

A Tale of Two Hash Tables

Vlad Roubtsov, 2017

A wise Chinese man once said:

“Everything is an Array”

(YY, Tolo Technologies)

Want:

- Want a fast hashtable for data compression, order books, etc
- Want memory compactness
- Want low latency for all ops, ideally all the way to load factors of ~100%
- Try not to require “special” K or V
- std::unordered_map is messed up

Books have written about (coalesced)
hashing...

INTERNATIONAL SERIES OF MONOGRAPHS
ON COMPUTER SCIENCE

Design and Analysis of Coalesced Hashing

JEFFREY S. VITTER

and

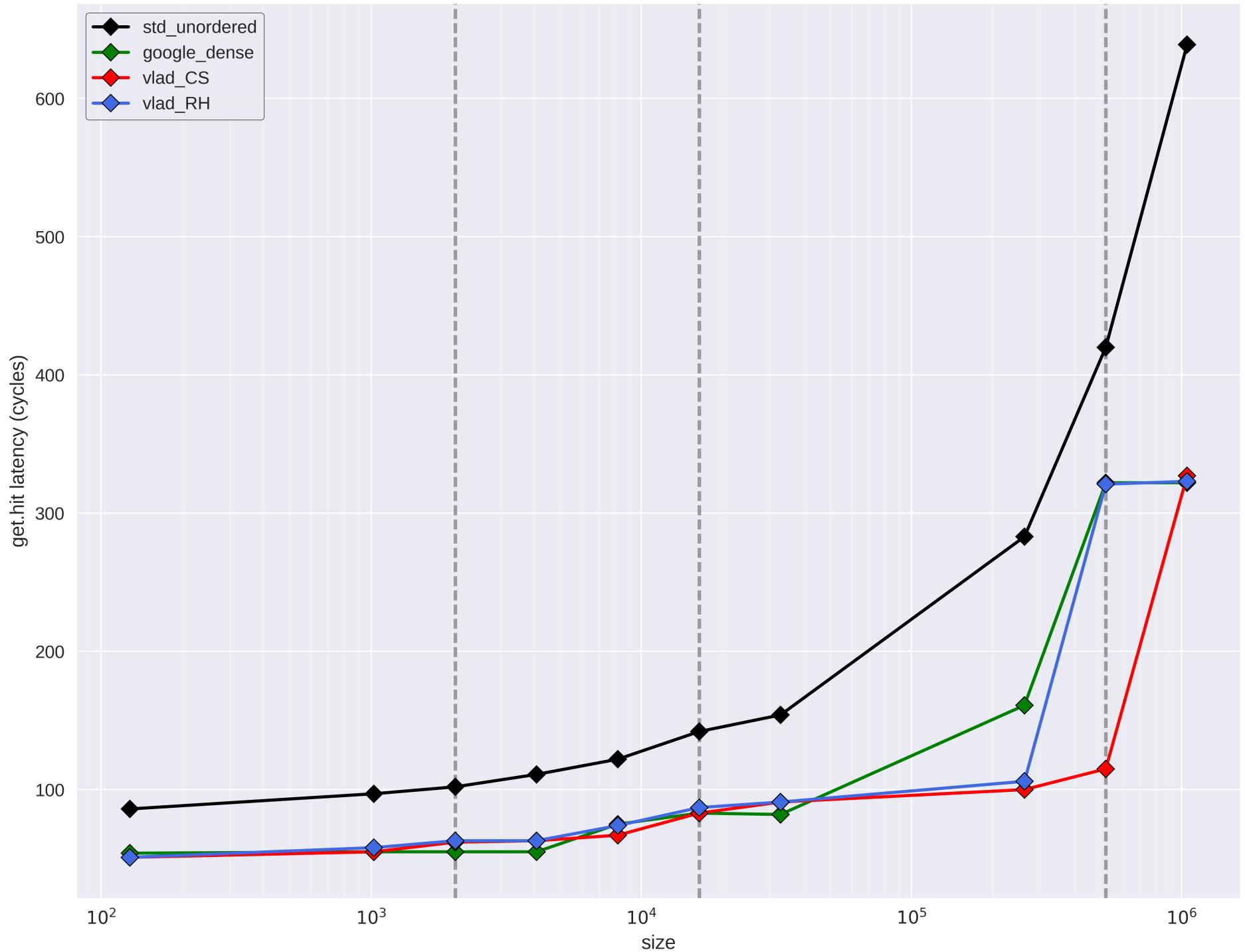
WEN-CHIN CHEN



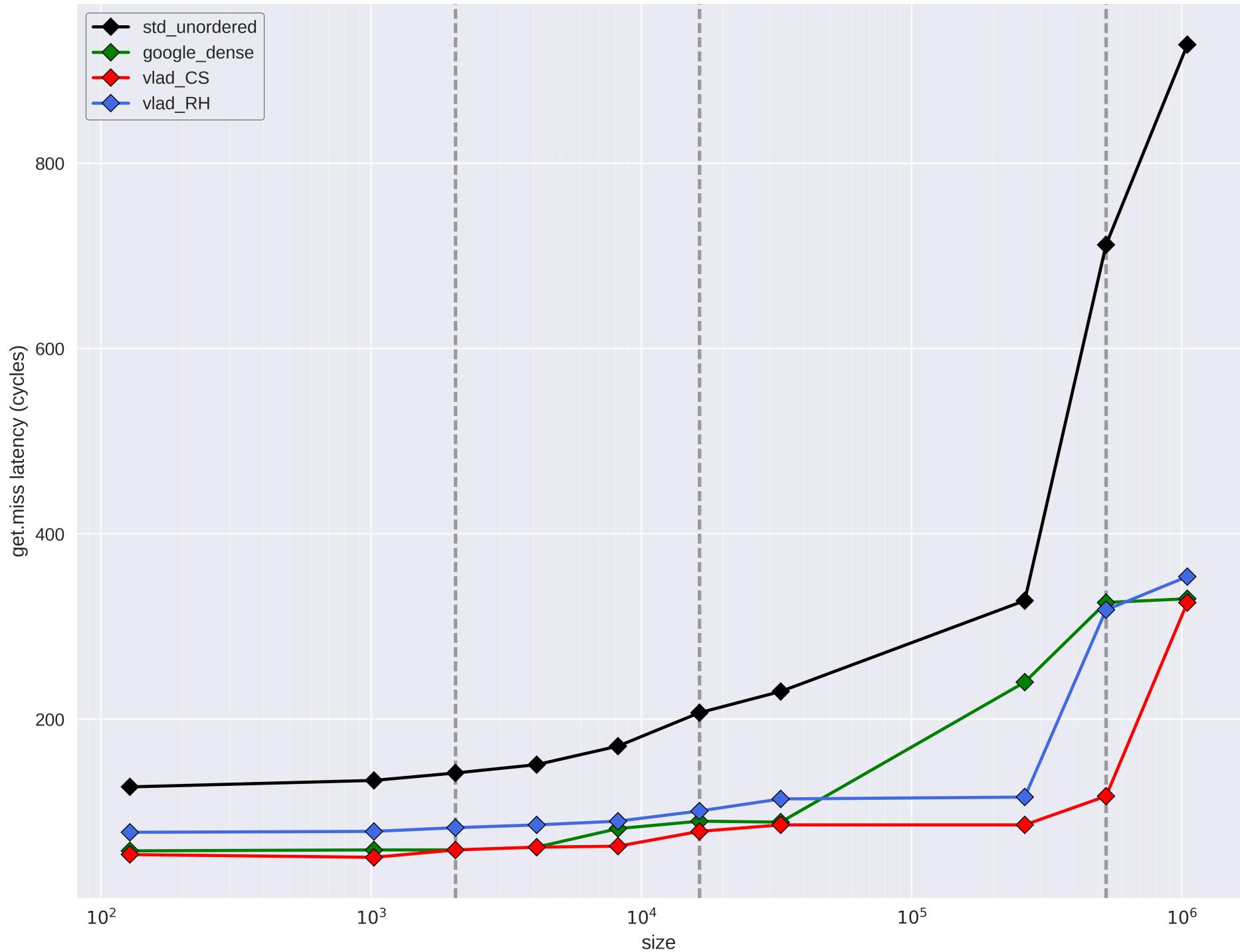
Keywords:

- Coalesced hashing
- Hopscotch hashing
- Brent's variant
- Robin Hood hashing
- Tombstone deletion
- Backshift deletion

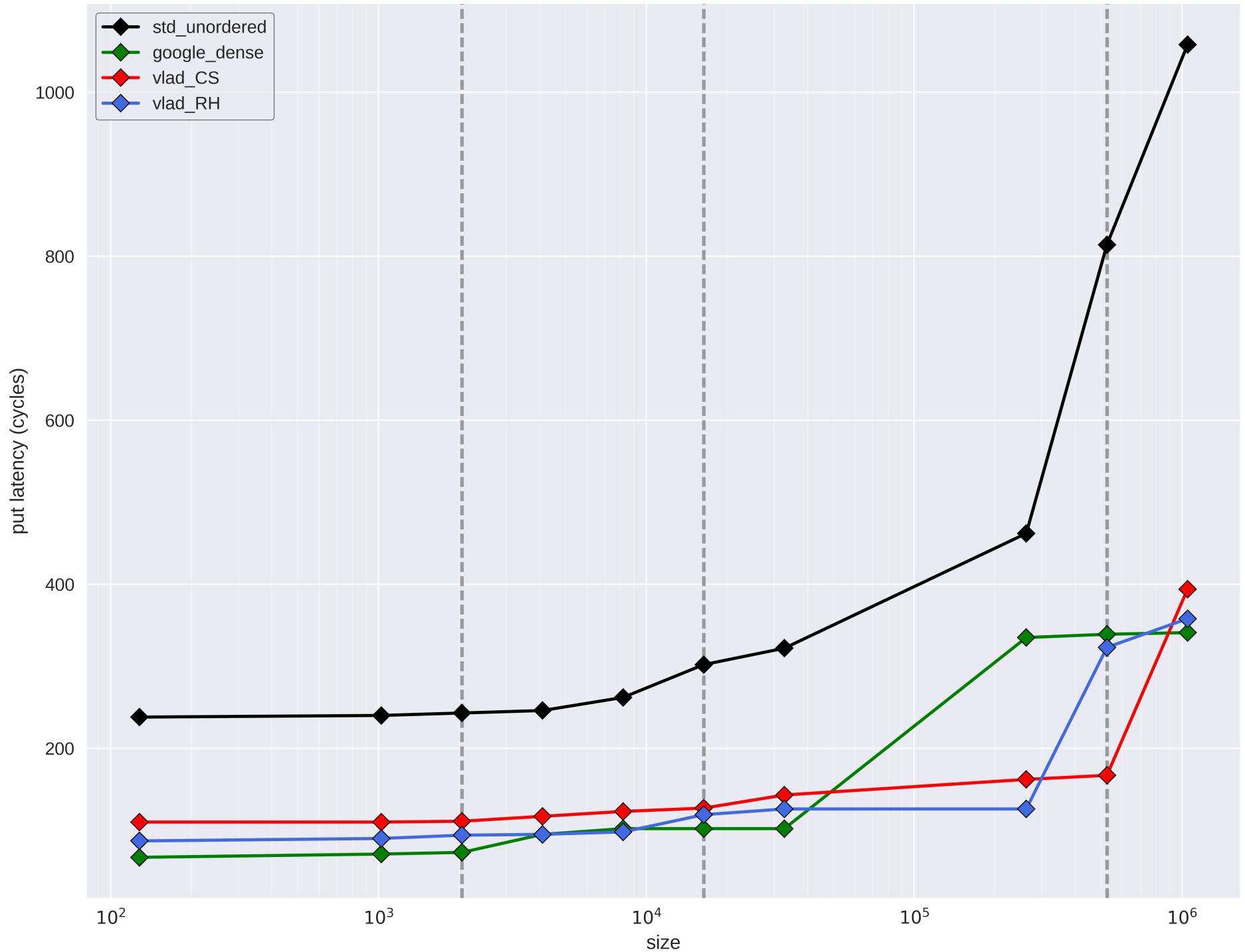
MEDIAN get.hit latency (cycles) vs table size



