




Výkonnostní archeologie

Tomáš Vondra, GoodData

tomas.vondra@gooddata.com / tomas@pgaddict.com
@fuzzycz, <http://blog.pgaddict.com>



PostgreSQL

A Practical Guide to the Advanced Open Source Database

PostgreSQL
Up and Running
O'REILLY
Regina Obe & Leo Hsu

Booking.com

OpenFest2013

O'REILLY

Graph Databases
Elasticsearch
Python and Django
The Way to Web 2.0

Jak se změnil výkon PostgreSQL za
několik posledních verzí?

7.4 vyšla 2003, tj. cca 10 let

(překvapivě) ošidná otázka

- během vývoje se většinou dělají “parciální” testy
 - srovnání dvou verzí / commitů
 - zaměřené na konkrétní část kódu / vlastnost
- komplexnější benchmarky pro srovnání dvou verzí
 - obtížné “zkombinovat” (různý hardware, ...)
- aplikační výkon (ultimátní benchmark)
 - podléhá (pravidelným) upgradům hardwaru
 - růst datových objemů, evoluce aplikace (nové fičury)

(poněkud) nefér otázka

- vývoj probíhá oproti aktuálně dostupnému hardwaru
 - Kolik RAM jste měli v serveru před 10 lety?
 - Kdo z vás měl před 10 lety SSD/NVRAM disky?
 - Kdo z nás měl stroje s 8 CPU?
- některé rozdíly jsou důsledkem těchto změn (algoritmy)

Každopádně vyšší výkon na aktuálním
hardwaru je fajn ;-)

Pojďme si zabenchmarkovat!

Pokud se bojíte čísel nebo grafů, asi byste
radši měli odejít hned.

Bude tu spousta obojího ...

<http://blog.pgaddict.com>

<http://planet.postgresql.org>

82,71% statistik na internetu je
vycucaných z prstu ...

... přísahám že ty moje to nejsou!

Benchmarky (přehled)

- **pgbench (TPC-B)**
 - “transakční” benchmark
 - operace pracují s malými počty řádek (přístup přes primární klíče)
- **TPC-DS (náhrada TPC-H)**
 - “warehouse” benchmark
 - dotazy drtící spousty dat (agregace, joiny, ROLLUP/CUBE, ...)
- **fulltext benchmark (tsearch2)**
 - primárně o vylepšeních GIN/GiST indexů
 - platí pro další aplikace používající GIN/GiST (geo, ...)

Použitý hardware

HP DL380 G5 (2007-2009)

- 2x Xeon E5450 (each 4 cores @ 3GHz, 12MB cache)
- 16GB RAM (FB-DIMM DDR2 667 MHz), FSB 1333 MHz
- 6x10k RAID10 (SAS) @ P400 with 512MB write cache
- S3700 100GB (SSD)

Workstation i5 (2011-2013)

- 1x i5-2500k (4 cores @ 3.3 GHz, turbo 3.9 GHz, 6MB cache)
- 8GB RAM (DIMM DDR3 1333 MHz)
- S3700 100GB (SSD)

pgbench

TPC-B “transakční” benchmark

pgbench

- tři velikosti datasetů
 - malý (150 MB)
 - střední (~50% RAM)
 - velký (~200% RAM)
- dva módy
 - read-only a read-write
- rozsah klientů (1, 2, 4, ..., 32)

pgbench

- tři velikosti datasetů
 - malý (150 MB) ← problémy se zámky, etc.
 - střední (~50% RAM) ← CPU bound
 - velký (~200% RAM) ← I/O bound
- dva módy
 - read-only a read-write
- rozsah klientů (1, 2, 4, ..., 32)

```
BEGIN;
```

```
    UPDATE accounts SET abalance = abalance + :delta  
    WHERE aid = :aid;
```

```
    SELECT abalance FROM accounts WHERE aid = :aid;
```

```
    UPDATE tellers SET tbalance = tbalance + :delta  
    WHERE tid = :tid;
```

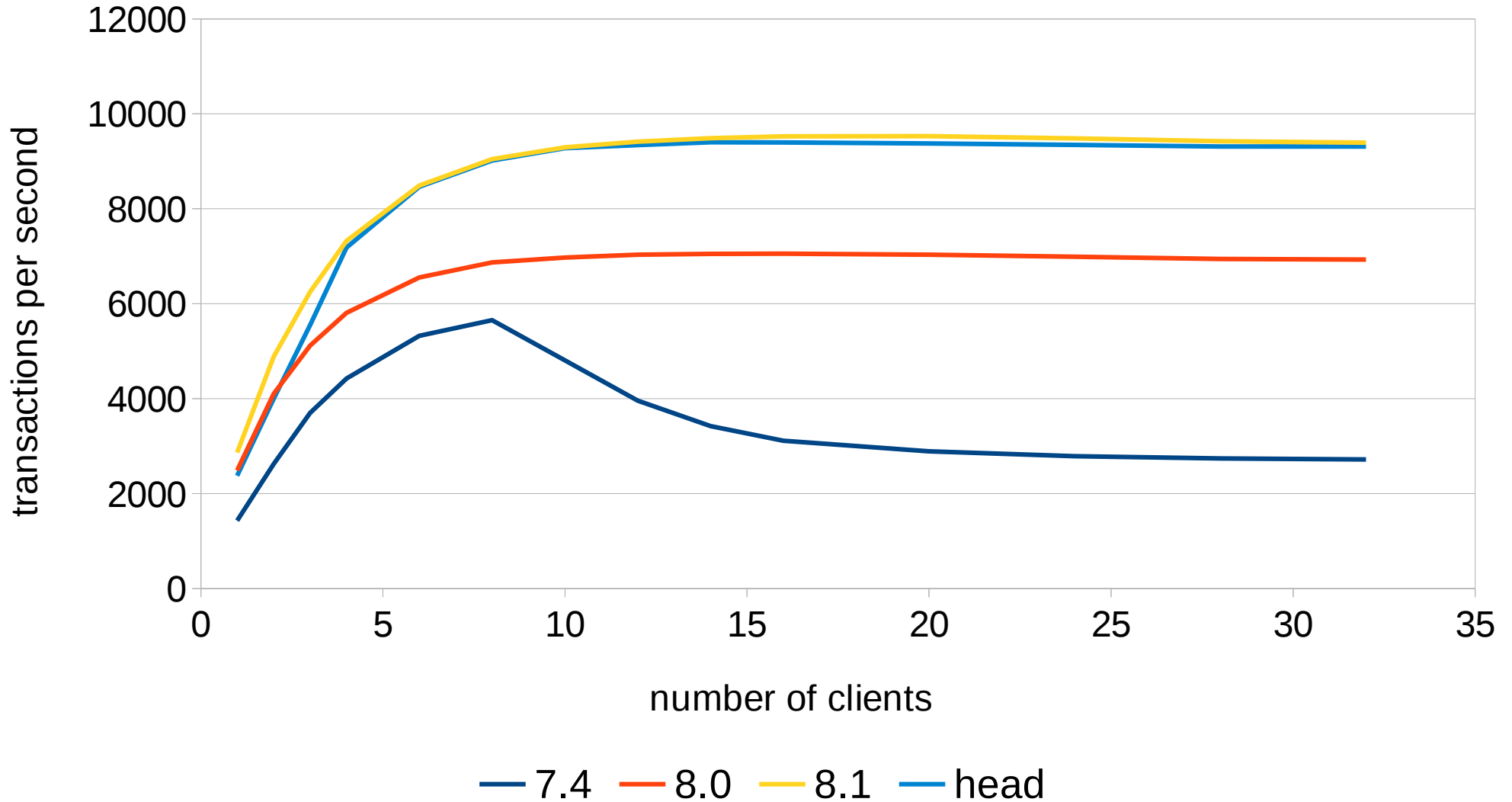
```
    UPDATE branches SET bbalance = bbalance + :delta  
    WHERE bid = :bid;
```

```
    INSERT INTO history (tid, bid, aid, delta, mtime)  
    VALUES (:tid, :bid, :aid, :delta, CURRENT_TIMESTAMP);
```

```
END;
```

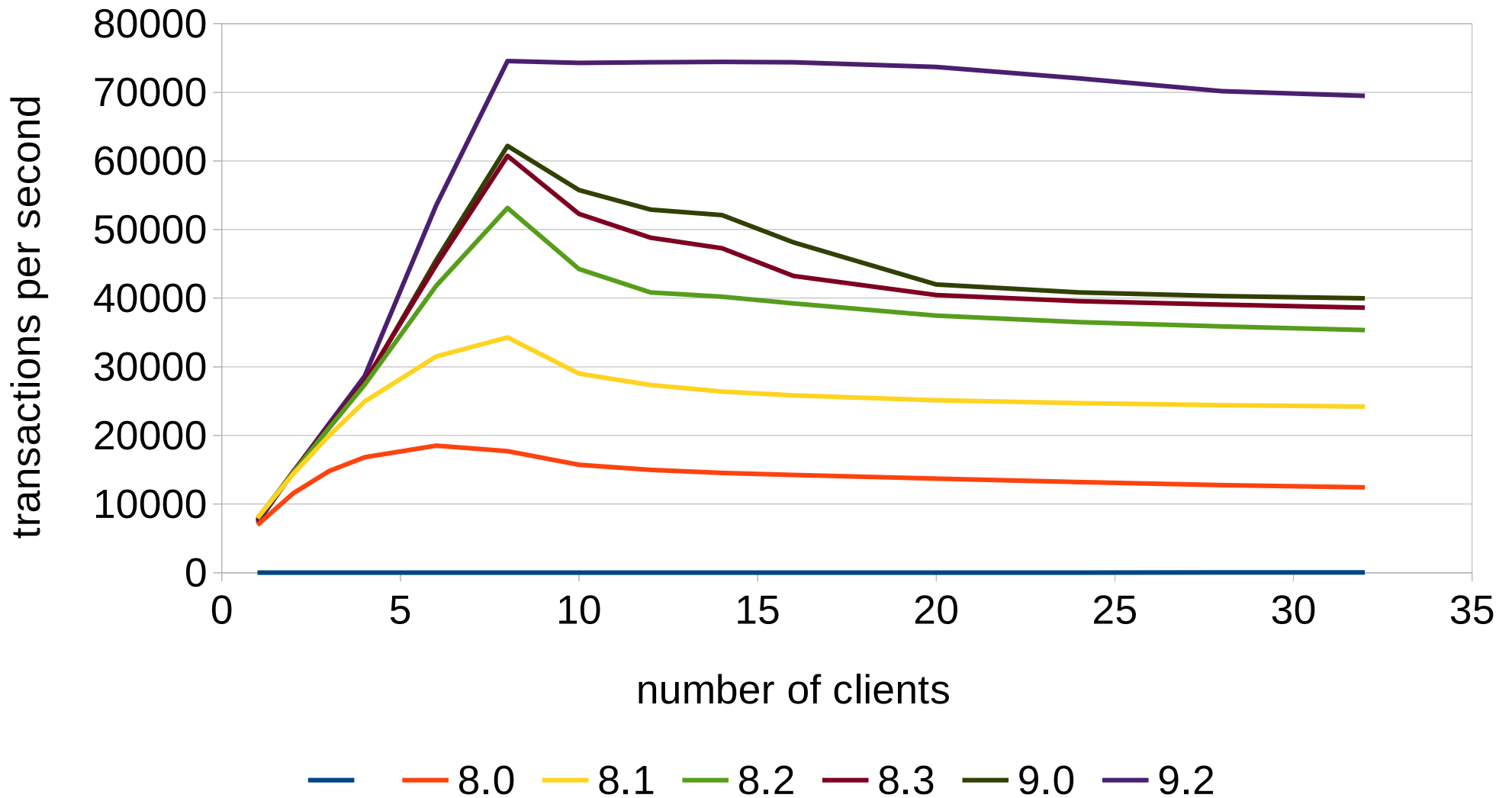

pgbench / large read-only (on SSD)

