# In MapReduce

# Evaluating Recursive Queries at Scale

*(compute the next generation of answers)*

*(remove the ones we've already seen)*

Join

Difference

$\Delta A_{i-1}$ → map

$R^{(0)}$ → map

$R^{(1)}$ → map

map → reduce → map

reduce → map

(a)

(b)

$A_i^{(0)}$ → map
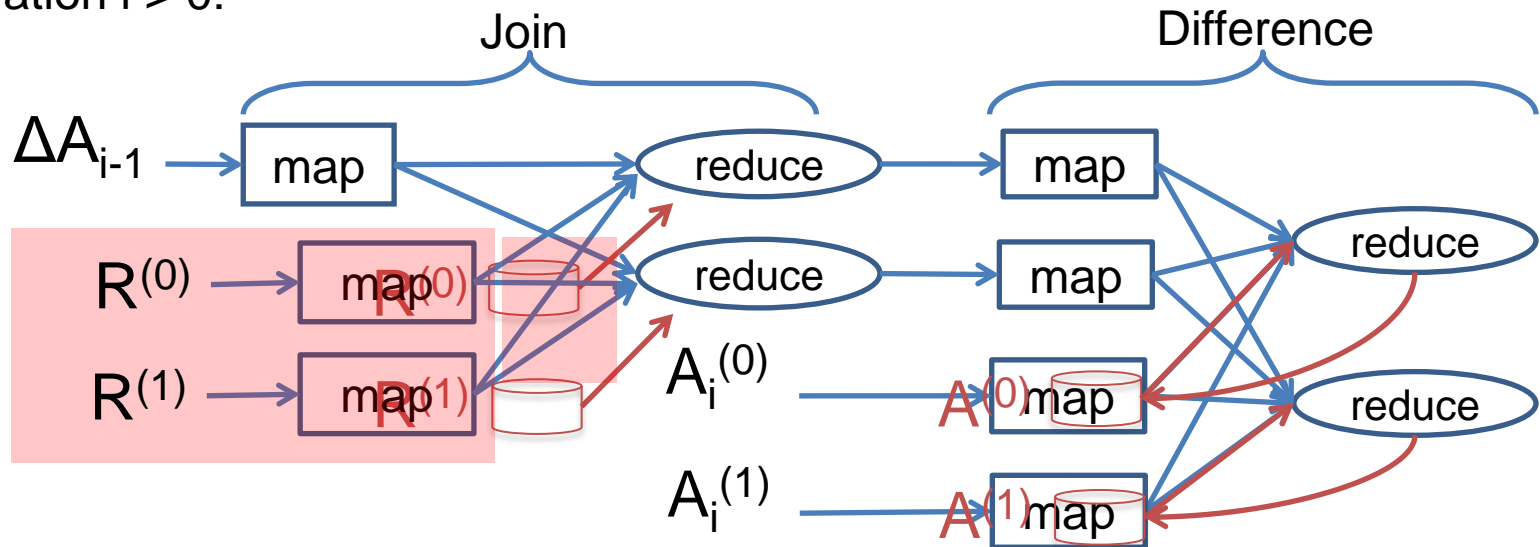
$A_i^{(1)}$ → map

reduce

reduce

*(a) R is loop invariant, but gets loaded and shuffled on each iteration*
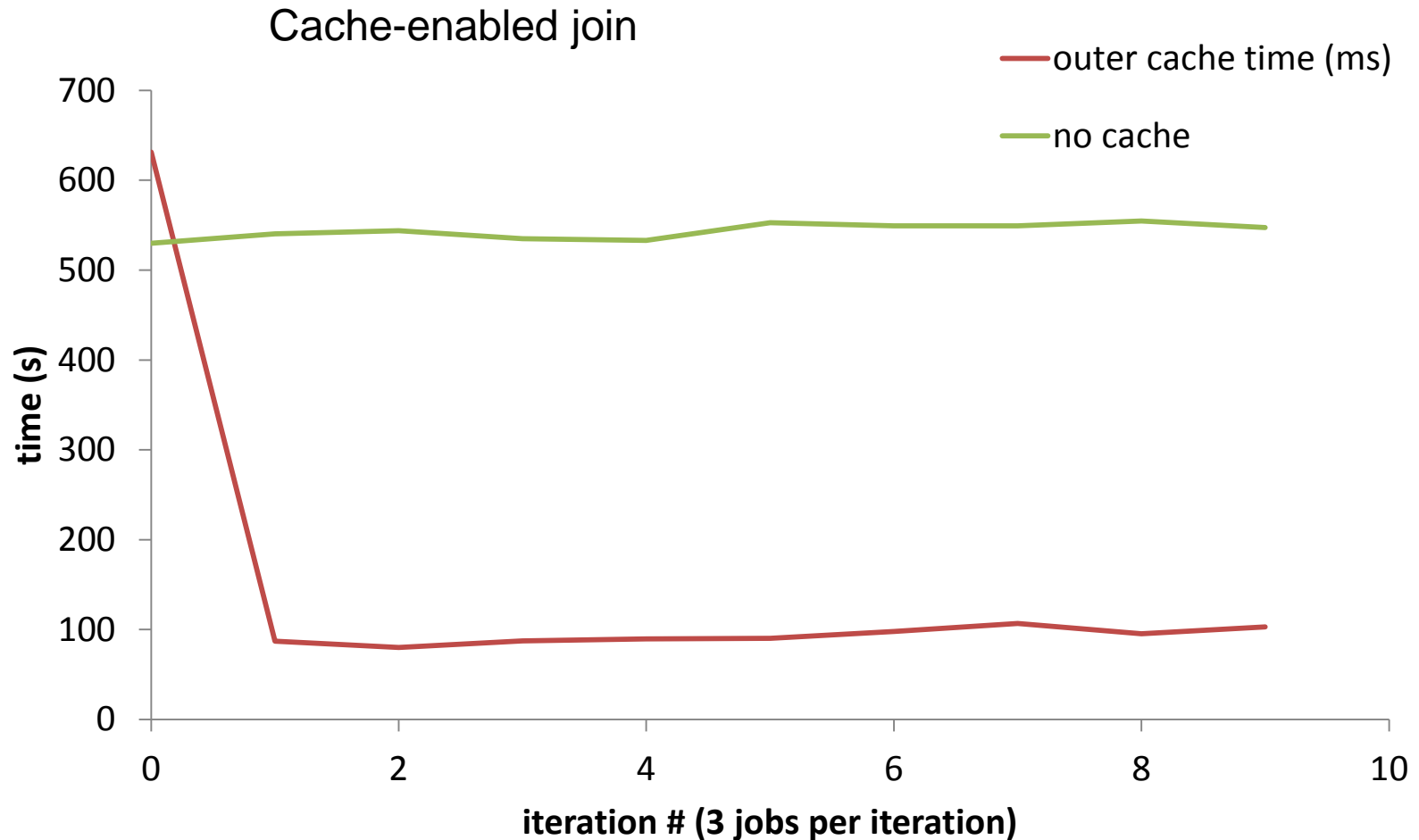
*(b) $A_i$ grows slowly and monotonically, but is loaded and shuffled on each iteration.*