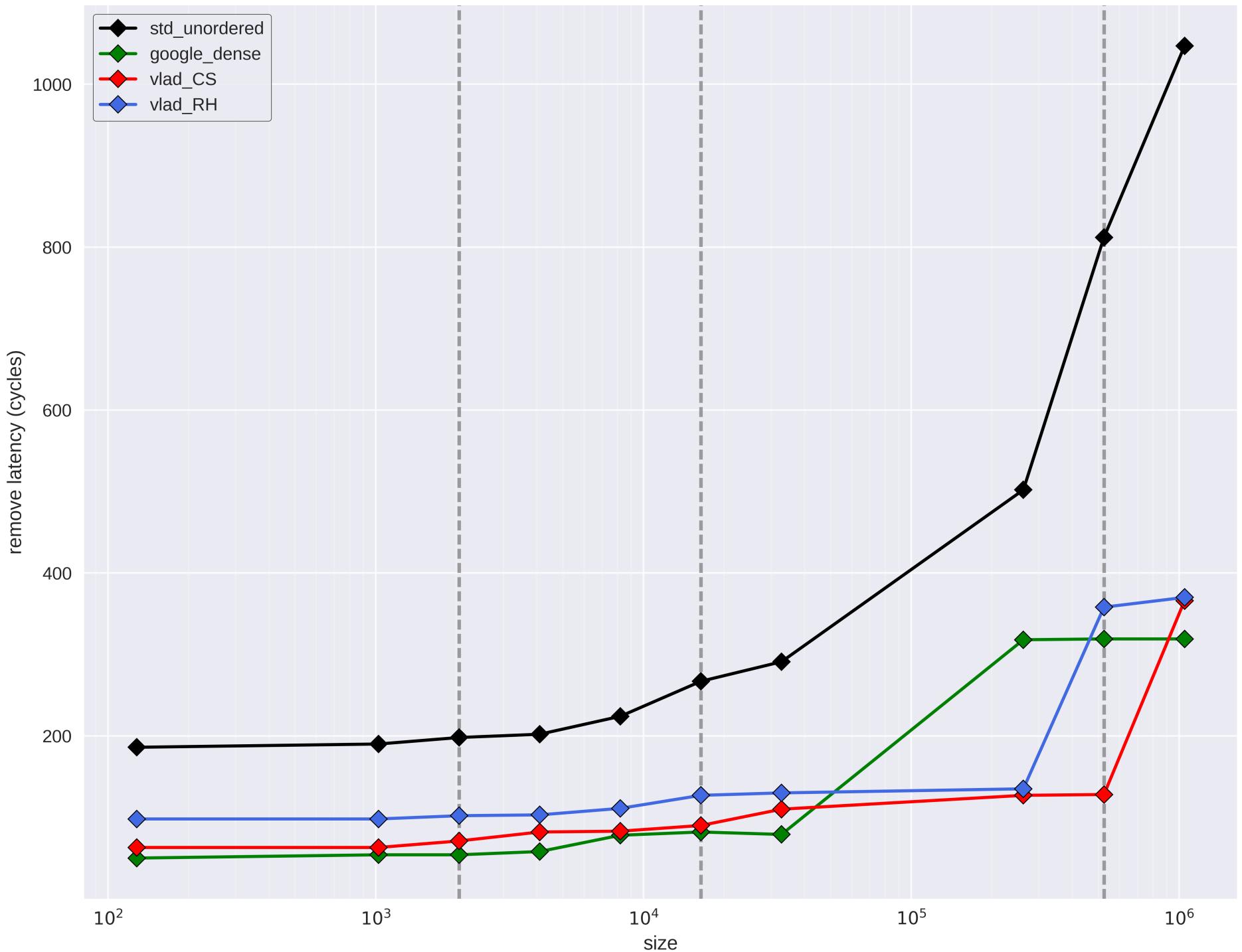
MEDIAN get.miss latency (cycles) vs table size