HP DL380 G5 (2x Xeon E5450, 16 GB DDR2 RAM), Intel S3700 100GB SSD



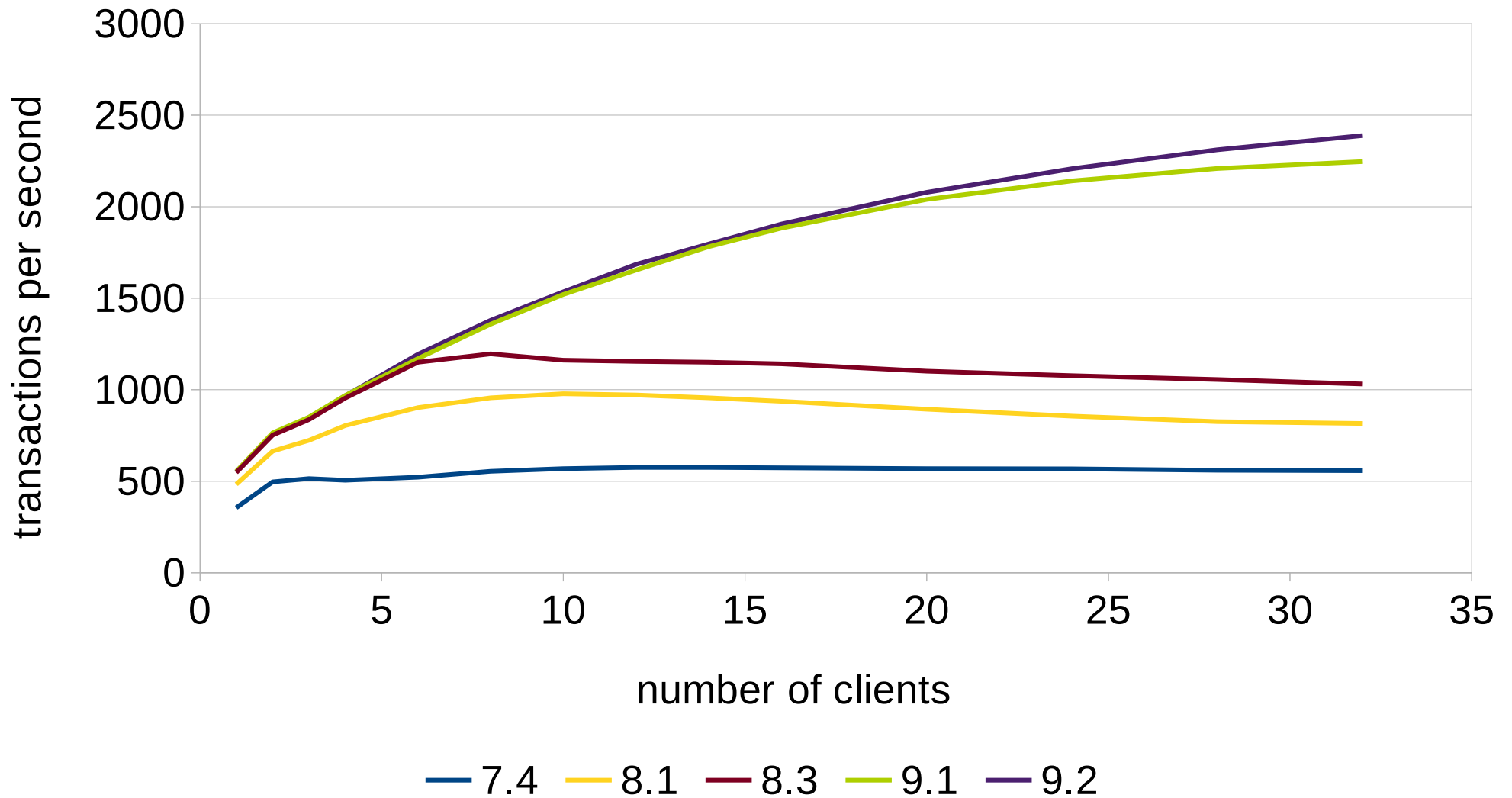
pgbench / medium read-only (SSD)

HP DL380 G5 (2x Xeon E5450, 16 GB DDR2 RAM), Intel S3700 100GB SSD



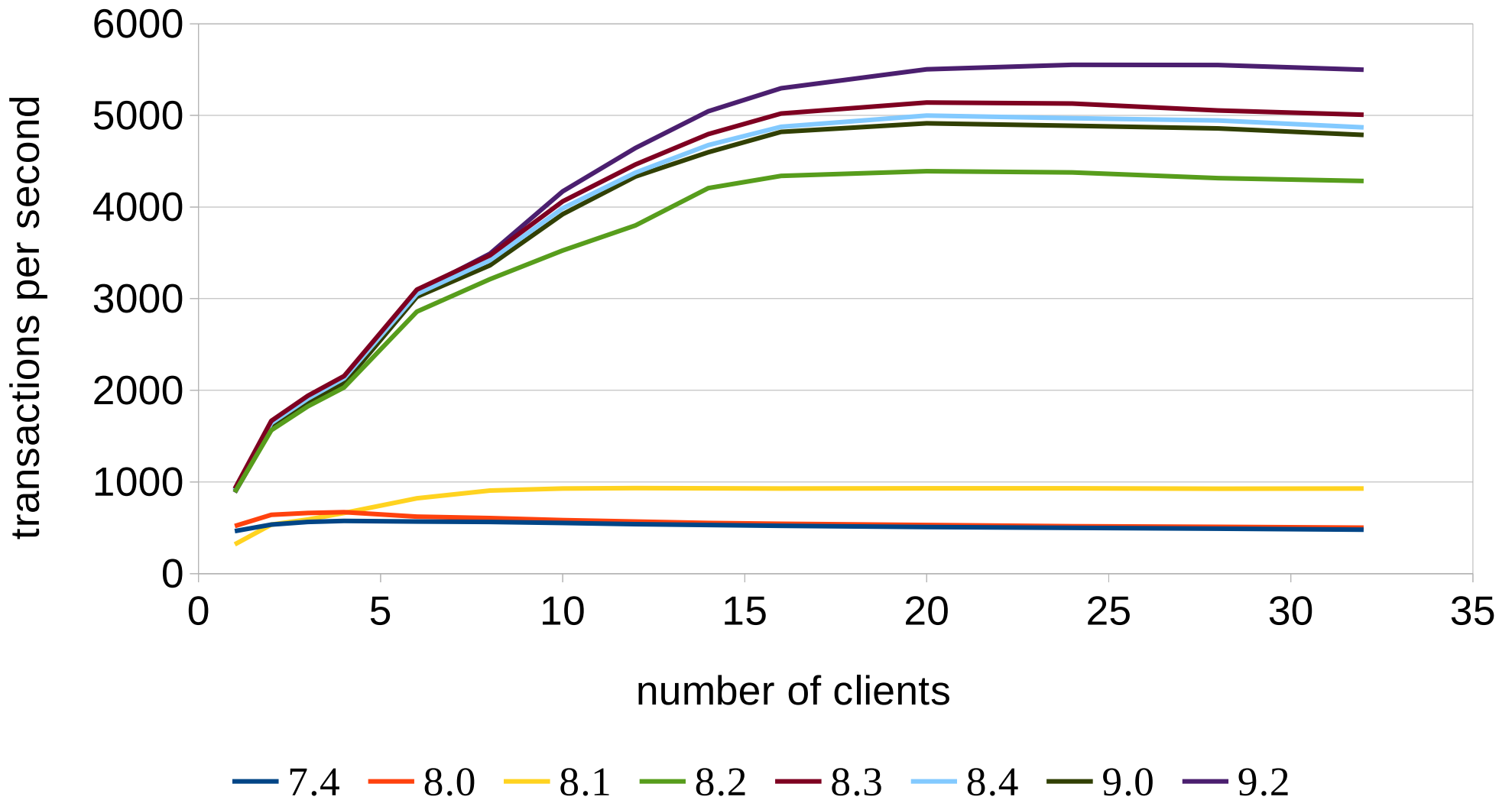
pgbench / large read-write (SSD)

HP DL380 G5 (2x Xeon E5450, 16 GB DDR2 RAM), Intel S3700 100GB SSD



pgbench / small read-write (SSD)

HP DL380 G5 (2x Xeon E5450, 16 GB DDR2 RAM), Intel S3700 100GB SSD



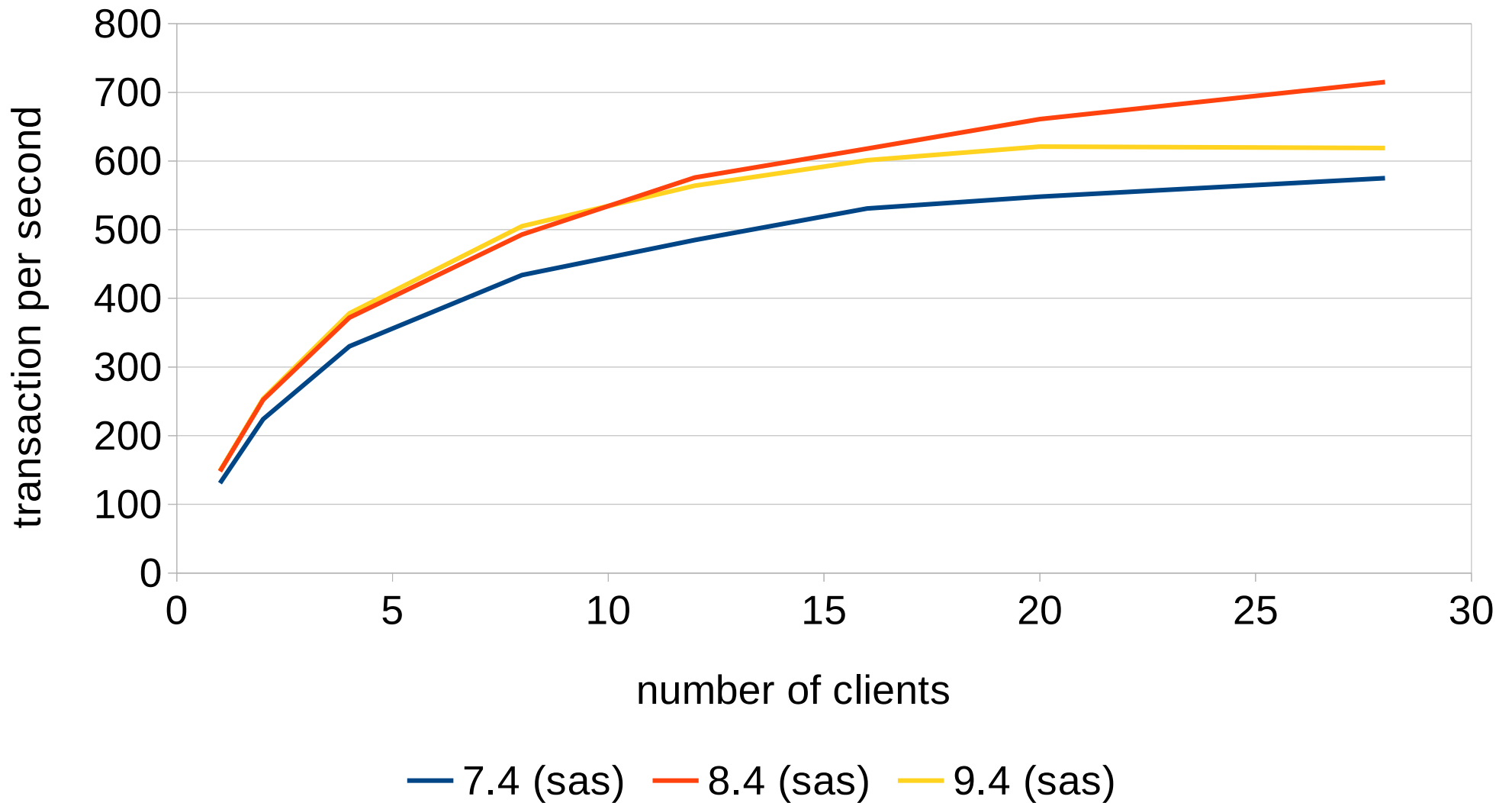
Co rotační disky?

