

Incremental Computation

Alex Degtiar, Ram Raghunathan, Thomas Marshall

Motivation

- Big data processing takes a lot of time/resources (e.g. PageRank of entire internet)
- Data can change just a little (one new page gets added), but updating result requires recomputation across entire data set.

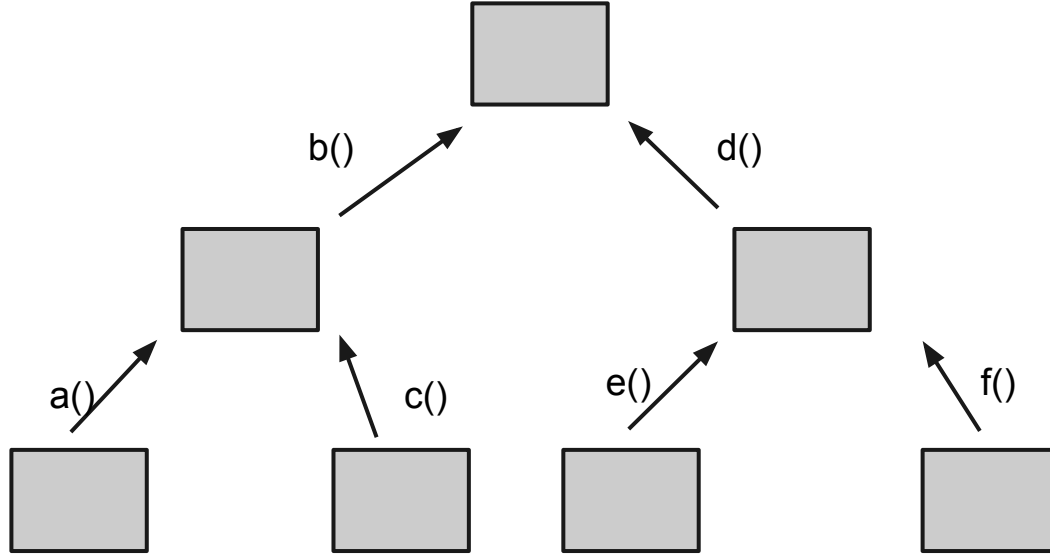
Prior Work

- Specialized systems - distributed data flow, materialized views
- In memory systems - Adaptive Functional Programming (Acar, 2006)

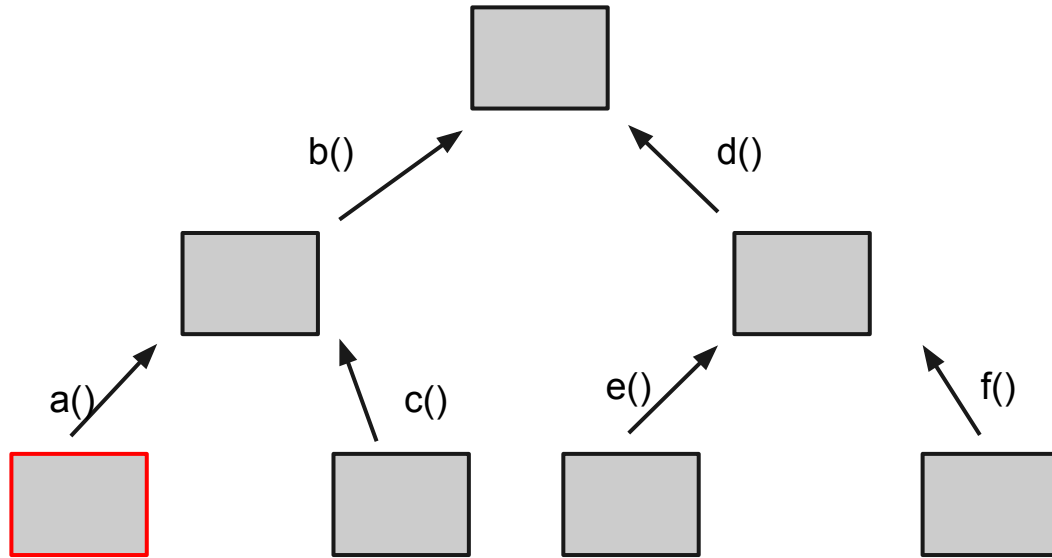
Theory

- Keep track of relationships between computation and data in a Dynamic Dependency Graph (DDG).
- When data changes, it can be propagated through the graph to minimize recomputation needed.