MEDIAN put latency (cycles) vs table size



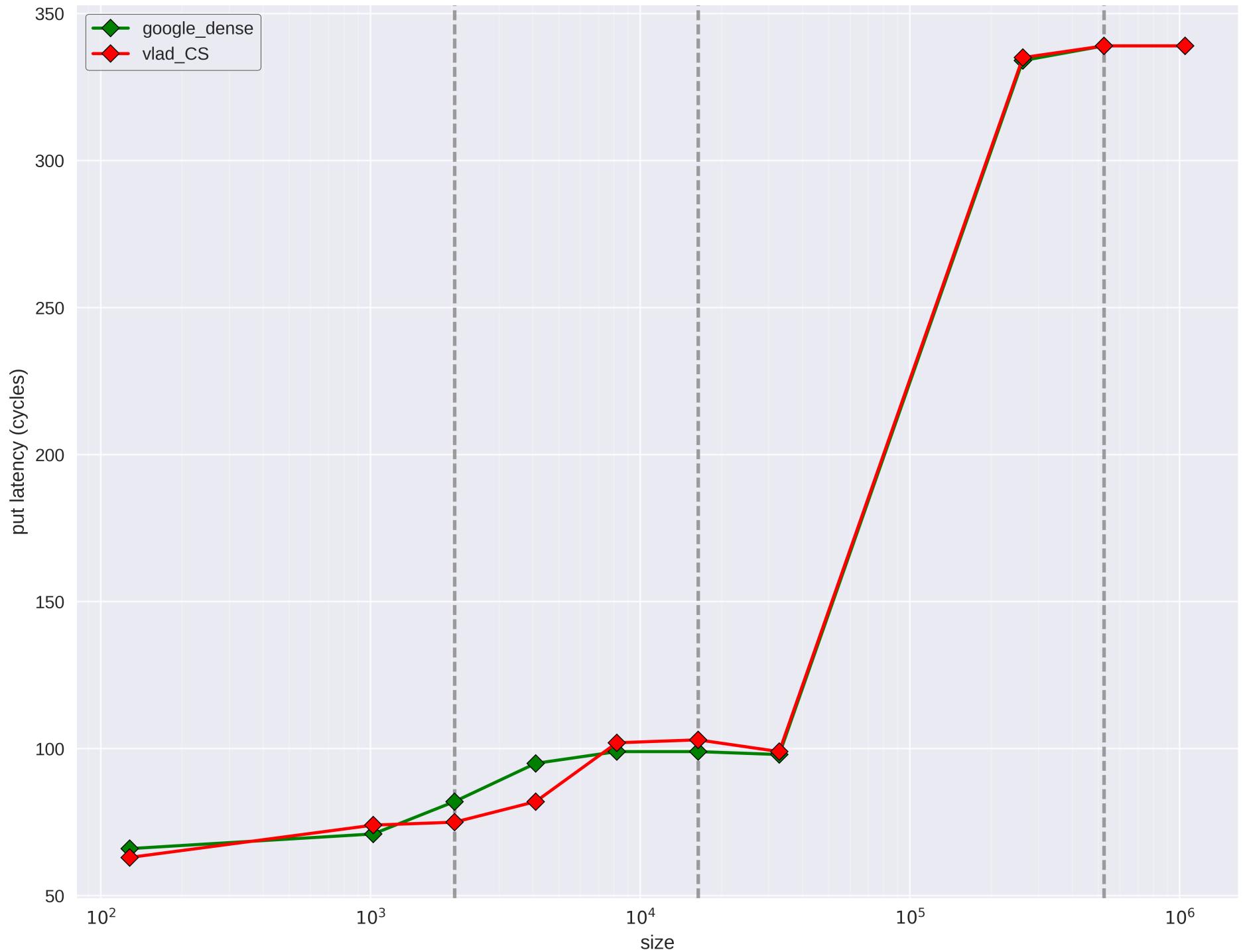
MEDIAN remove latency (cycles) vs table size