6 x 10k SAS drives (RAID 10)

P400 with 512MB write cache

pgbench / large read-write (SAS)

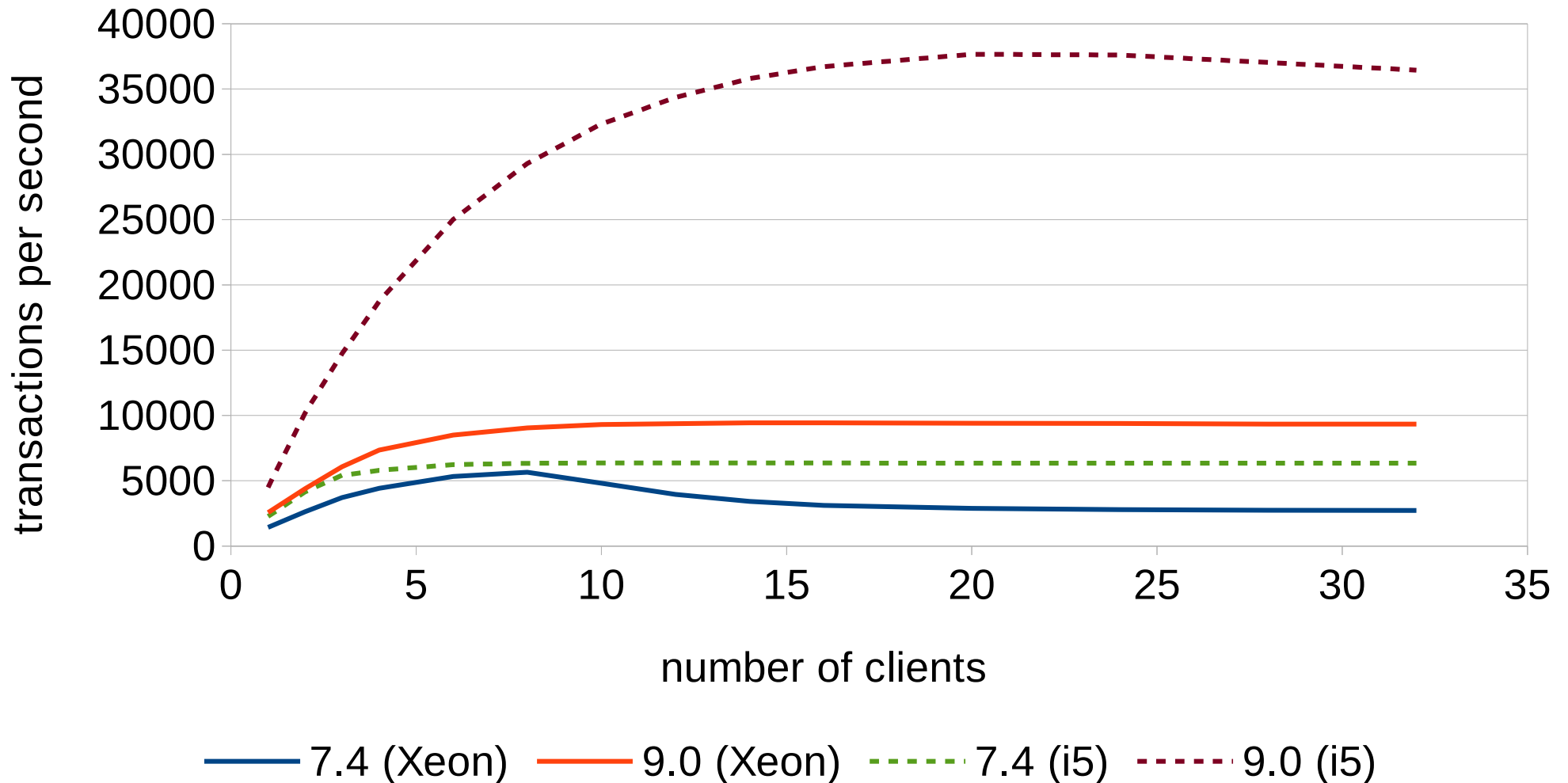
HP DL380 G5 (2x Xeon E5450, 16 GB DDR2 RAM), 6x 10k SAS RAID10



No a co ta i5-2500k mašina?

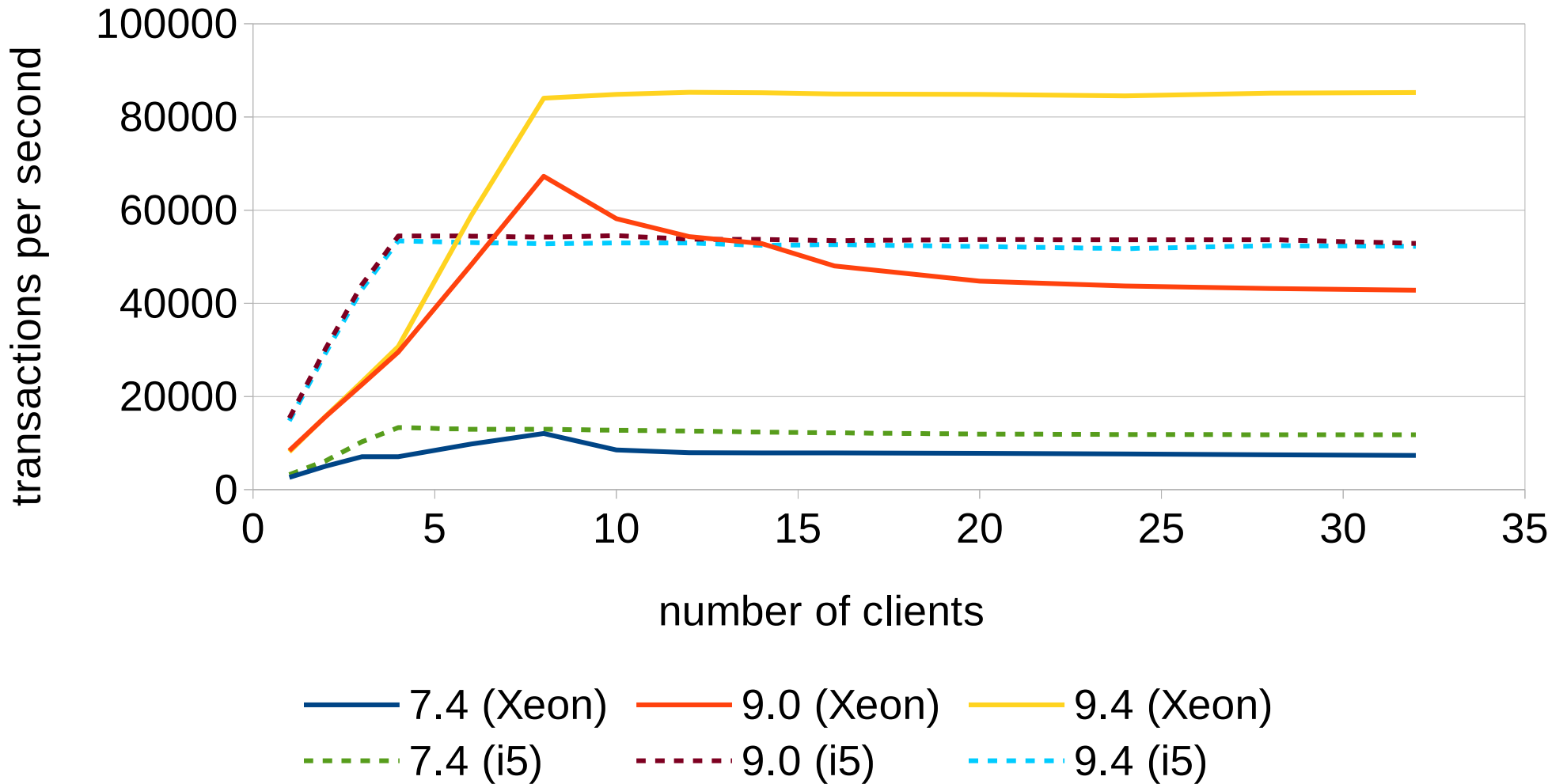
pgbench / large read-only (Xeon vs. i5)

2x Xeon E5450 (3GHz), 16 GB DDR2 RAM, Intel S3700 100GB SSD
i5-2500k (3.3 GHz), 8GB DDR3 RAM, Intel S3700 100GB SSD



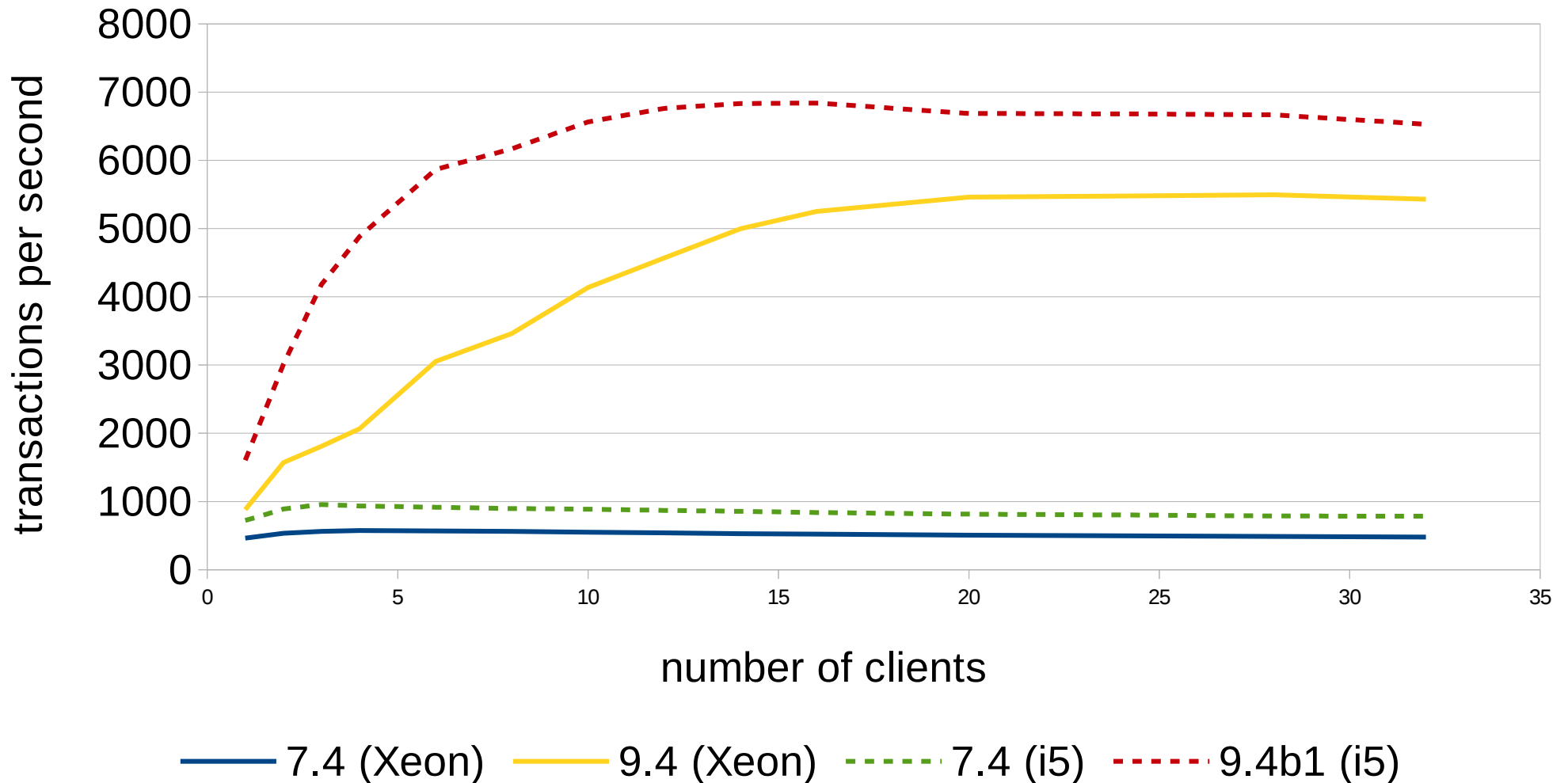
pgbench / small read-only (Xeon vs. i5)