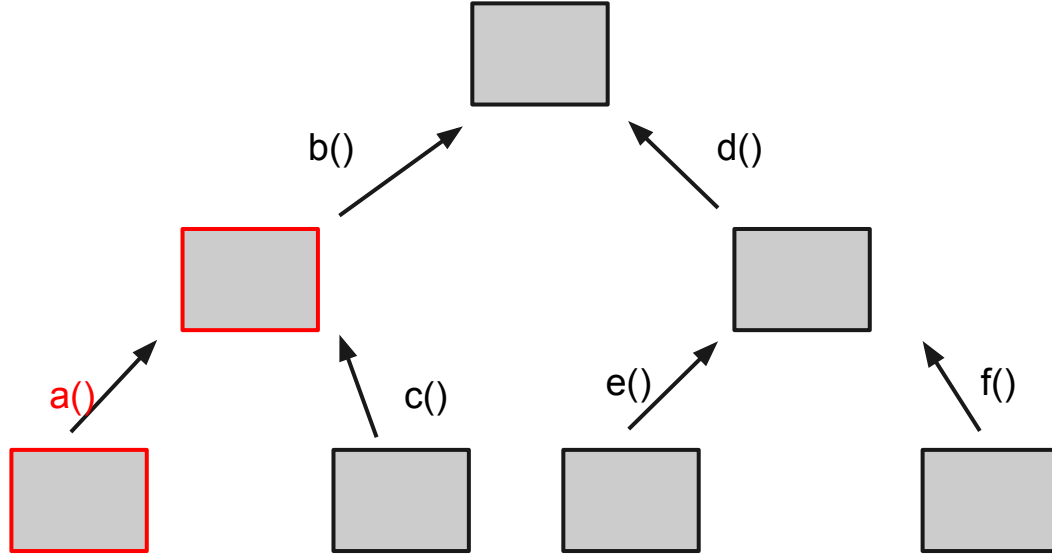
DDG



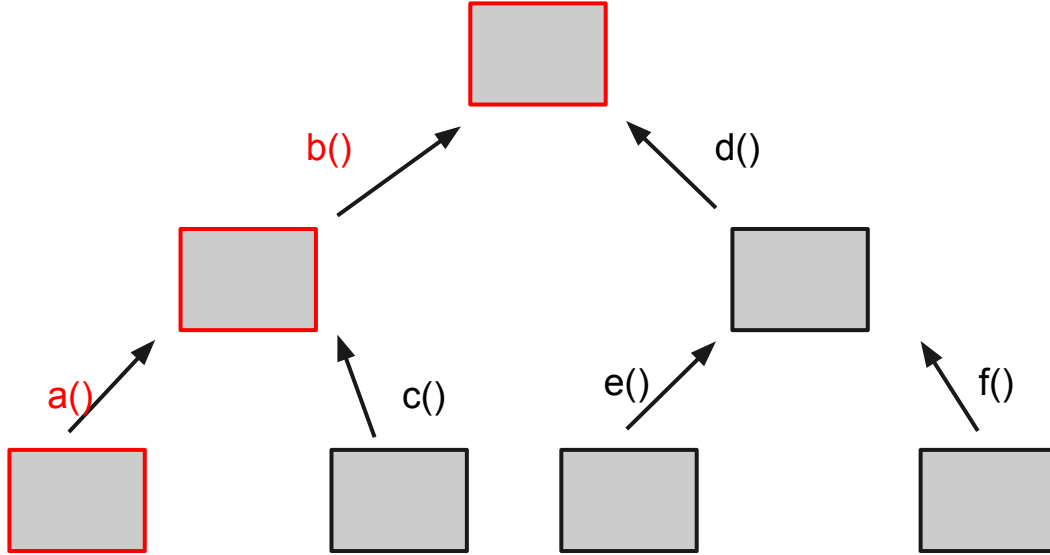
DDG



DDG



DDG



Implementation

- Back-end datastore - Postgres
 - Contains input data along with intermediate results and DDG.
- Python library
 - Contains primitives for interacting with data in a DDG-aware manner.

