

Выбор оптимальной базы данных для работы с социальным графом

И.В Гражданкин, email: grazhdankin.i.i@yandex.ru

С.В Власов , email: svv@cs.vsu.ru

Цель работы

- Цель работы – тестирование и анализ скорости выполнения запросов для реализаций графовых баз данных на примере социального графа.

Графы

Применение:

- Маршруты перевозки грузов
- Взаимосвязи пользователей социальных сетей

Особенности применения:

- Очень большой объем данных
- Требуется быстрый анализ графа

Реализации базы данных

Neo4j

- Самая распространённая база данных
- Обширный функционал
- Есть бесплатная версия

Sparksee

- Заявлена высокая производительность
- Есть бесплатная версия

Методология тестирования

Тестовый стенд:

- Intel Xeon X6550 2.0Gz
- 80Gb DDR3
- 2Tb hard drive

Настройка Neo4j:

- Version 3.5.11
- Ubuntu 18.04 LTS
- Cache size 60Gb

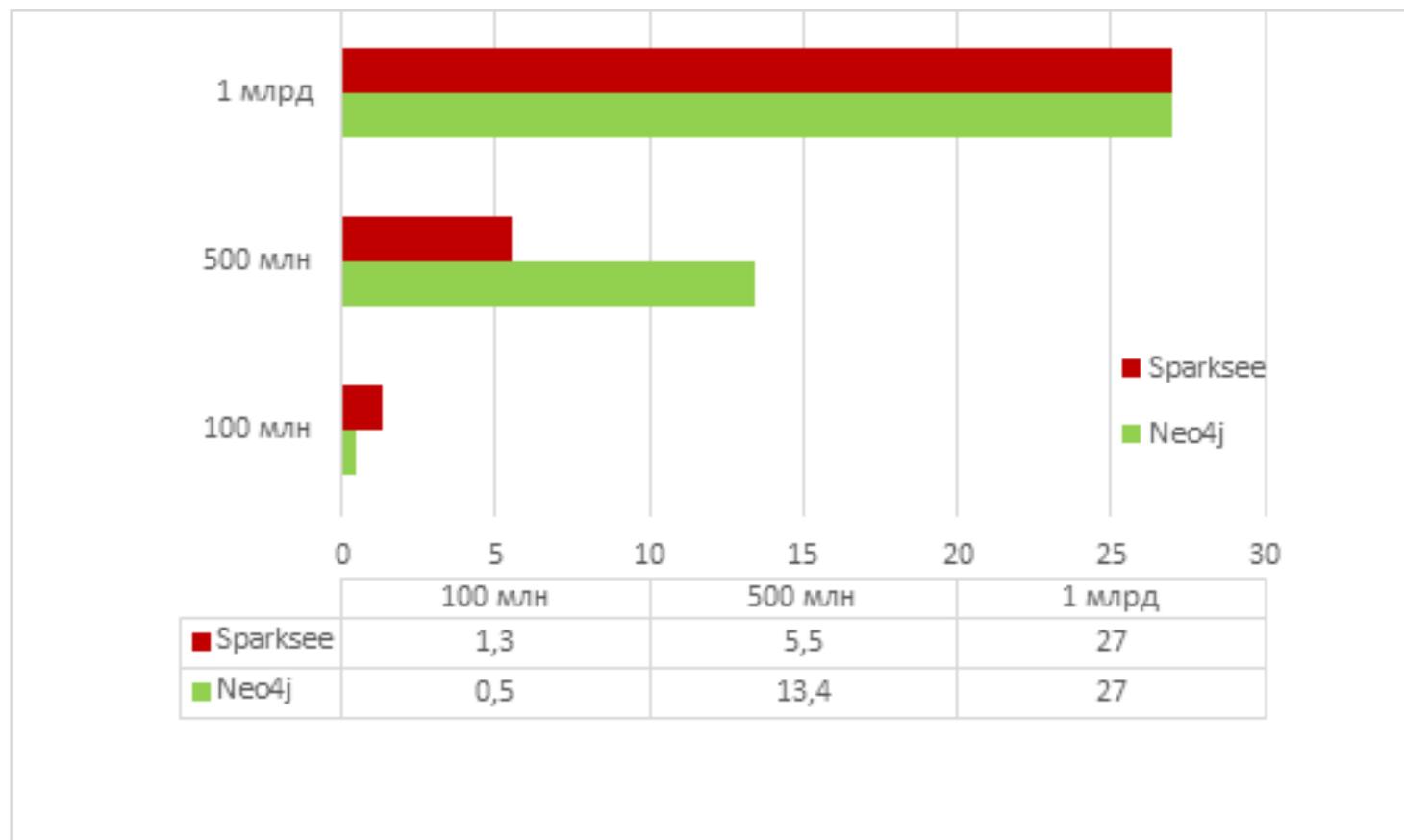
Настройка Sparksee:

- Version 5.2
- Ubuntu 18.04 LTS
- Cache size 60Gb

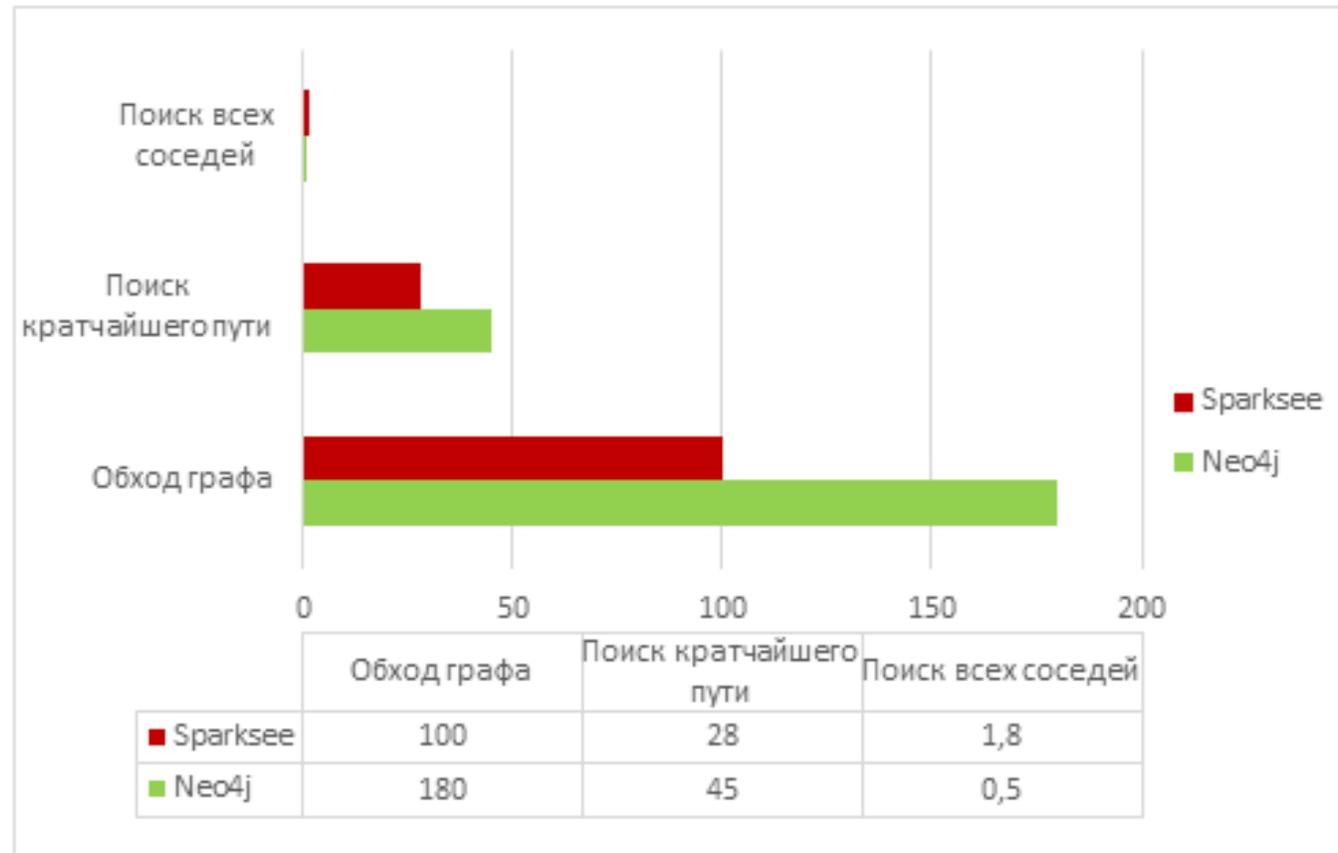
Аналитические запросы:

- Простой запрос – получить всех соседей вершины
- Более сложный запрос - найти кратчайший путь
- Сложный запрос – выполнить обход графа