2x Xeon E5450 (3GHz), 16 GB DDR2 RAM, Intel S3700 100GB SSD
i5-2500k (3.3 GHz), 8GB DDR3 RAM, Intel S3700 100GB SSD



pgbench / small read-write (Xeon vs. i5)

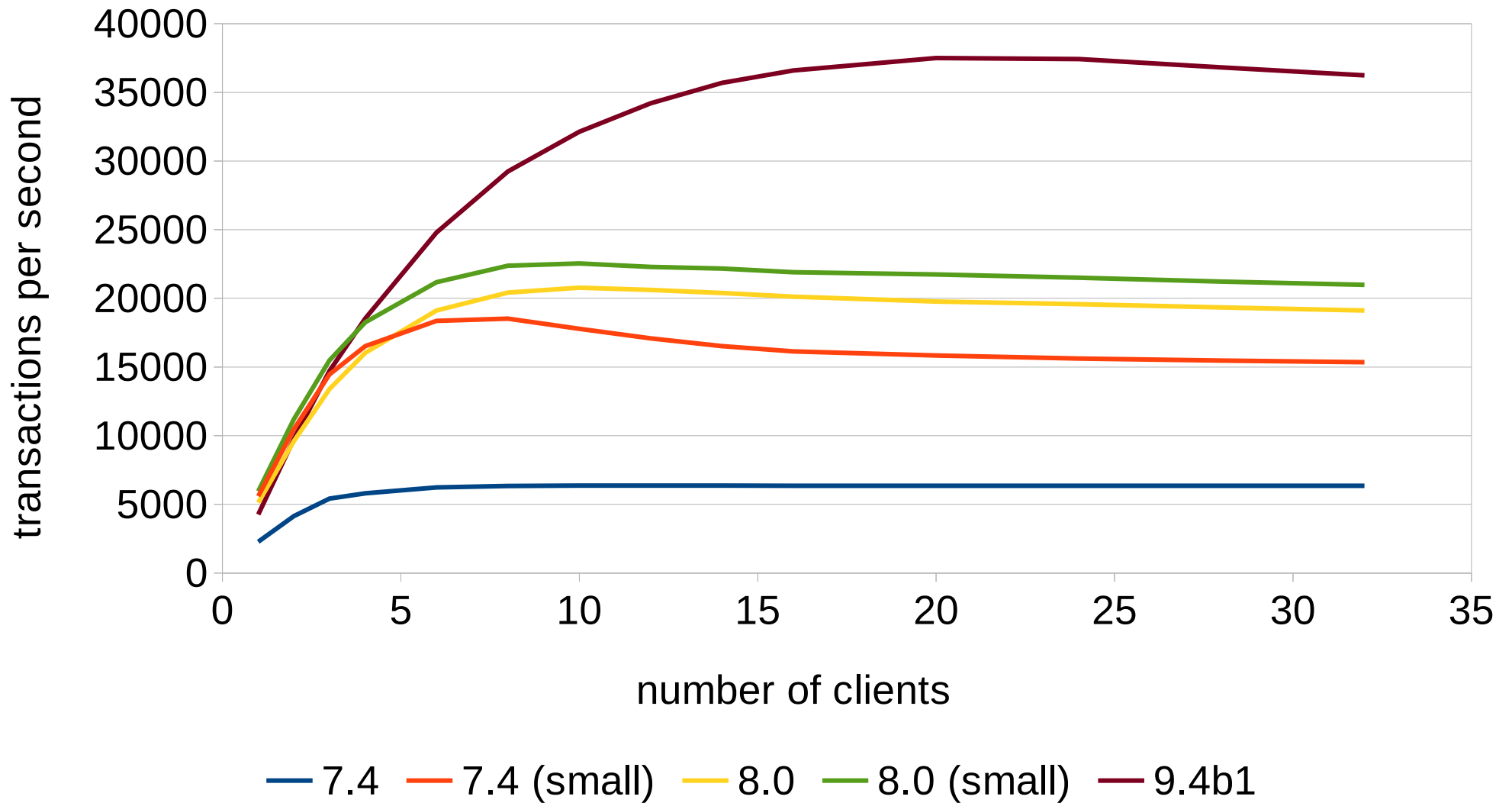
2x Xeon E5450 (3GHz), 16 GB DDR2 RAM, Intel S3700 100GB SSD
i5-2500k (3.3 GHz), 8GB DDR3 RAM, Intel S3700 100GB SSD



Legendy říkají že starší verze lépe fungují s nižšími paměťovými limity
(shared_buffers etc.)

pgbench / large read-only (i5-2500)

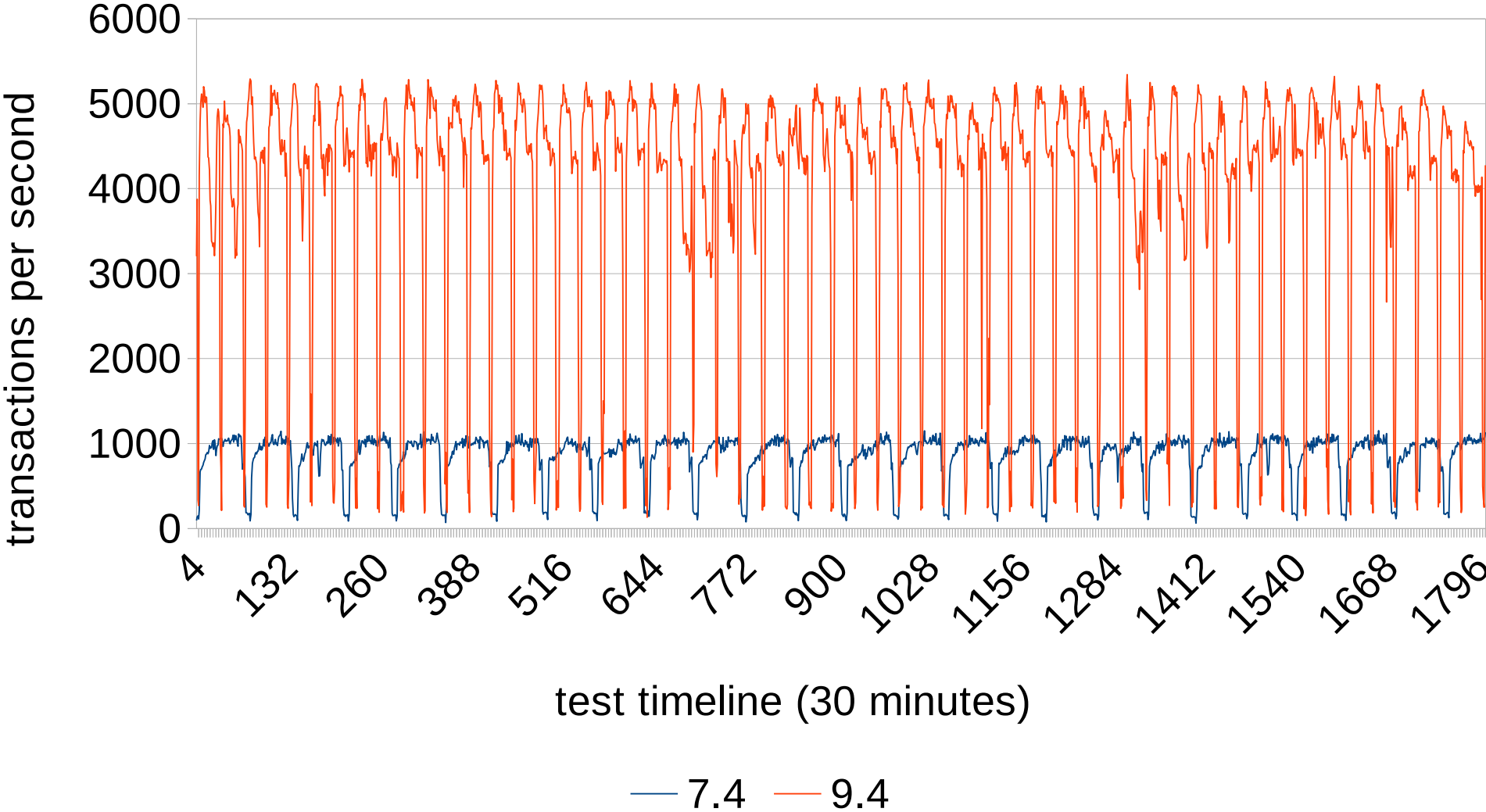
different sizes of shared_buffers (128MB vs. 2GB)



transactions per second vs. latence

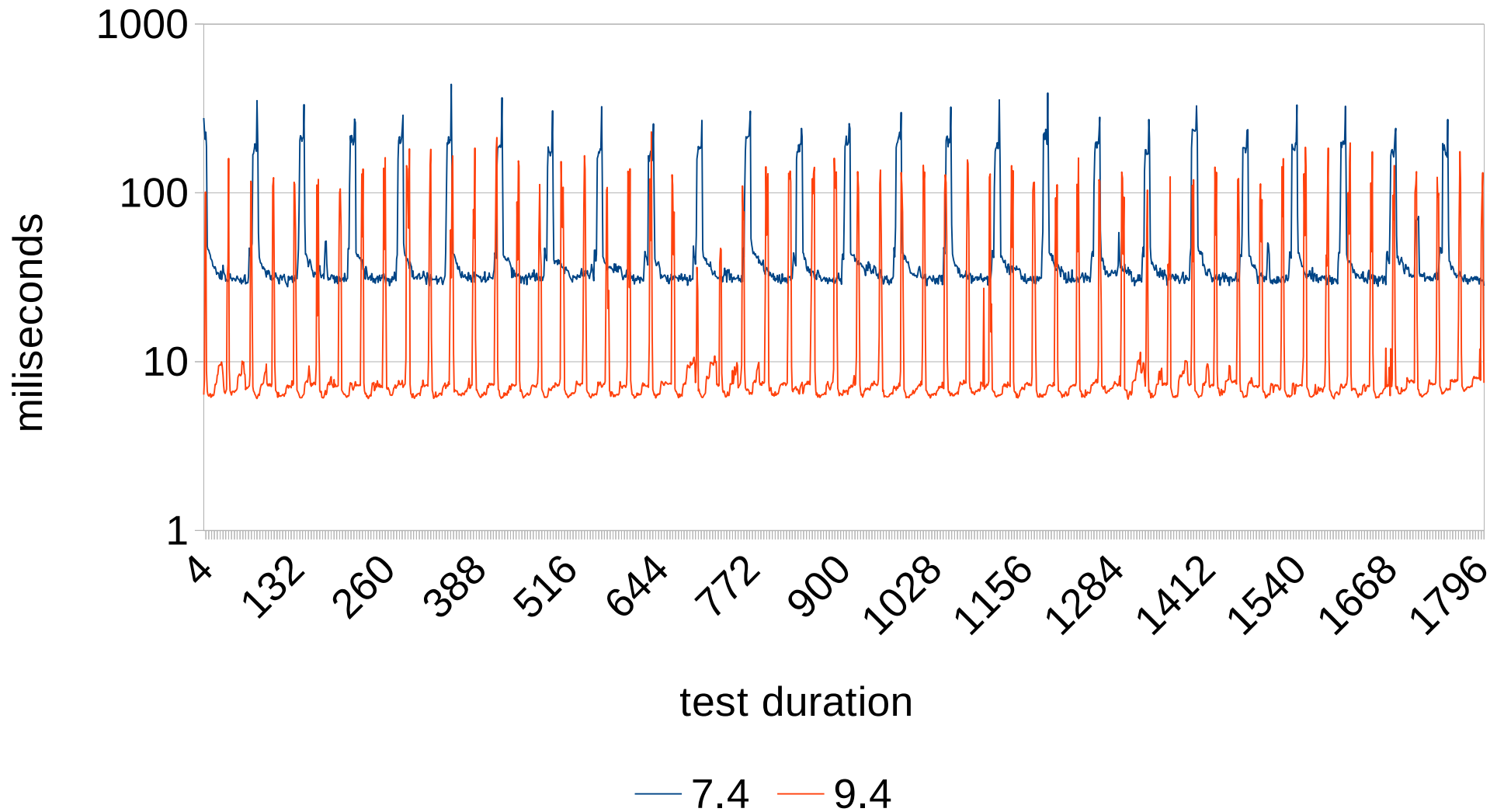
pgbench transaction rate / 7.4 vs. 9.4

transactions per second / large read-write dataset (32 clients)