IncLib

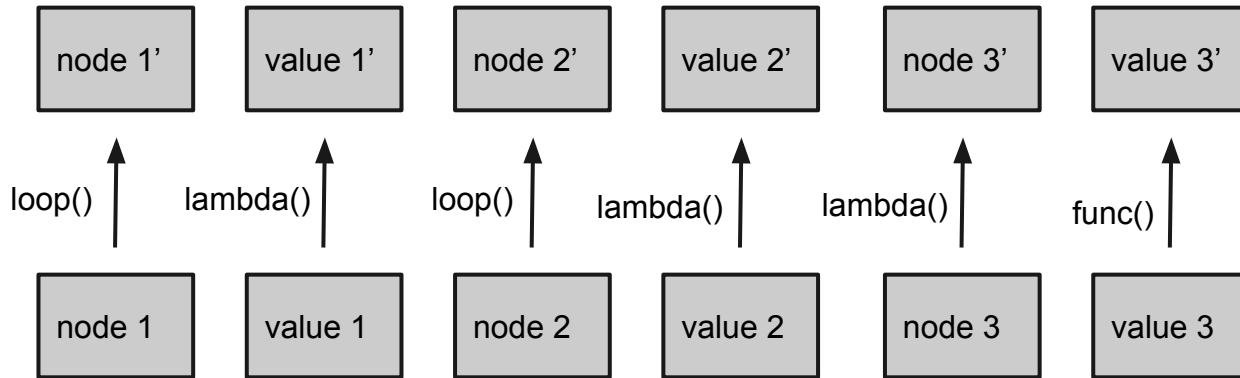
- How to track memory accesses?
 - modifiabiles
 - continuation passing style

```
def mod(initializer) -> mod
def write(mod, value)
def read(mod, func)
def select(pred) -> mod_list
```

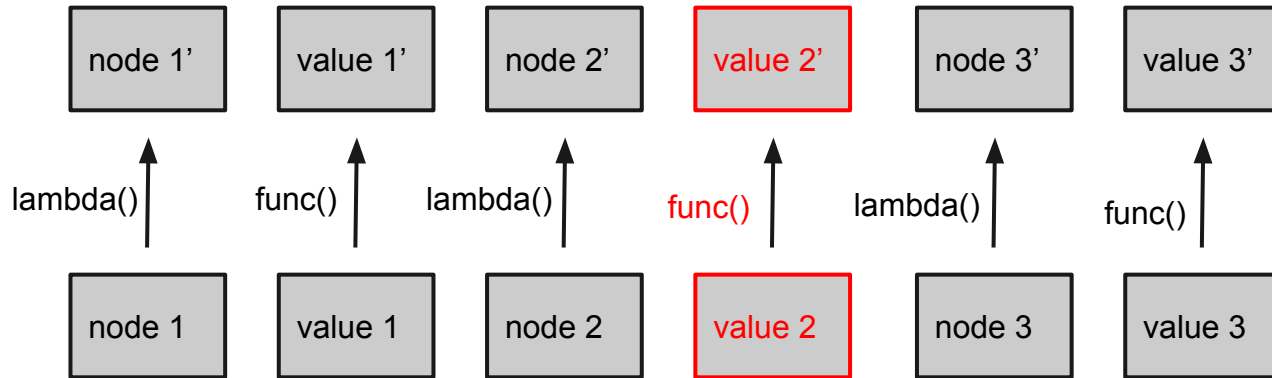
Continuation Passing Style

```
def map(mod_list):
    def m(lst, d):
        read(lst, lambda l:
            (val, nxt) = l
            if nxt is not None:
                write(d, (mod(lambda d: read(val, lambda v: write(d, func(v)))),
                           mod(lambda d: write(d, m(nxt, d))))
            else:
                write(d, (mod(lambda d: read(val, lambda v: write(d, func(v)))),
                           None)))
    return mod(lambda d: m(mod_list, d))
```