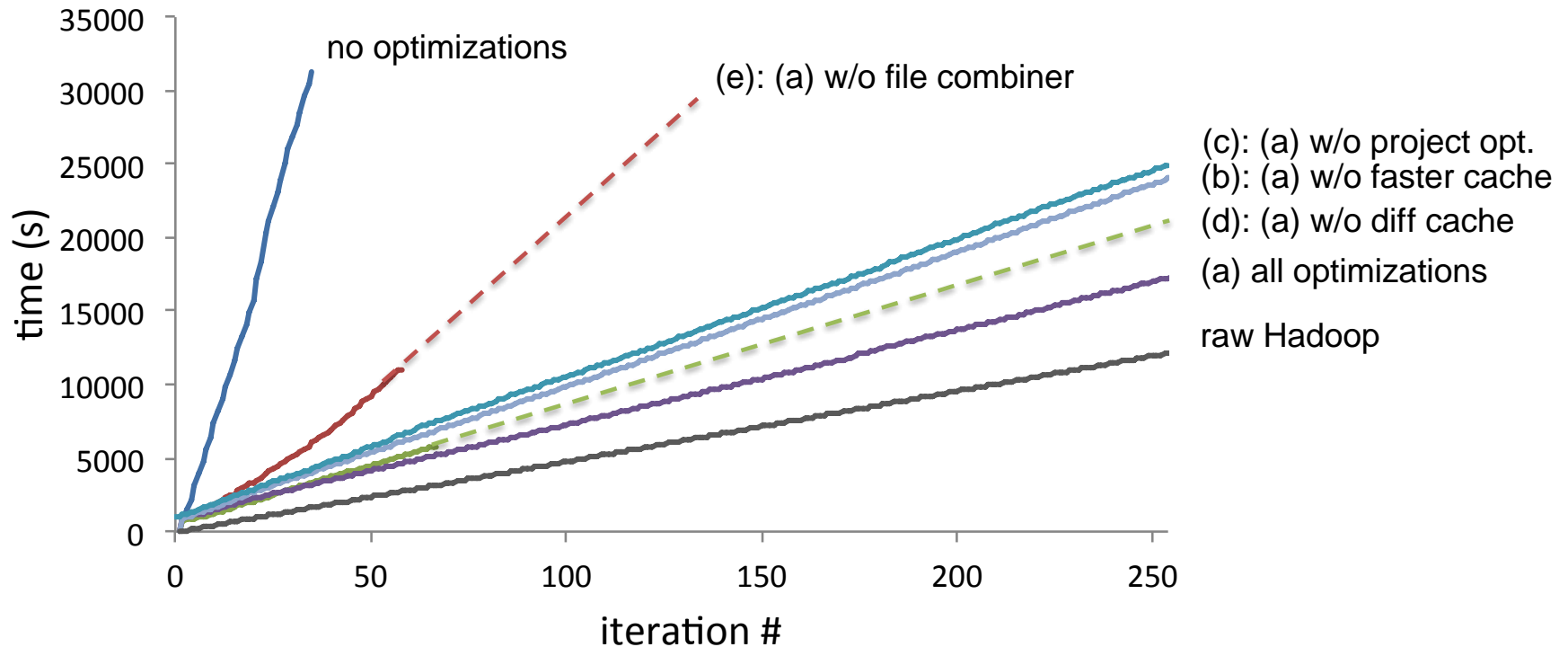
# Idea: Cache Loop-Invariant Data

Iteration i = 0: Load a distributed cache

Iteration i > 0:

UNIVERSITY *of* WASHINGTON



Cache-enabled join

# Effect of various optimizations for a recursive graph query on BTC 2010
(query: transitive reachability from 7 nodes)



Takeaways:
- 10x improvement over no optimizations.
- All optimizations are useful
- We're approaching the raw overhead of Hadoop (bottom gray line)

## New tuples discovered by iteration number