pgbench transaction latency / 7.4 vs. 9.4

average latency (milliseconds) / large read-write



pgbench / shrnutí

- značná vylepšení
- vylepšené zamykání
 - lepší škálování na velké počty jader (64 ...)
- mnoho dalších optimalizací
 - výrazné zrychlení i na malém počtu klientů
- lessons learned
 - frekvence procesoru není míra výkonu
 - počet jader není míra výkonu

TPC-DS

“Decision Support” benchmark
(aka “Data Warehouse” benchmark)

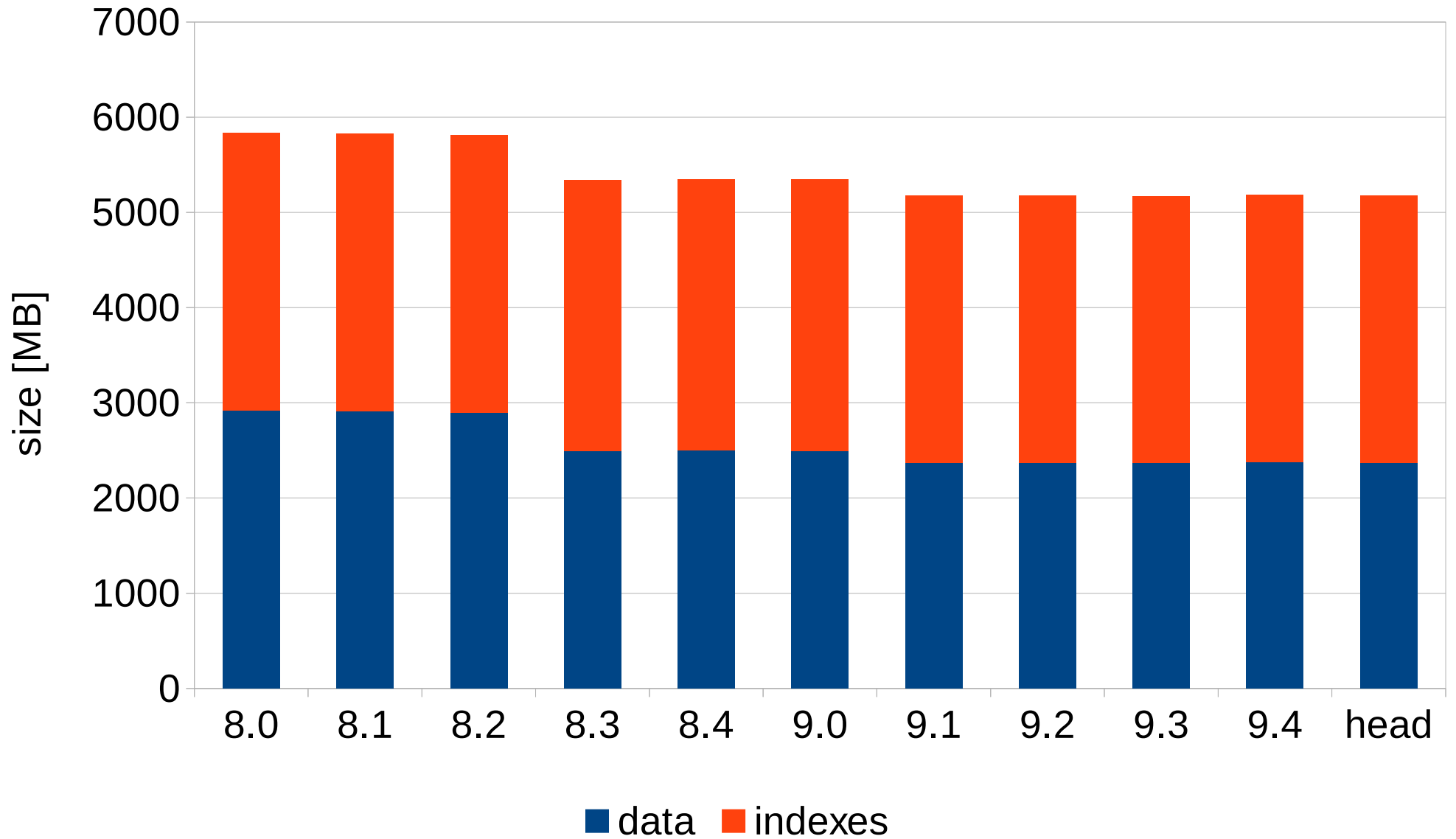
TPC-DS

- orientováno na analytiku / warehousing
 - dotazy drtící velké objemy dat (GROUP BY, JOIN)
 - neuniformní rozdělení dat (realističtější než TPC-H)
- definováno 99 šablon dotazů (TPC-H jen 22)
 - některé rozbité (padá generátor)
 - některé zatím nepodporované (ROLLUP/CUBE)
 - 41 dotazů pro ≥ 7.4
 - 61 dotazů pro ≥ 8.4 (CTE, Window functions)

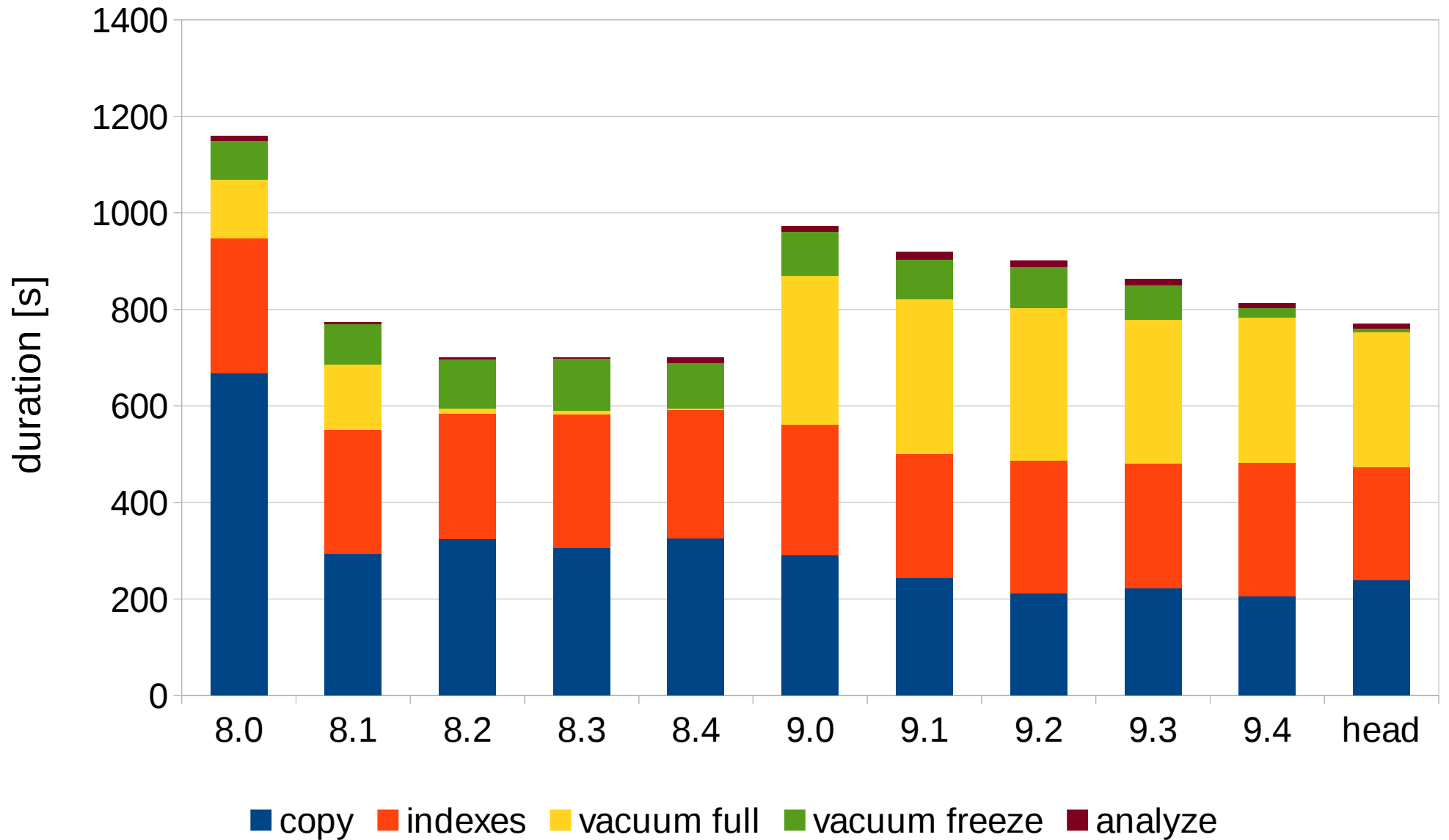
TPC-DS

- **1GB and 16GB datasets (raw data)**
 - 1GB nedostatečný pro publikaci, 16GB je nestandardní (dle TPC)
- **ale i tak je to zajímavé ...**
 - většina produkčních databází se vejde do 16GB
 - ukazuje to trendy (do jisté míry aplikovatelné na větší DB)
- **schéma**
 - víceméně defaultní (standard compliance FTW!)
 - stejné pro všechny verze (indexy na FK/join keys, pár dalších)
 - nepochybně možno dál zoptimalizovat