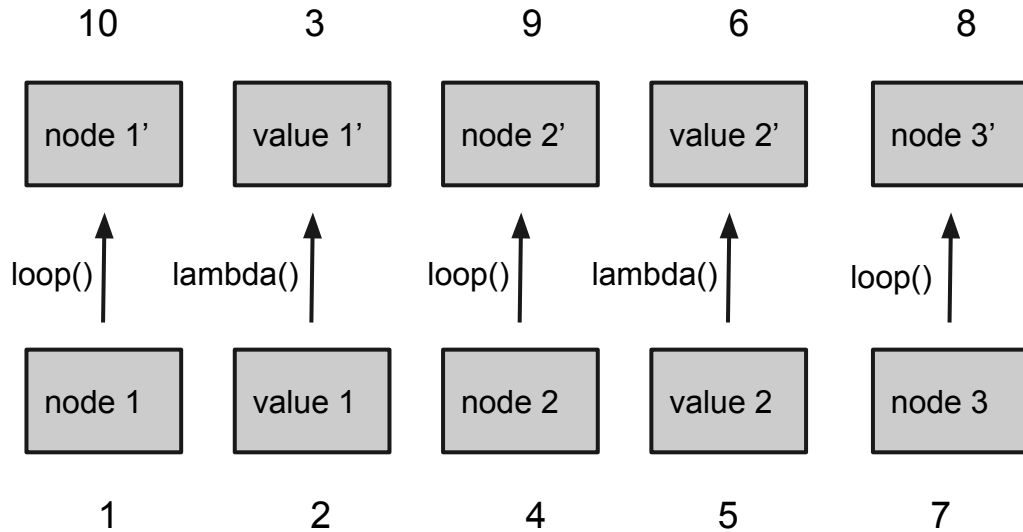
DDG of Map



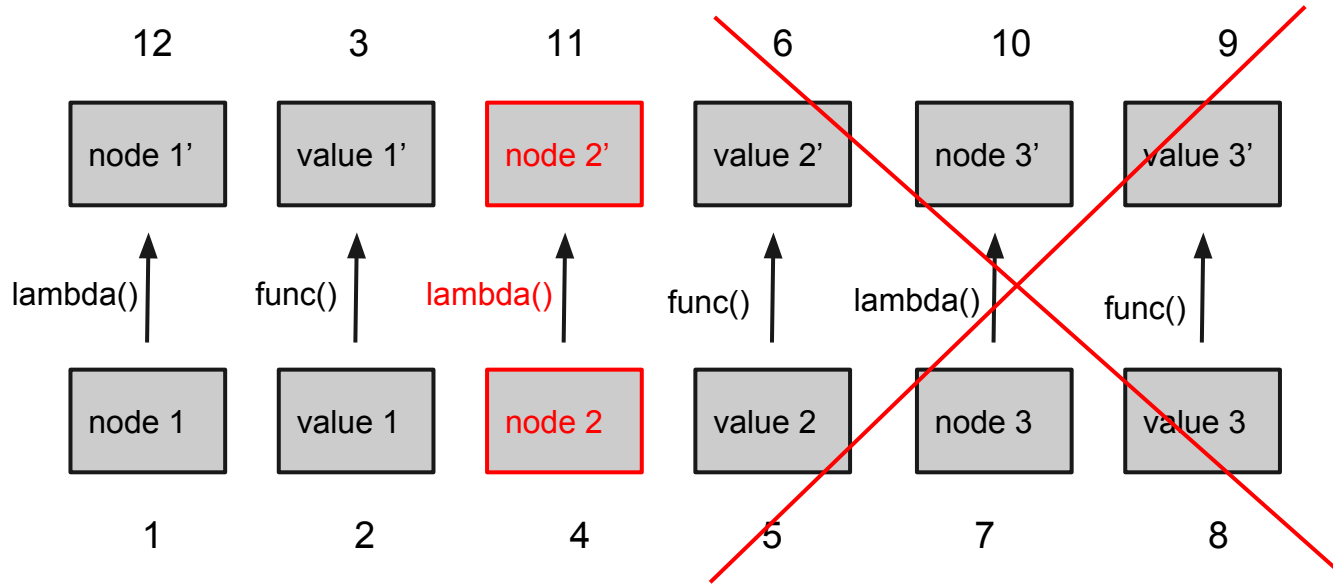
DDG of Map



DDG of Map



DDG of Map



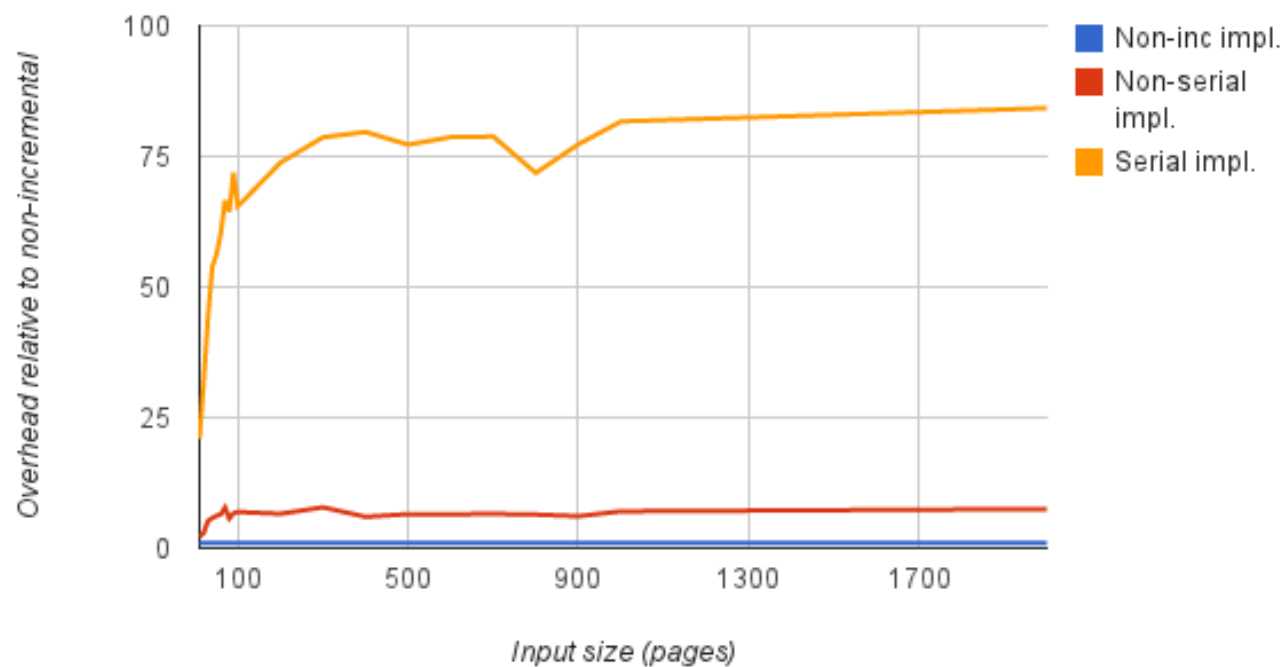
Experiments

- Word Frequency calculation across “pages”
- Used the Wikipedia corpus as input
- Compared against typical non-incremental implementation
- Varied input size and the number of inserts/updates

