# stackoverfloooooooow

Questions | Tags | Users | Badges | Unanswered | Ask Question

Ten. Million. Questions. Let's celebrate all we've done together.

Stack Overflow is a question and answer site for professional and enthusiast programmers. It's 100% free.          Take the 2-minute tour      ✕

## Seeking Neo4J Cypher query for long but (nearly) unique paths

▲
**3**
▼
★

We have a Neo4J database representing an evolutionary process with about 100K nodes and 200K relations. Nodes are individuals in generations, and edges represent parent-child relationships. The primary goal is to be able to take one or nodes of interest in the final generation, and explore their evolutionary history (roughly, "how did we get here?").

The "obvious" first query to find all their ancestors doesn't work because there are just too many possible ancestors and paths through that space:

```
match (a)-[:PARENT_OF*]->(c {is_interesting: true})
return distinct a;
```

So we've pre-processed the data so that some edges are marked as "special" such that *almost* every node has at most one "special" parent edge, although occasionally both parent edges are marked as "special". My hope, then, was that this query would (efficiently) generate the (nearly) unique path along "special" edges:

```
match (a)-[r:PARENT_OF* {special: true}]->(c {is_interesting: true})
return distinct a;
```

This, however, is still unworkably slow.

This is frustrating because "as a human", the logic is simple: Start from the small number of "interesting" nodes (often 1, never more than a few dozen), and chase back along the almost always unique "special" edges. Assuming a very low number of nodes with two "special" parents, this should be something like O(N) where N is the number of generations back in time.

In Neo4J, however, going back 25 steps from a unique "interesting" node where *every* step is unique, however, takes 30 seconds, and once there's a *single* bifurcation (where both parents are "special") it gets worse much faster as a function of steps. 28 steps (which gets us to the first bifurcation) takes 2 minutes, 30 (where there's still only the one bifurcation) takes 6 minutes, and I haven't even thought to try the full 100 steps to the beginning of the simulation.

Some similar work last year seemed to perform better, but we used a variety of edge labels (e.g., `(a)-[:SPECIAL_PARENT_OF*]->(c)` as well as `(a)-[:PARENT_OF*]->(c)` ) instead of using data fields on the edges. Is querying on relationship field values just not a good idea? We have quite a few different values attached to a relationship in this model (some boolean, some numeric) and we were hoping/assuming we could use those to efficiently limit searches, but maybe that wasn't really the case.

Suggestions for how to tune our model or queries would be greatly appreciated.

**Update** I should have mentioned, this is all with Neo4J 2.1.7. I'm going to give 2.2 a try as per Brian Underwood's suggestion and will report back.

neo4j    cypher

share  improve this question

edited Apr 6 at 14:59              asked Apr 6 at 5:41
                                   👤 Nic McPhee
                                      16 ● 5

> I'd really like to see an answer to this myself. I've definitely had some Cypher queries with open-ended path lengths which seemed like they shouldn't be traversing very much but which take a while. –
> Brian Underwood Apr 6 at 7:24

add a comment

## 2 Answers

active     oldest     **votes**

### Related

0   neo4j not reusing existing vertex in cypher create unique query

4   Cypher Query in Neo4j Returns 'undefined'

0   Neo4j crashes on 4th degree Cypher query

0   Slow neo4j Cypher Queries, using Baconator query for less-short paths than shortest path

0   Creating unique path in Neo4J using Cypher and SDN @Query annotation

2   Neo4j Cypher alternative paths

2   Cypher MATCH query speed

0

After exploring things with the profiling tools in Neo4J 2.2 (thanks to Brian Underwood for the tip) it's pretty clear that (at the moment) Neo4J doesn't do any pre-filtering on edge properties, which leads to nasty combinatorial explosions with long paths.

For example the original query:

```
match (a)-[r:PARENT_OF* {special: true}]->(c {is_interesting: true})
return distinct a;
```

finds *all* the paths from `a` to `c` and *then* eliminates the ones that have edges that aren't `special`. Since there are many millions of paths from `a` to `c`, this is totally infeasible.

If I instead add a `IS_SPECIAL` edge wherever there was a `PARENT_OF` edge that had `{special: true}`, then the queries become really fast, allowing me to push back around 100 generations in under a second.

This query creates all the new edges:

```
match (a)-[r:PARENT_OF {special: true}]->(b)
create (a)-[:IS_SPECIAL]->(b);
```

and takes under a second to add 91K relationships in our graph.

Then

```
match (c {is_interesting: true})
with c
match (a)-[:IS_SPECIAL*]->(c)
return distinct a;
```

takes under a second to find the 112 nodes along the "special" path back from a unique target node `c` . Matching `c` first and limiting the set of nodes using `with c` seems to also be

important, as Neo4J doesn't appear to pre-filter on node properties either, and if there are several "interesting" target nodes things get a lot slower.

share  improve this answer

add a comment

## Your Answer

Post Your Answer

*By posting your answer, you agree to the privacy policy and terms of service.*

Not the answer you're looking for? Browse other questions tagged `neo4j` `cypher` or ask your own question.

question feed