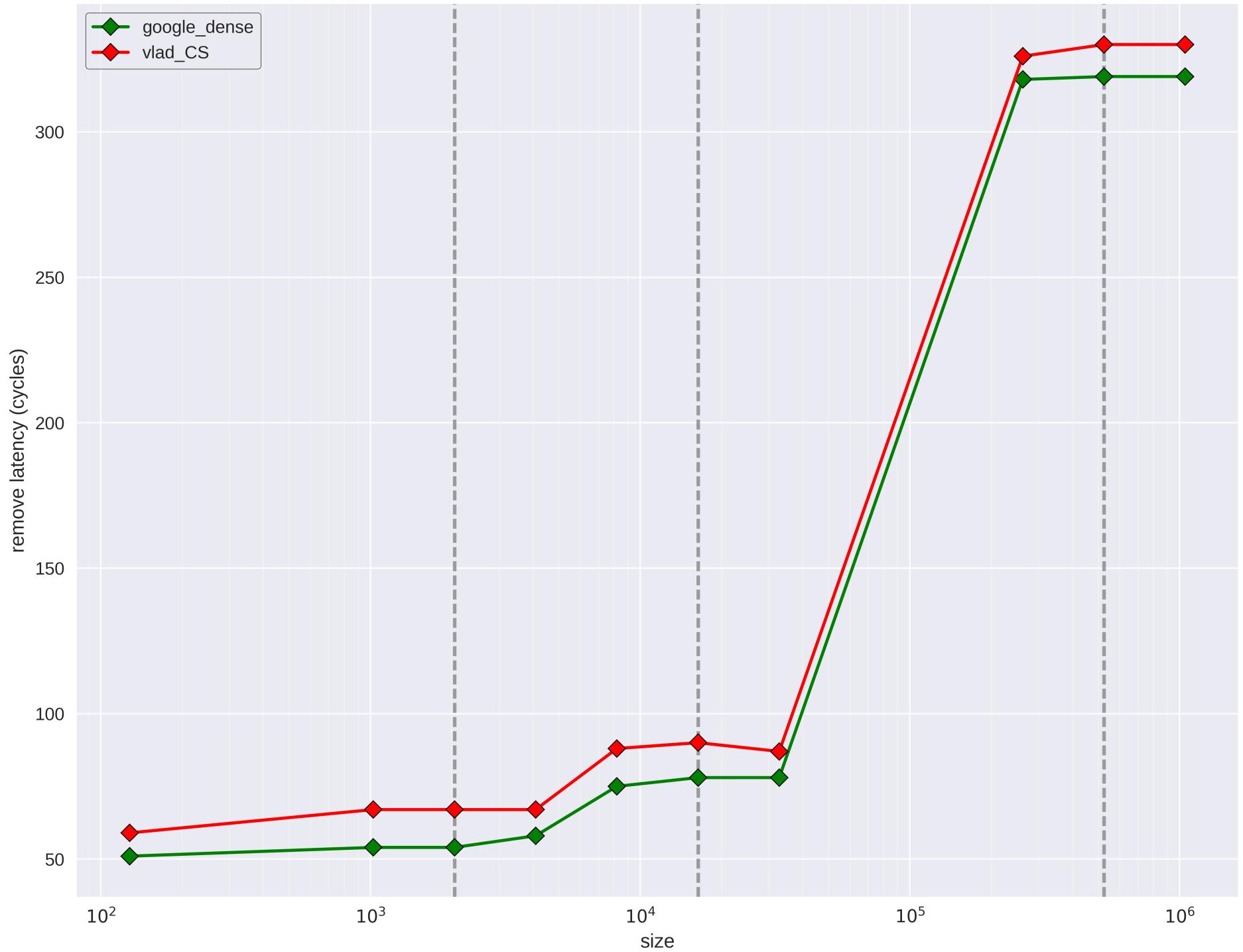
But, but, but! Is google_dense_map cheating?