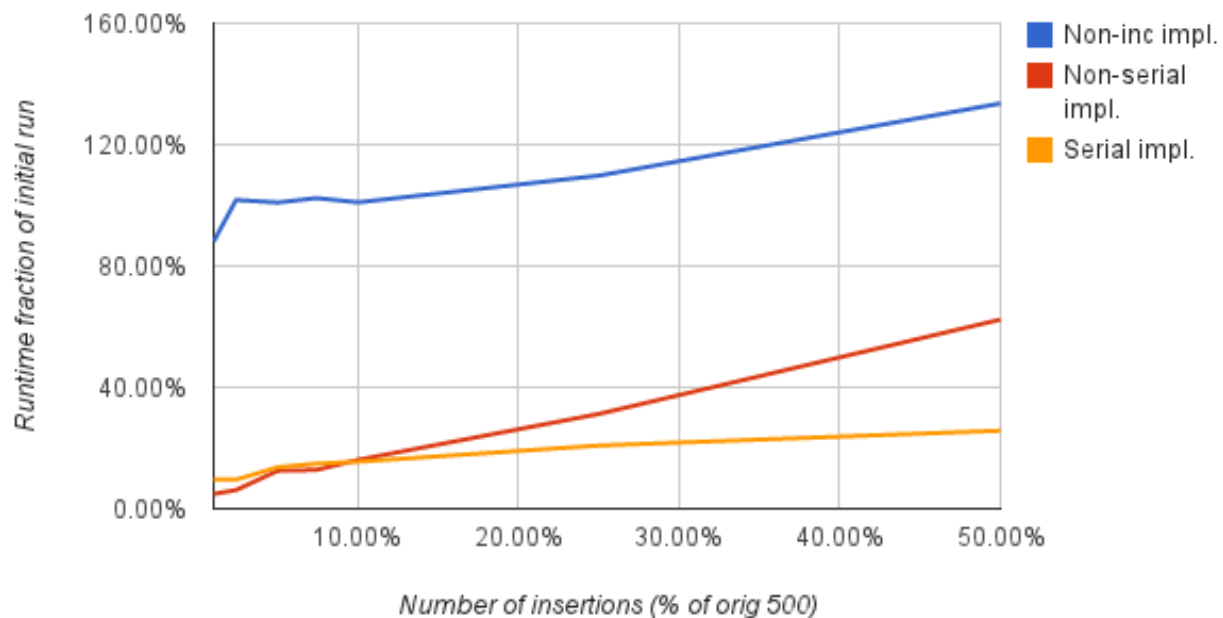
Results

- Overhead of maintaining dependencies
 - Lots of CPU processing (esp. serialization)
 - High Postgres overhead, partially addressed by tuning. More substantially addressed by avoiding serialization.
 - Some space overhead in postgres
- Improvements with re-running after modifying data