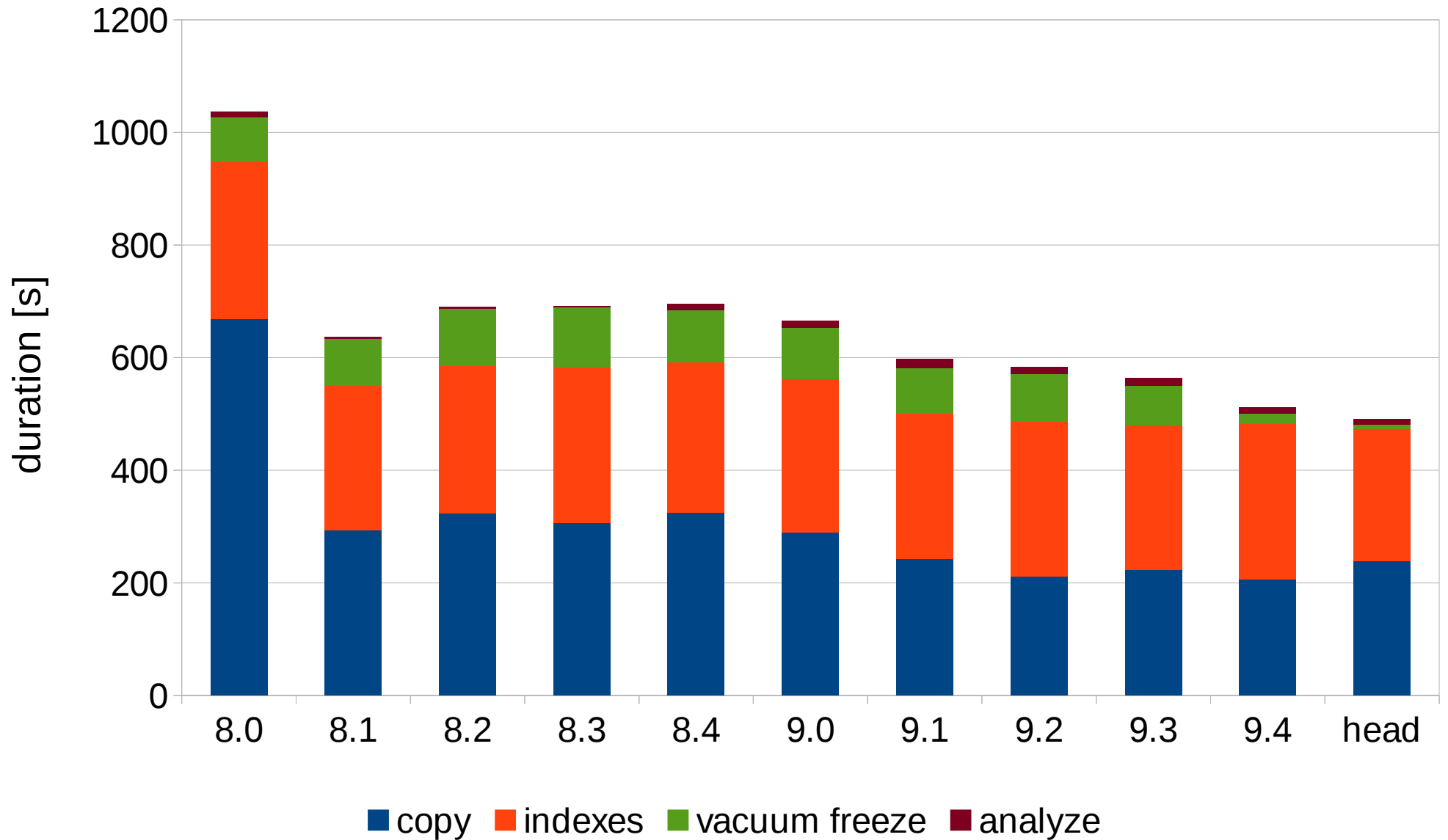
TPC DS / database size per 1GB raw data



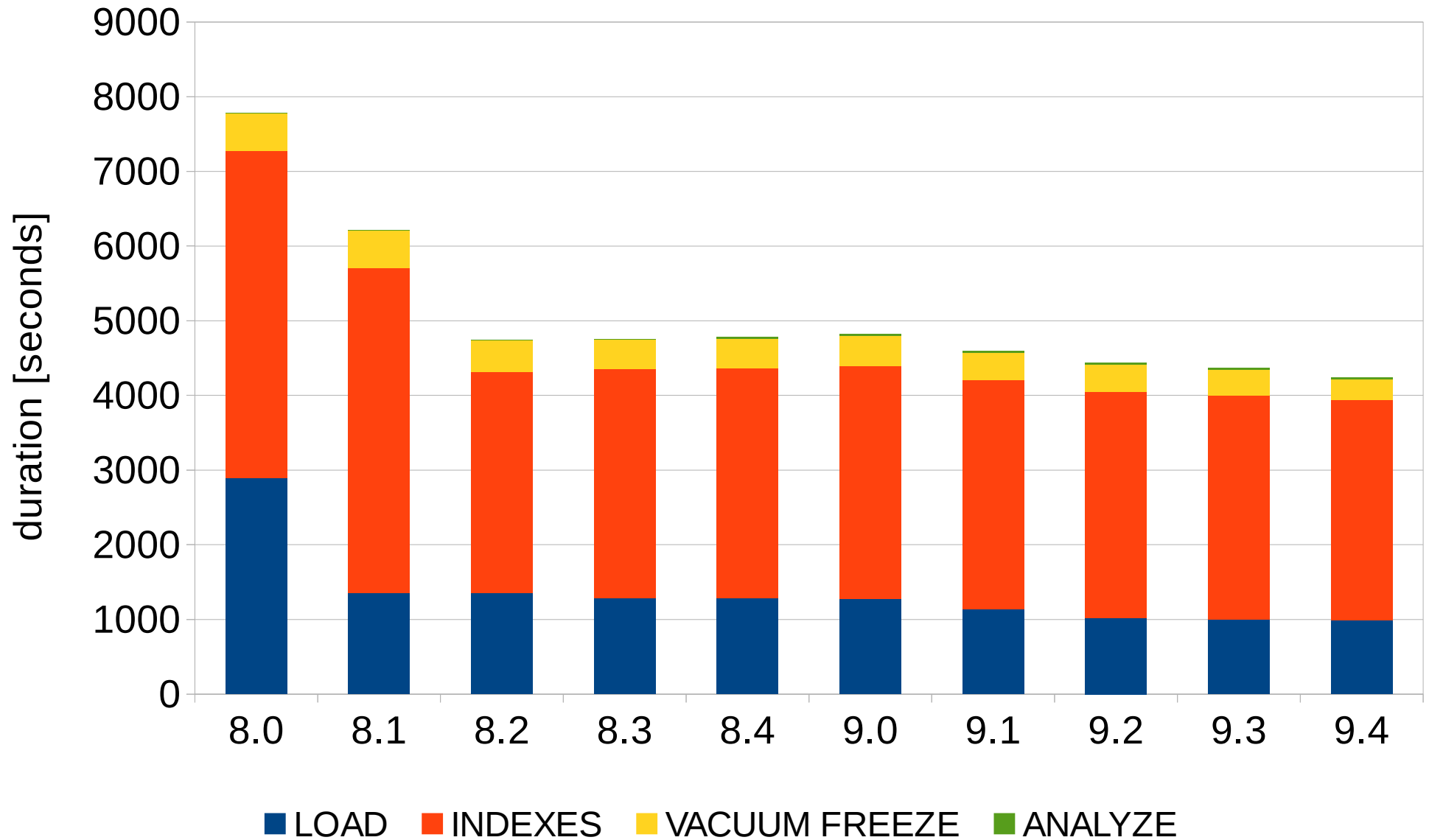
TPC DS / load duration (1GB)



TPC DS / load duration (1GB)

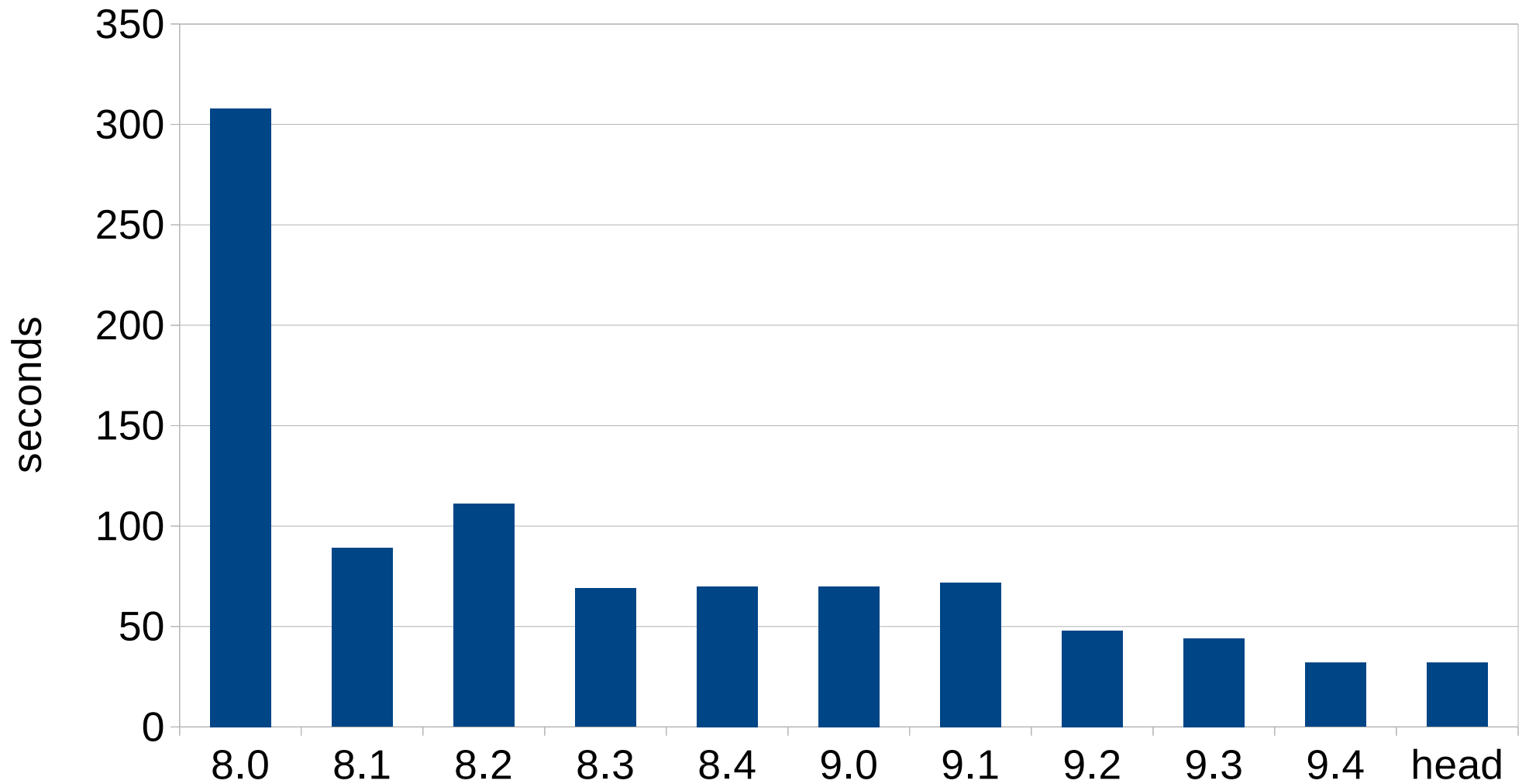


TPC DS / load duration (16 GB)