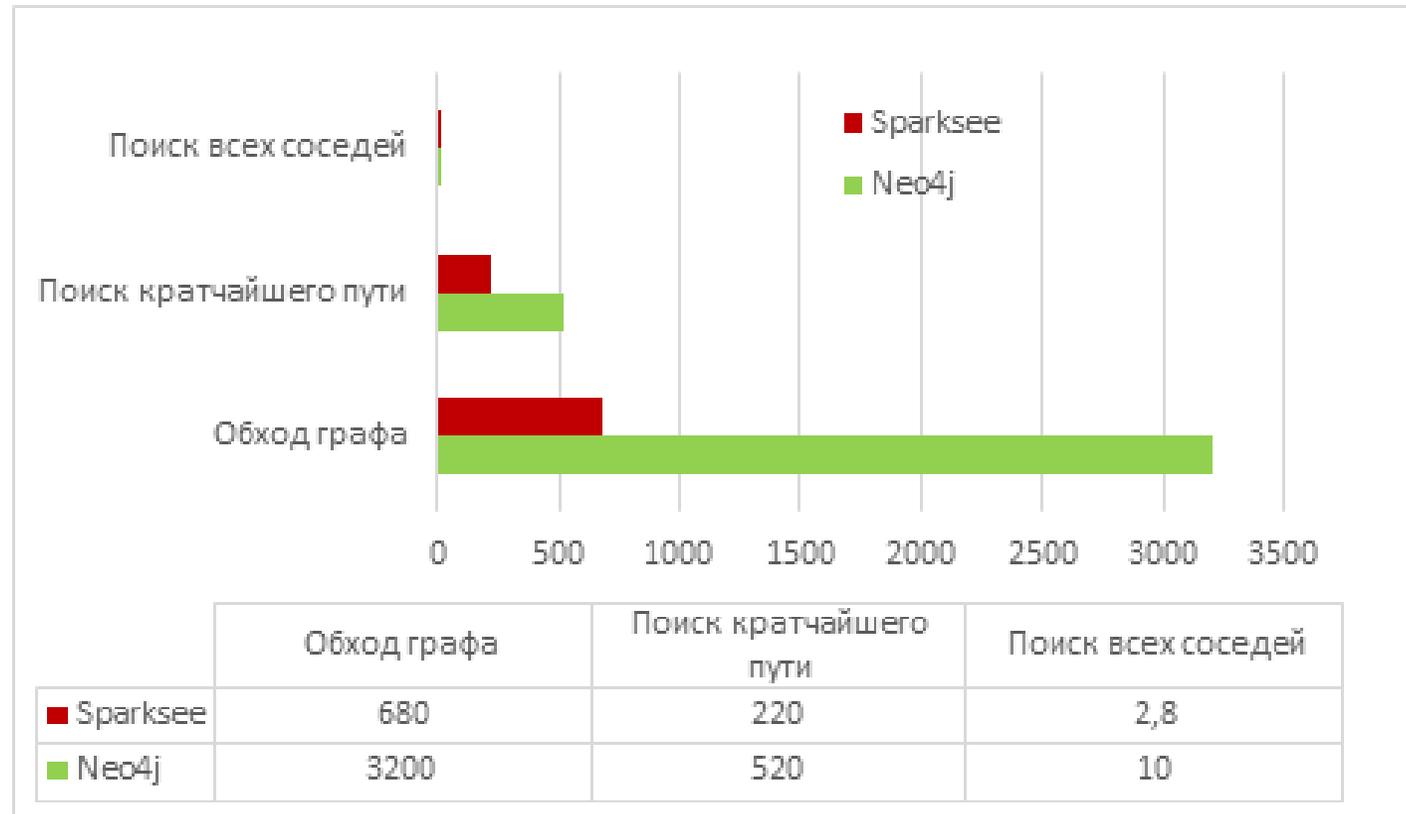
Время импорта данных



Время обработки аналитических запросов 100 млн



Время обработки аналитических запросов 500 млн



Выводы

- Время импорта занимает разумное время если, количество ребер не превышает одного миллиарда
- Sparksee предпочтительнее, чем Neo4j по производительности. Это очень хорошо видно при выполнении сложных аналитических запросов и при использовании большей части графа
- Главной характеристикой для высокой производительности графовых баз данных работы с графом, у которого количество ребер превышает один миллиард необходимо использовать кластерные решения