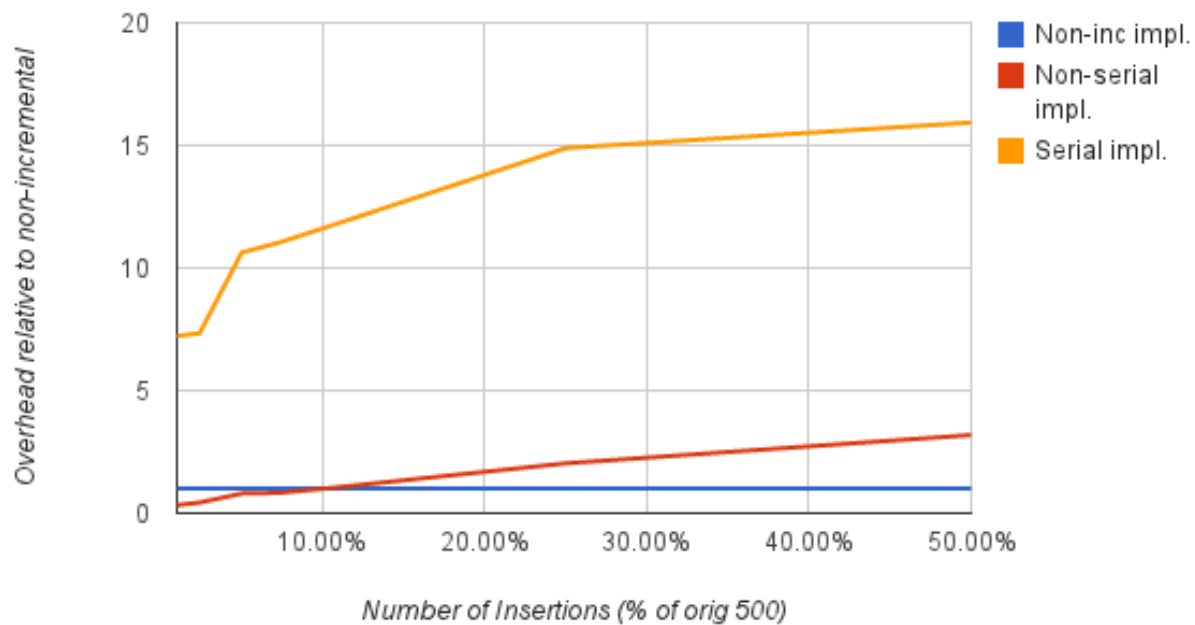
Overhead of Single Run



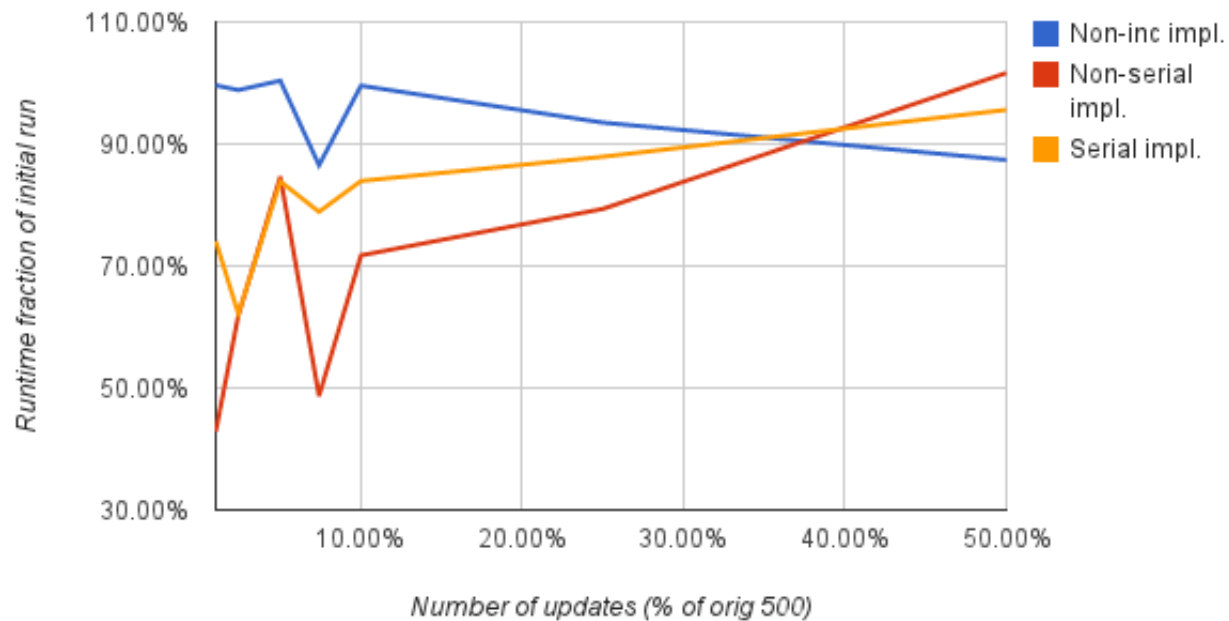
Initial-run vs. re-run after insertions



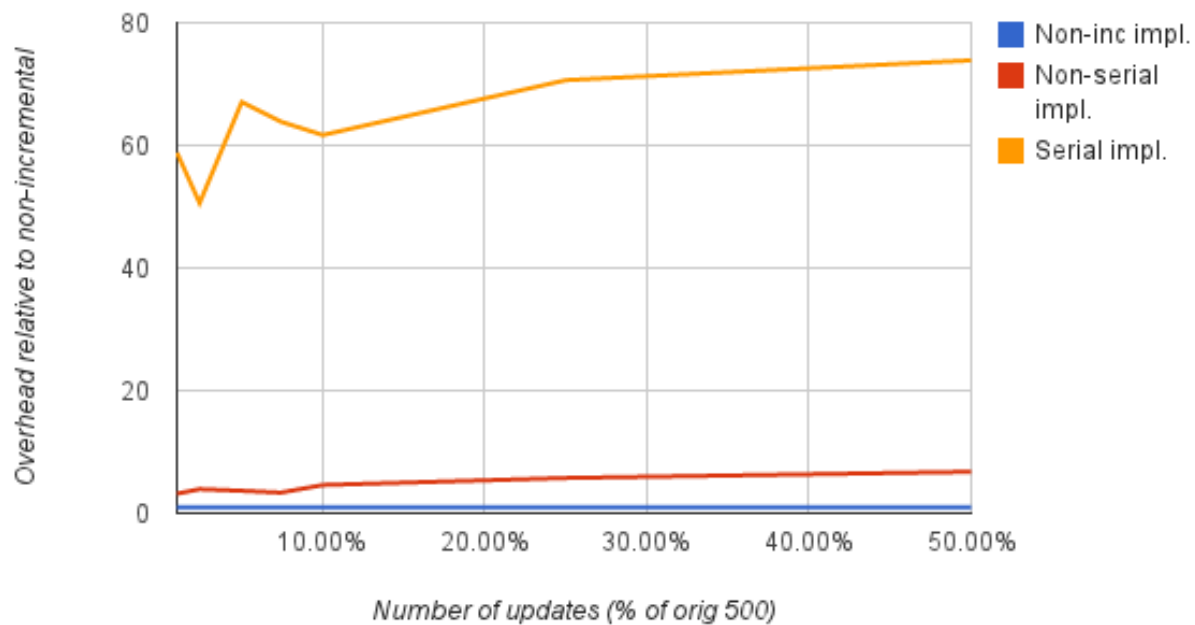
Re-run overhead vs. number of insertions



Initial-run vs. re-run after updates



Re-run overhead vs. number of updates



Conclusions

- Avoid serialization overhead
- Don't need many Postgres features
- Utility of incremental computation depends on workload