- Given expected size, chosen load factors:
 - std::unordered_map ~100%
 - google_dense_map **~25%**
 - vlad_CS **~100%**
 - vlad_RH **~50%**
- Let's try mano a mano comparison that is fair (same load factors):

MEDIAN put latency (cycles) vs table size

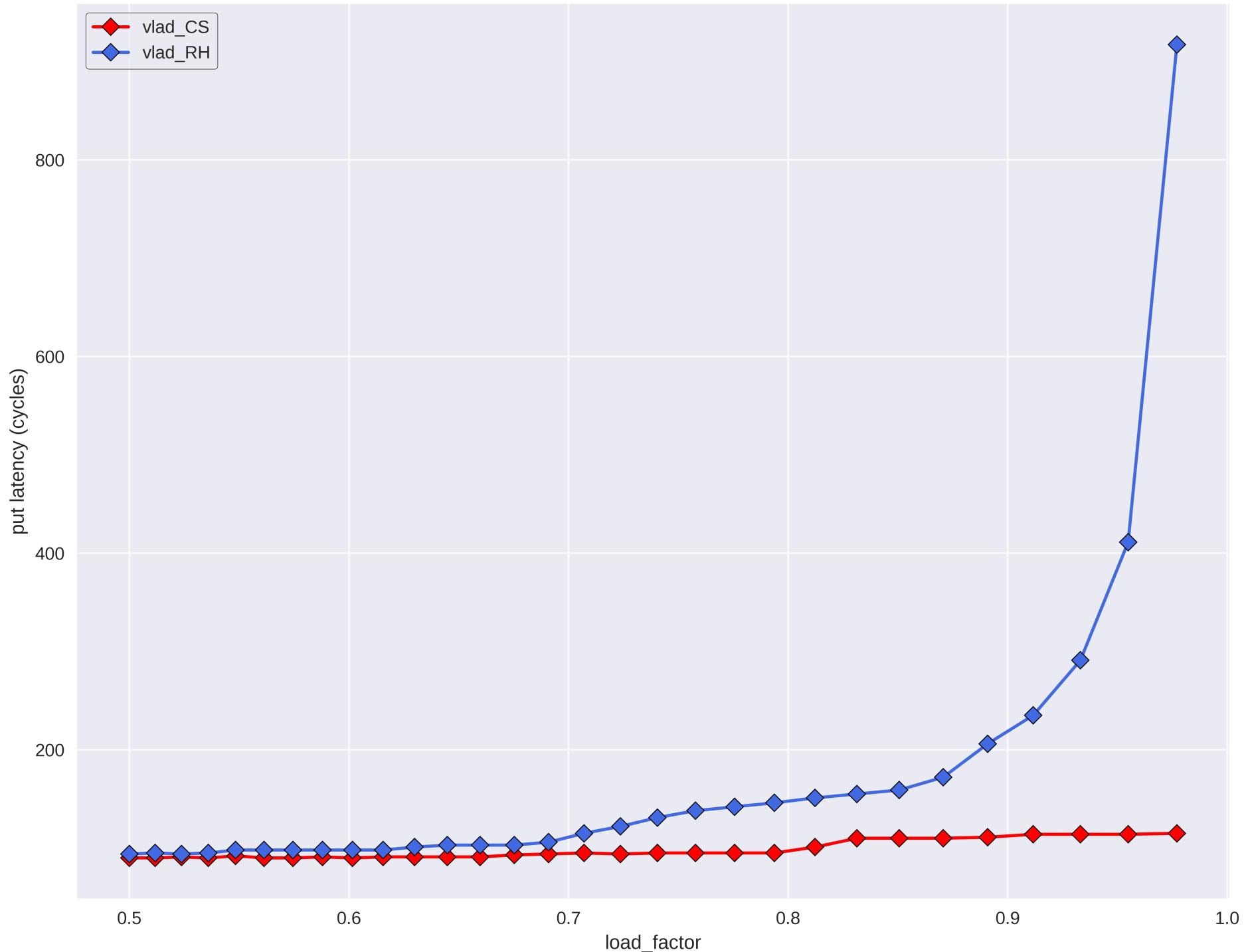


MEDIAN remove latency (cycles) vs table size

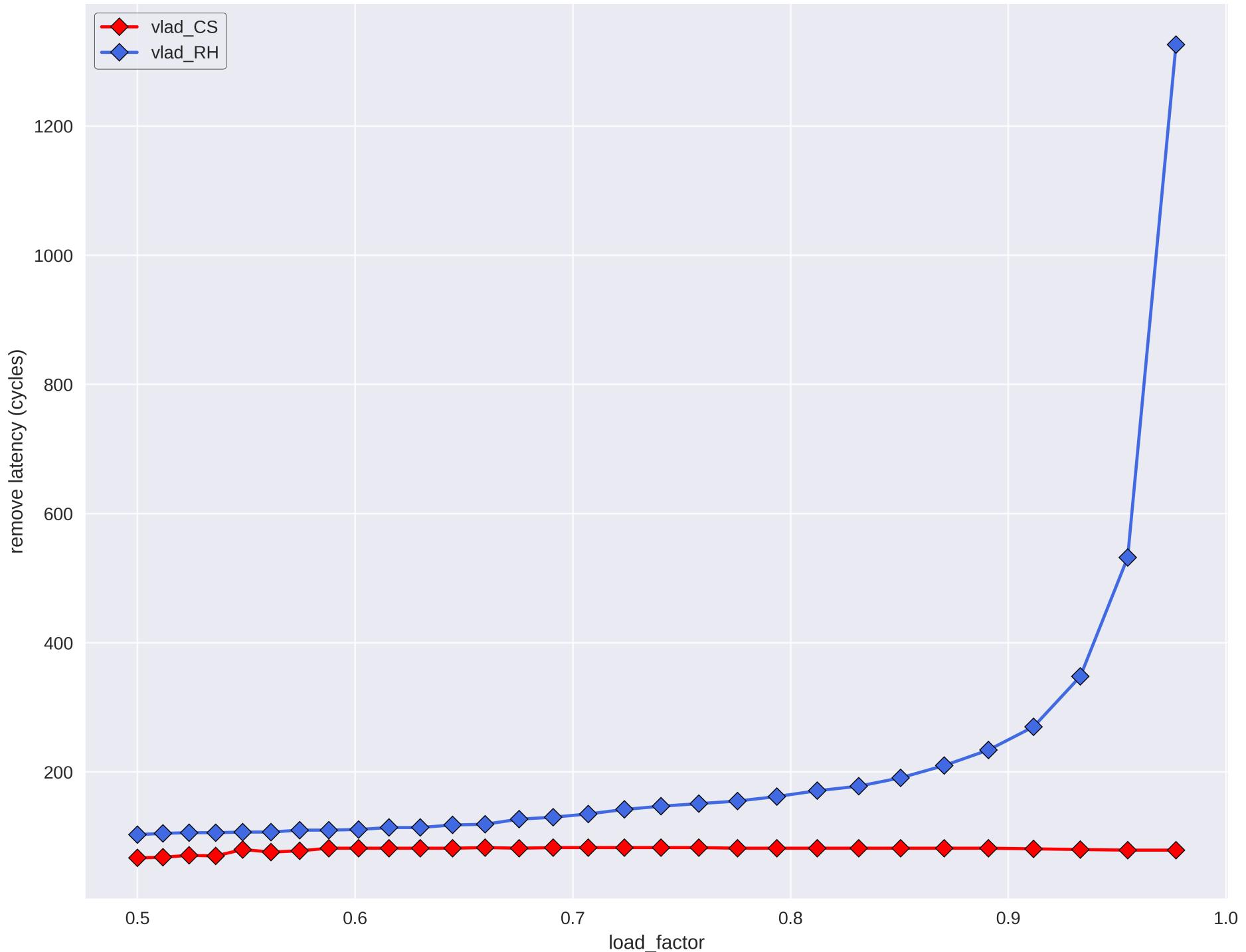


Robin Hood scheme seems to have issues with
nearly full tables...

MEDIAN put latency (cycles) vs table load_factor



MEDIAN remove latency (cycles) vs table load_factor



Lessons:

- Main benefit of **chained scatter** *without* coalescing: simpler and faster deletions
- “Short pointers” and sentinels rock
- Reordering structs for memory compactness rocks
- Power-of-2 sizing still rules unless have a “fast mod” alternative
- Linear probing still sucks
- Back or front shifting not so great in practice – consider tombstones instead