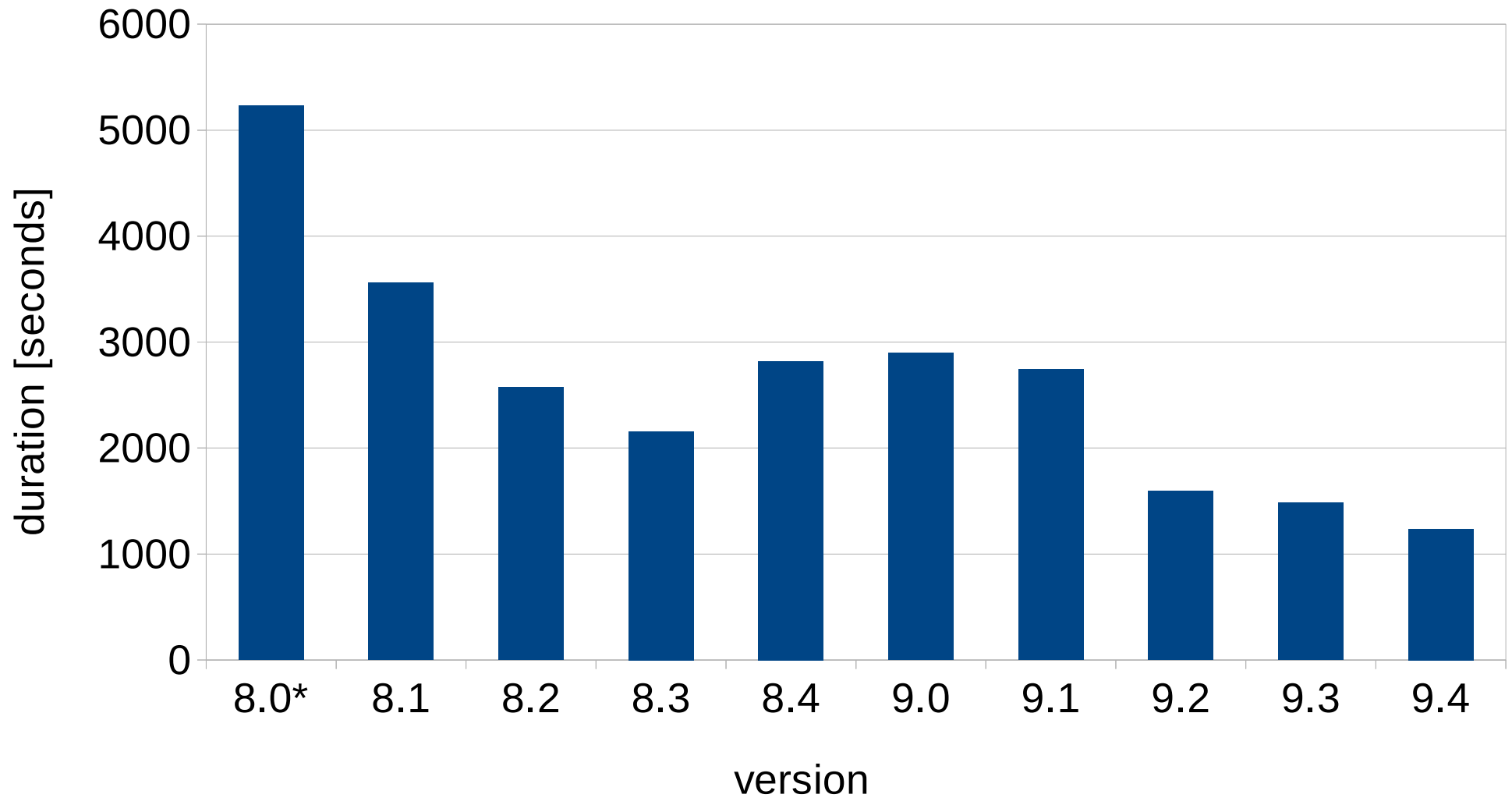
TPC DS / duration (1GB)

average duration of 41 queries



TPC DS / duration (16 GB)

average duration of 41 queries



TPC-DS / shrnutí

- **výrazně rychlejší load dat**

- většina času se vytváří indexy (paralelizace, RAM)
- pokud ignorujeme VACUUM FULL (změna implementace v 9.0)
- mírné snížení velikosti

- **výrazně rychlejší dotazy**

- značné zrychlení dotazů (~6x)
- širší použití indexů, index only scany

Fulltext Benchmark

testování GIN a GiST indexů
prostřednictvím fulltextu

Fulltext benchmark

- prohledávání archivů pgsql mailing listů
 - ~1M zpráv, ~5GB dat
- ~33k reálných dotazů (z postgresql.org)
 - syntetické dotazy dávají cca stejné výsledky

```
SELECT id FROM messages
```

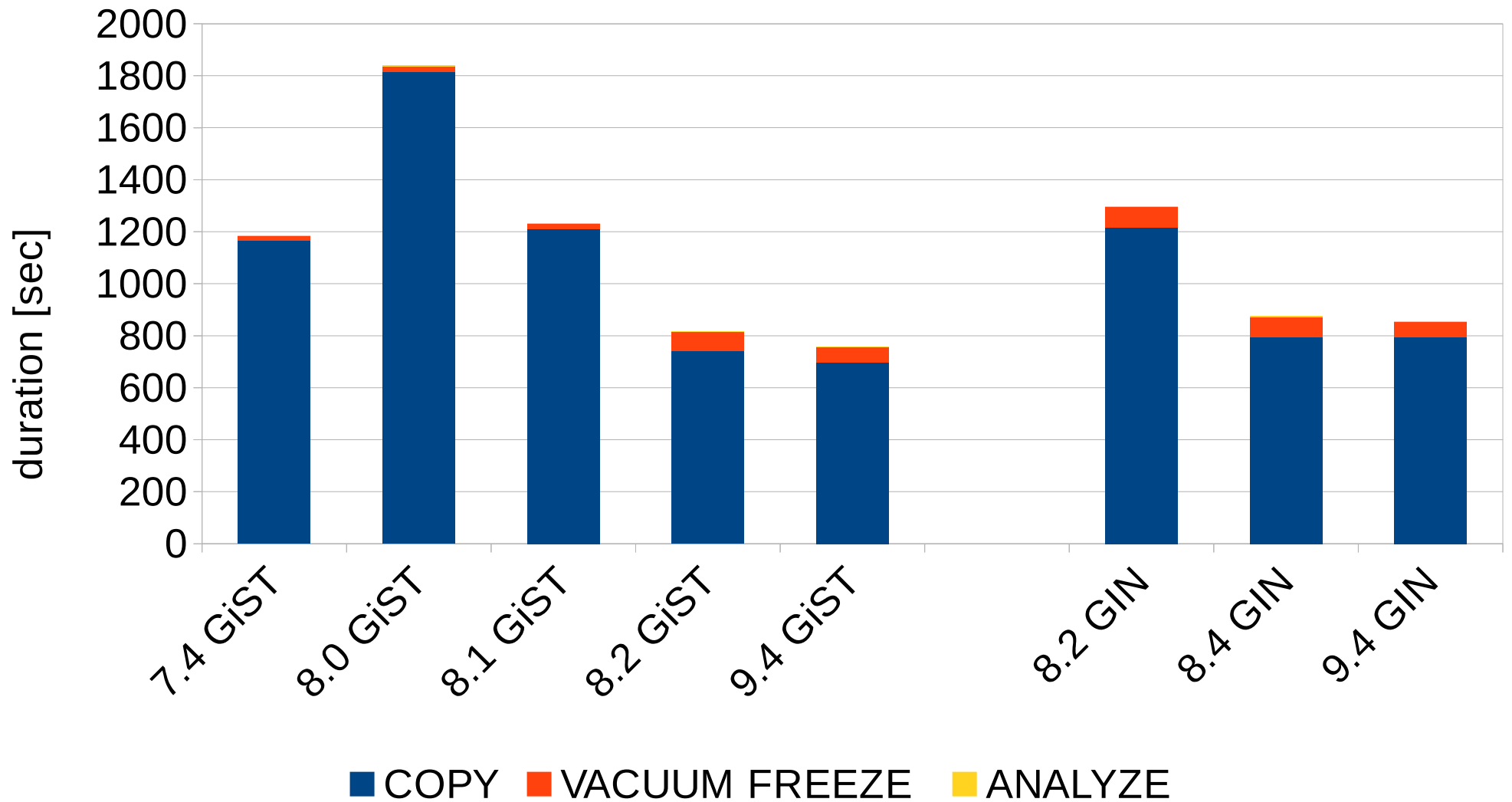
```
WHERE body @@ ('high & performance')::tsquery
```

```
ORDER BY ts_rank(body, ('high & performance')::tsquery)
```

```
DESC LIMIT 100;
```

