Vynucování schématu a uzamykání na úrovni řádků jako v relačních databázích může překomplikovat tyto aplikace. Absence některých vlastností ACID pak dovoluje význačné zrychlení a decentralizaci NoSQL databází.

Existuje mnoho diskusí o roli NoSQL databází při poskytování informačních služeb. ProNoSQL tábor tvrdí, že tato technologie je budoucností databází. Na druhé straně prorelační databázový tábor tvrdí, že databáze NoSQL mají velkou nevýhodu v tom, že neposkytují korektní zacházení s integritou dat. To souvisí s nedostatkem sémantiky způsobeným jejich základní vlastností – nemají schéma. Nedostatek metadata zabraňuje aplikačnímu systému vědět, která data jsou uložena a jak jsou vzájemně propojena.

V databázovém světě však NoSQL zaujímají významné místo. V hodnotícím DB-Engines Ranking se v květnu 2017 sledovalo 328 systémů. V prvních 10 místech se objevují MongoDB (5. místo), Cassandra (8. místo) a Redis (9. místo).

## 4 Dualita mezi SQL a NoSQL

V práci [4] autoři argumentují, že NoSQL databáze jsou spíše komplementem tradičních transakčních databází. Neměly by se spíše jmenovat „co-relational“<sup>4</sup>? Možná přirozenější je říkat coSQL místo NoSQL. V Tab. 1 je uvedeno devět takových rozdílů.

---

<sup>2</sup> <http://cassandra.apache.org/>

<sup>3</sup> <https://www.mongodb.com/>

<sup>4</sup> Pozor, v češtině pojem „korelační“ znamená něco jiného.

Důležité jsou rozdíly 1 a následně 8. Díky normalizaci mohou být data o jednom objektu v relační databázi rozložena do více relací. Např. data o zákazníkovi jsou v jedné tabulce, data o bankách, kde má zákazník účet, jsou v druhé tabulce. Propojení je realizováno přes cizí klíče. V NoSQL databázi je toto možné realizovat tak, že každý „řádek“ od banky může obsahovat pro každého zákazníka jeho data i čísla účtu. NoSQL jsou denormalizované, tj. ukládají na místě objektu nikoliv objekt, ale kopii objektu. Ten může být dokonce kompozicí řádků (hnízděná data), což v klasickém SQL není možné, není totiž kompozitní. To vede k horším možnostem aktualizace dat.

Fundamentálním rozdílem je nedostatek schémat dat (sémantiky) u NoSQL databáze. Ten brání analytikům porozumění struktuře dat a tím i vytváření seriózních analýz. Tendencí je proto vytvářet víceúrovňové modelovací přístupy zahrnující relační i NoSQL architektury včetně jejich integrace v jedné infrastruktuře. Vytvářejí se tak metody společného návrhu pro relační a NoSQL databáze založené na modifikaci 3-úrovňového ANSI/SPARC přístupu (konceptuální, logický, fyzický návrh) [2]. Z hlediska uložení a přístupu k datům se pak hovoří o *polyglotní perzistenci* či o *polyglotních databázích*. Počítá se i s vývojem konceptuálního/databázového schématu v celkové infrastruktuře. Konceptuální návrh ovšem předpokládá korektnost současné znalosti aplikační domény. Silnou motivací proto je i fakt, že při návrhu databáze je třeba uvažovat vzory pro DM/ML, shlukování některých atributů pro zajištění výkonu systému apod.

V praxi se objevují i další možnosti, jako model konverze schématu, ve kterém se schéma SQL database konvertuje do schématu NoSQL databáze [11].

**Tab. 1:** Duální vlastnosti SQL a NoSQL databázi

	SQL	NoSQL
1	data závislých relací ukazují k rodičům (přes cizí klíče)	od dat rodičů se ukazuje k dětem
2	uzavřený svět	otevřený svět
3	entitty mají identitu (primární klíč)	identitu určuje prostředí
4	data jsou silně typovaná	potenciálně dynamicky typovaná
5	synchrónní (ACID) aktualizace přes více řádků	asynchronní (BASE) aktualizace v jednotlivých hodnotách
6	změny (transakce) koordinuje prostředí	entitty odpovědně reagovat na změny (případná konzistence)
7	referenční integrita založená na hodnotách	slabá referenční integrita založená na výpočtu
8	není kompozitní	je kompozitní
9	optimalizátor dotazů	vývojář/vzor

## 5 Závěry a výzvy

Klíčové problémy pro budování infrastruktury zpracování dat jsou v lidských rozhodnutích týkajících se NoSQL databázi. Zahrnují zejména volbu správných produktů a návrh vhodné databázové architektury pro danou třídu aplikací.