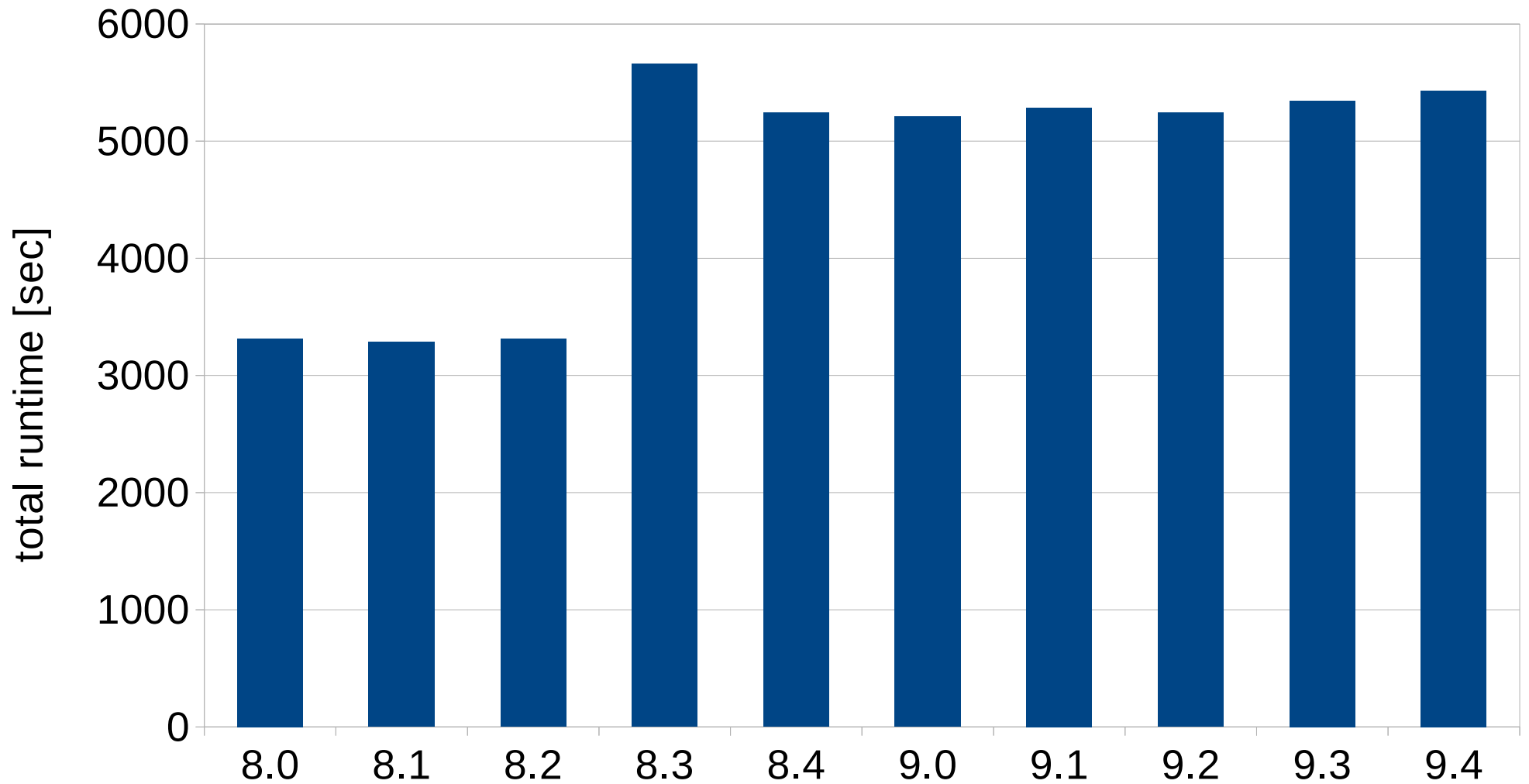
Fulltext benchmark / load

COPY / with indexes and PL/pgSQL triggers



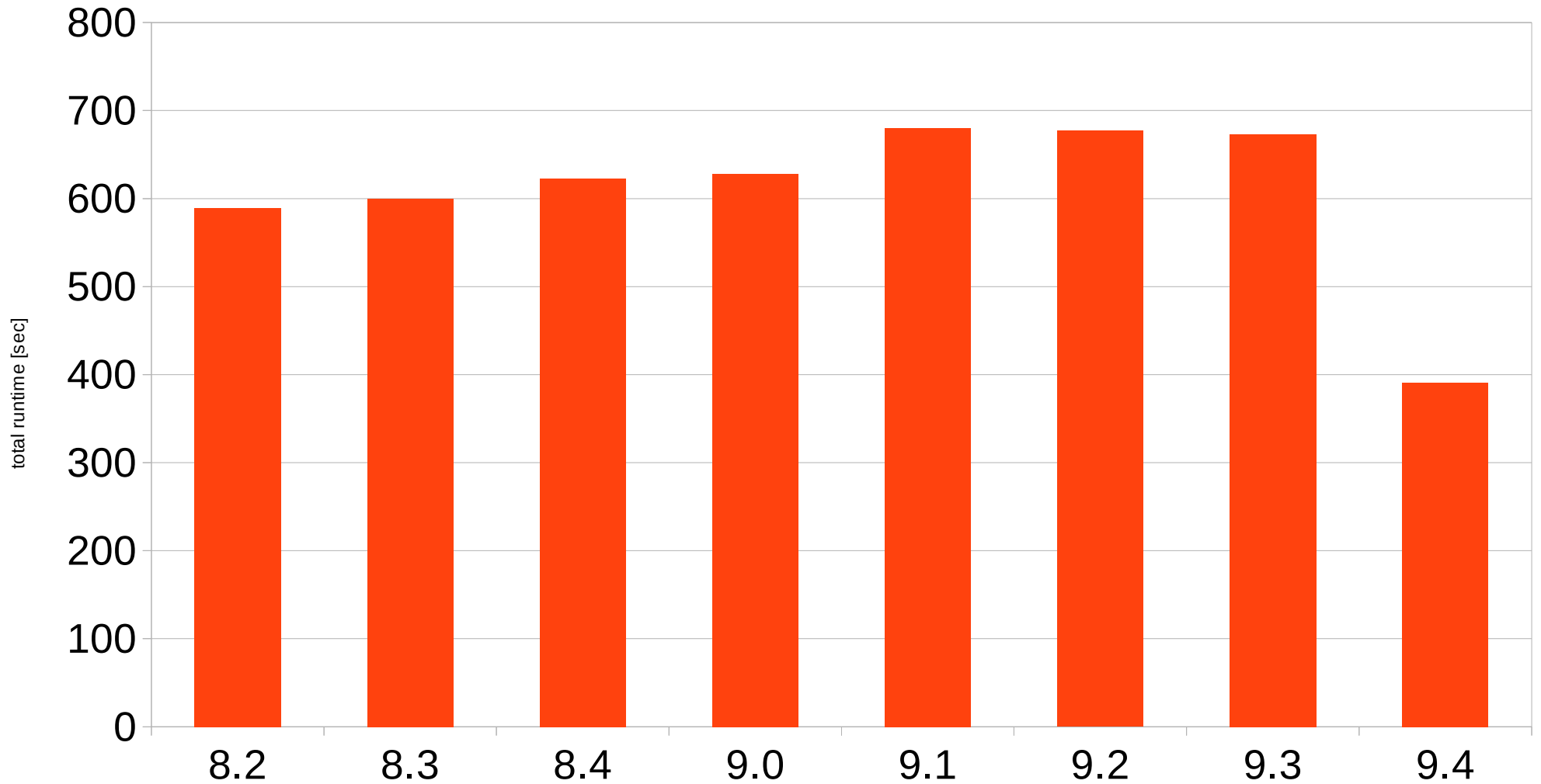
Fulltext benchmark / GiST

33k queries from postgresql.org [TOP 100]



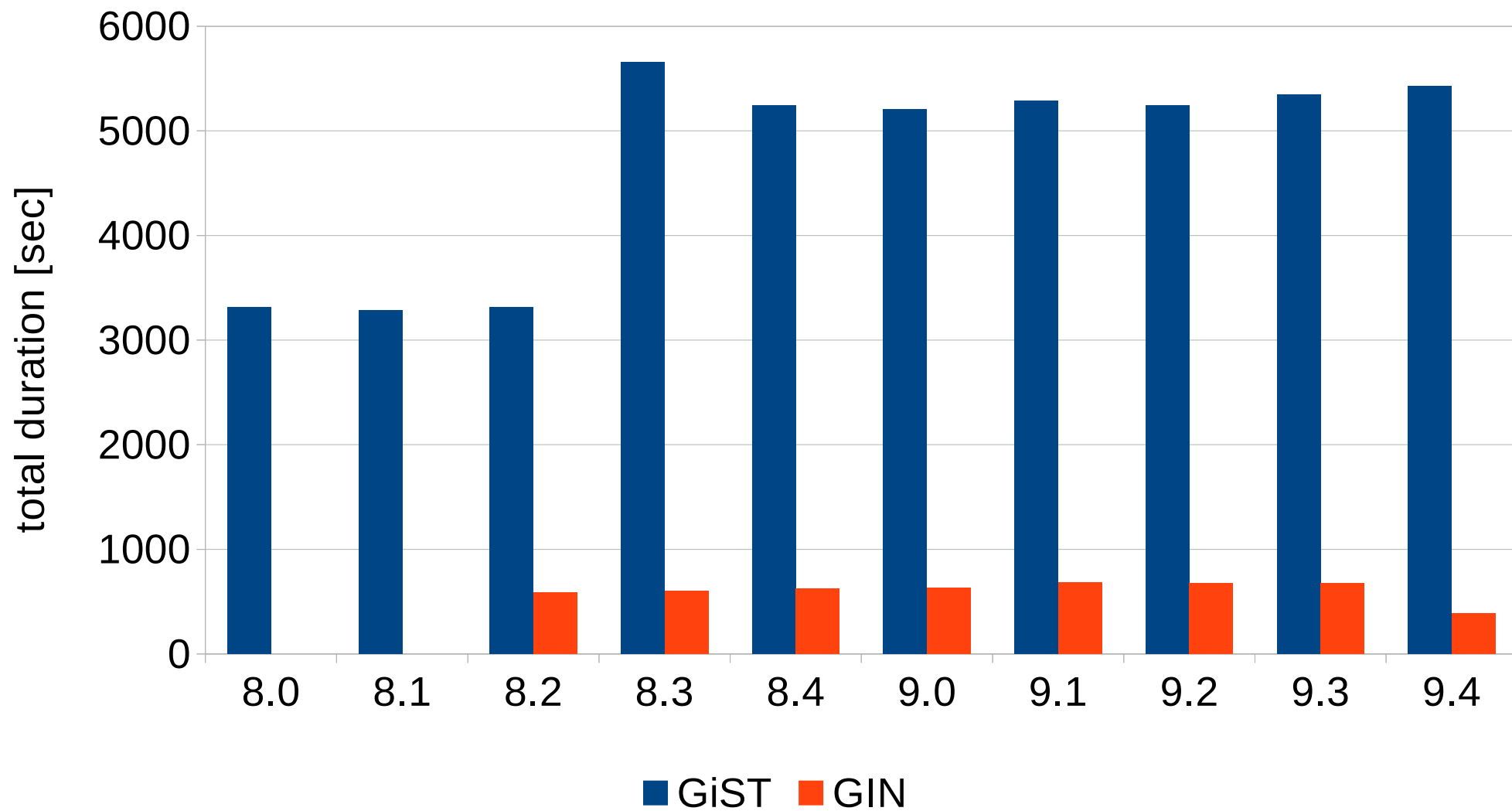
Fulltext benchmark / GIN

33k queries from postgresql.org [TOP 100]



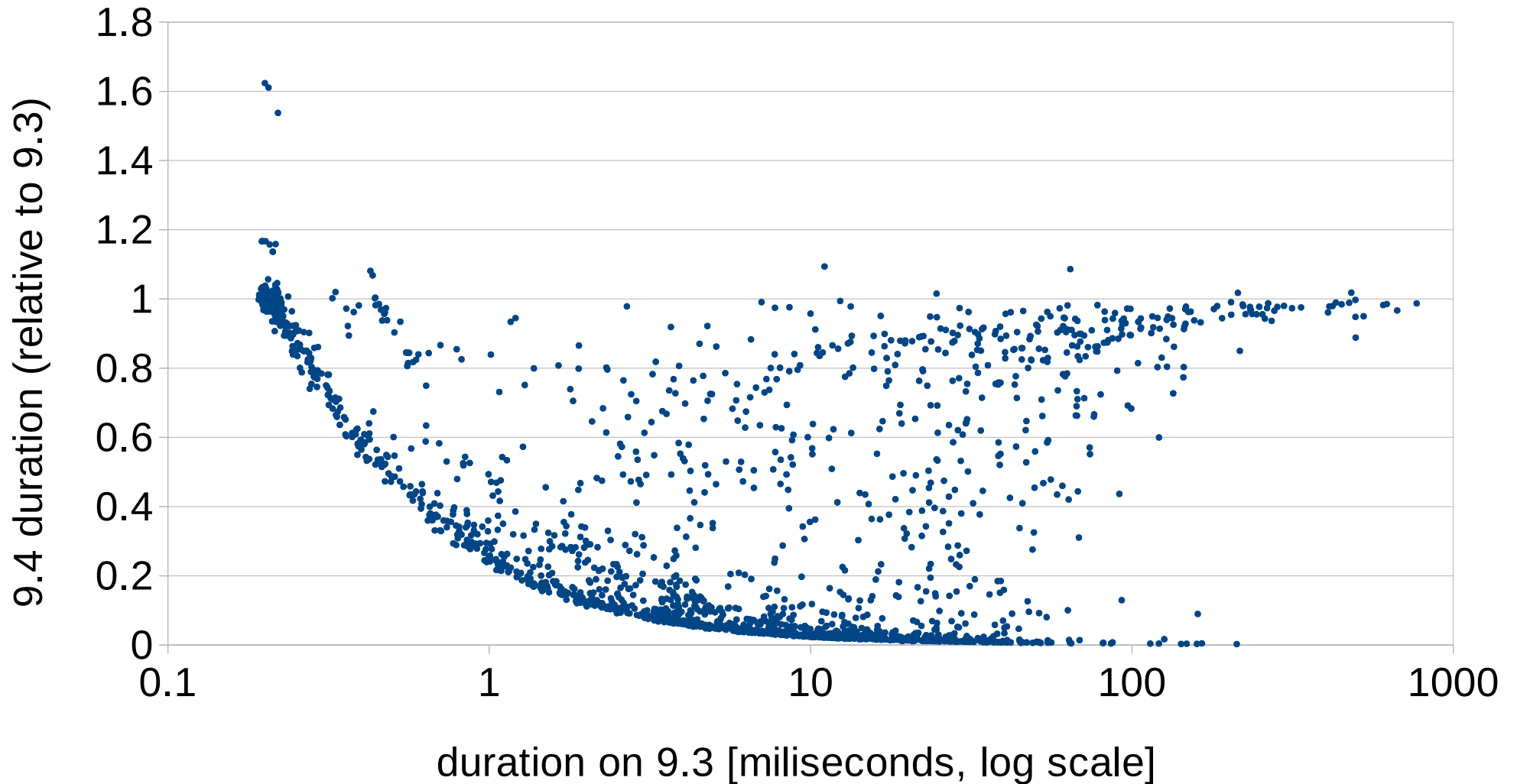
Fulltext benchmark - GiST vs. GIN

33k queries from postgresql.org [TOP 100]



Fulltext benchmark / 9.3 vs. 9.4 (GIN fastscan)

9.4 durations, divided by 9.3 durations (e.g. 0.1 means 10x speedup)



Fulltext / shrnutí

- GIN fastscan
 - dotazy kombinující “časté & vzácné”
 - 9.4 skenuje “vzácné” seznam první
 - exponenciální zrychlení pro tyto dotazy
 - ... to je celkem fajn ;-)
- jenom ~5% dotazů se zpomalilo
 - víceméně dotazy pod 1ms (chyby měření)