Role člověka je významná i v Big Analytics. Dnes je proces DM řízen analytikem či datovým vědcem. V závislosti na aplikačním scénáři ten určuje část dat, odkud mohou být např. užitečné vzory extrahovány. Lepší řešení by bylo mít k dispozici automatický proces DM s cílem získat přibližné syntetické informace jak o struktuře, tak o obsahu velkého množství dat.

Současné výzvy pro databázový výzkum zahrnují:

- modelování polyglotních databází (relační i NoSQL v jedné infrastruktuře) [1].
- zlepšení kvality a škálovatelnosti metod DM. Formulace dotazu - zejména při absenci schématu - a prezentace a interpretace odpovědí může být netriviální.
- transformaci obsahu do strukturovaného formátu pro pozdější analýzu, protože mnoho dat není nativně strukturovaných. Filtrací lze i zmenšit objem dat.
- vývoj smysluplného a použitelného formalismus pro modelování NoSQL databází a následně silný a použitelný uživatelský dotazovací jazyk.

Vztah SQL a coSQL databází lze charakterizovat v pojmech jin a jang [4]. V čínské filosofii jde o dva související a protikladné pojmy, pomocí nichž se dá popsat vzájemný poměr rovnováhy jak v těle, tak i v dějích okolo nás (např. noc a den). Rovněž coSQL a SQL nejsou v konfliktu. Jsou to dva protiklady, které koexistují v harmonii, vzájemně se doplňují, podporují a mohou se vzájemně měnit. Kéž by tomu tak bylo i v praxi.

## **Literatura**

1. Abelló, A.: Big Data Design. In: Proc. of DOLAP (2015) 35-38.
2. Herrero, V., Abelló, A., Romero, O.: NOSQL Design for Analytical Workloads: Variability Matters. Proc. of ER Conf. (2016) 50-64.
3. Holubová, I., Kosek, J., Minařík, K., Novák, D.: Big Data a NoSQL databáze. Grada, 2015.
4. Meijer, E., Bierman, G.M.: A co-relational model of data for large shared data banks. Commun. ACM 54(4) (2011) 49-58.
5. O'Neil, P. E., Cheng, E., Gawlick, D., O'Neil, E.: The Log-Structured Merge-Tree (LSM-Tree). Acta Inf., 33(4) (1996) 351-385.
6. Pokorný, J.: NoSQL databáze. In: Proc. of the Annual Database Conf. DATAKON'2011, J. Zendulka, M. Rychlý (eds.), Mikulov, VUT Brno (2011) 71-82.
7. Pokorný J.: NoSQL Databases: a step to databases scalability in Web environment. International Journal of Web Information Systems, 9 (1) (2013) 69-82.
8. Pokorný, J.: Big Data: jejich ukládání, zpracování a použití. Proc. of the 34th Ann. Database Conference DATAKON'2014, P. Šaloun, D. Chlapek (Eds.), VŠB Ostrava (2014) 3-16.
9. Pokorný, J.: Big Data Storage and Management: Challenges and Opportunities. In: Proc. of 12th IFIP WG 5.11 Int. Symp. on Environmental Software Systems, IFIP AICT 507, Springer (2017)
10. Stonebraker, M., Deng, D., Brodie, M.L.: Database decay and how to avoid it. In: Proc. of 2016 IEEE Int. Conference on Big Data, IEEE Explore (2016) 7-16.
11. Zhao, G., Lin, Q., Li, L., Li, Z.: Schema Conversion Model of SQL Database to NoSQL. In: Proc. of the 3PGCIC, IEEE (2014) 355-362.

*Relační a NoSQL databáze: dvě strany téže mince?*

**Poděkování:** Práce byla podpořena projektem Q48 programu Progres na UK, Praha.

**Annotation:**

*Relational and NoSQL databases: two sides of the same coin?*

The analysis of relational and NoSQL databases leads to the conclusion that these data processing systems are to some extent complementary. In current Big Data applications, especially where extensive analyses are needed, it turns out that it is non-trivial to design an infrastructure involving software of both types. In terms of performance, it may even be beneficial to transform the SQL database schema into NoSQL or to perform double-sided data migration between relational and NoSQL databases. The aim of the article is to discuss these possibilities and some new methods of designing such database architectures standing on the legacy of the three-level ANSI/SPARC architecture.