

Tree Cache Learning  
Or, What I Did this Summer

Jack Weinstein  
Argonne National Laboratories

# Normal Cache Behavior

- No Caching
  - Each basket request is a separate file transaction
- Caching
  - Cache misses are file transactions
  - No cache fills until after learn phase
    - Basket requests are separate file transactions while learning

# Motivation

- Current best for learn phase is  $N$  file transactions for each of  $N$  branches used
- Can't make good guesses at branch usage
- Few large reads are less expensive than many small reads
- A single large read is not much more expensive than a single smaller read
  - Latency is the dominating factor
- Goal: reduce file read calls for learn phase

# Testing

- group.test.hc.NTUP\_TOPJET
- ~4000 branches, flat NTuples
- “Large” clusters
- Rewritten
  - Auto-flush 666 entries
  - Baskets sorted by branch
  - Baskets sorted by entry

# Testing

- Files on NFS storage
- ROOT macro reads all entries of tree
- Reads a subset of branches
- Learn Entries left as default 100 (far below first cluster boundary)

# Changes already in ROOT Trunk

- Added TTreeCache::Enable() and Disable()
  - Duplicate / extraneous calls to TTreeCache::ReadBuffer
  - TFile::fReadCache
- Extraneous cache clear / fill after learn phase

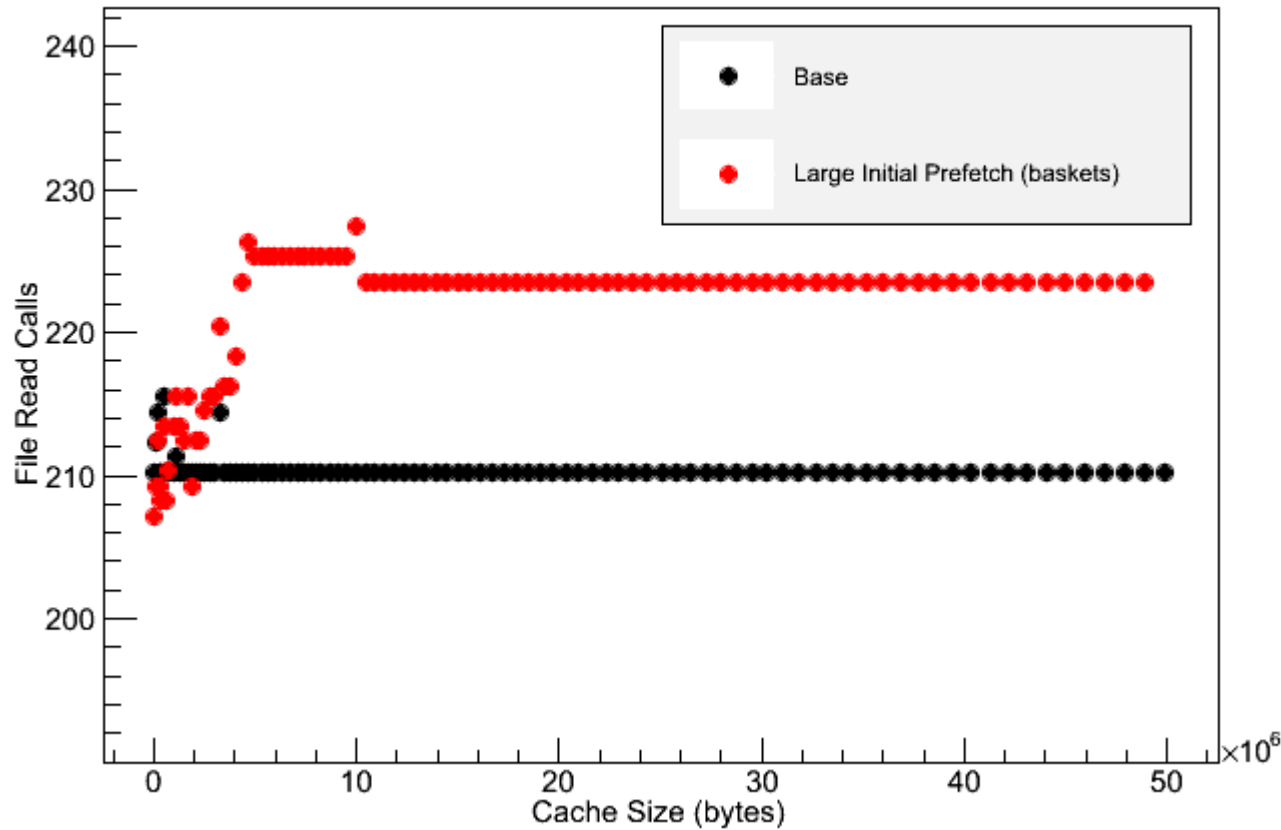
# Learning Phase Strategies

- Large Initial Prefetch
  - Large, single read
  - Data from beginning of Tree
- Neighboring Data Prefetch
  - On basket request, prefetch adjacent data on disk
  - Exploit physical locality of related branches
- By baskets
  - Add baskets similarly to cache fill
- By raw data blocks
  - Read blocks from disk, basket or not
  - On block request, check contained in read block

# Prefetching by Baskets

- Iterate over baskets of tree branches, add to cache
- Works well for cache fill – but not for the learn phase, wide in branches and shallow in baskets
- Small cache compared to branches and cluster concerning
- Too many fragmented reads
- Looks like: raw block size = cache size





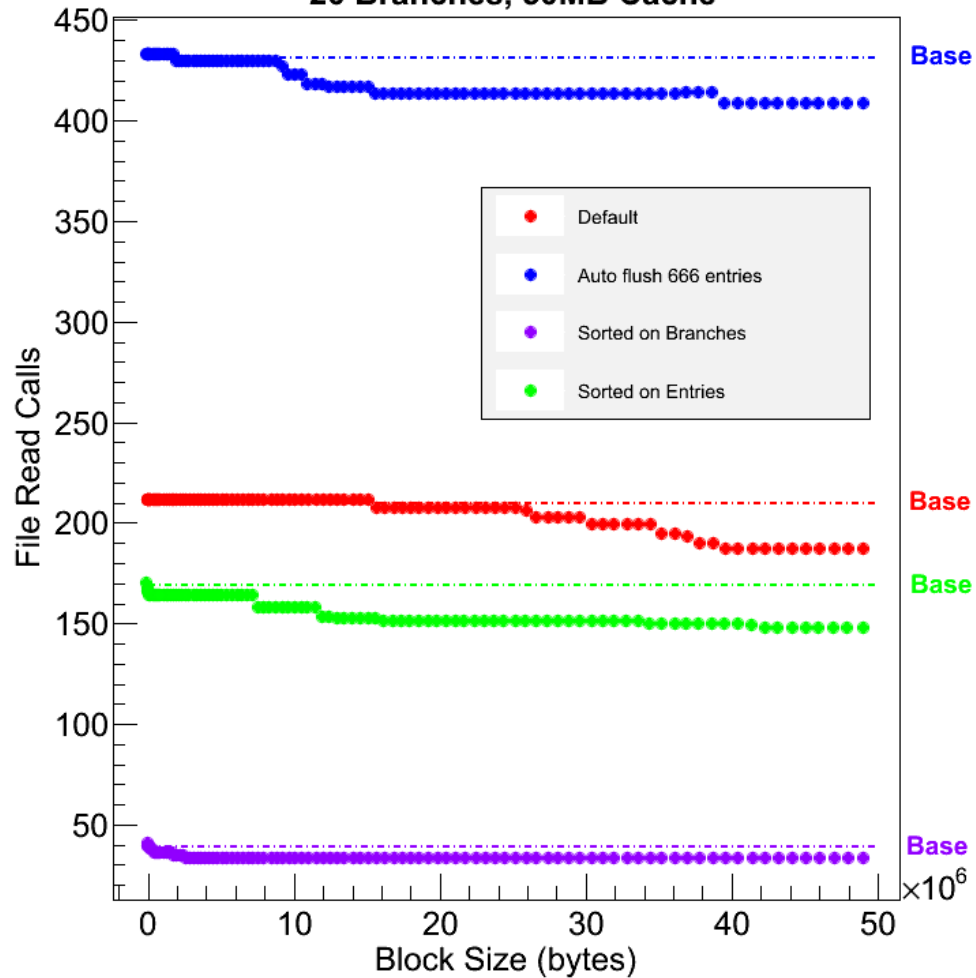
- 20 branches (not random)
- Default basket arrangement
- Base (no changes)
- Large initial prefetch, selecting baskets

# Large Initial Prefetch as a Raw Block

- Read a large block of data from the beginning of tree data
- No sorting, guaranteed single read
- Dealing with “nice” files. Trees are not entangled on disk
- Block size compared to cluster
  - Benefits from small initial cluster
- Possible to grab data beyond learn phase

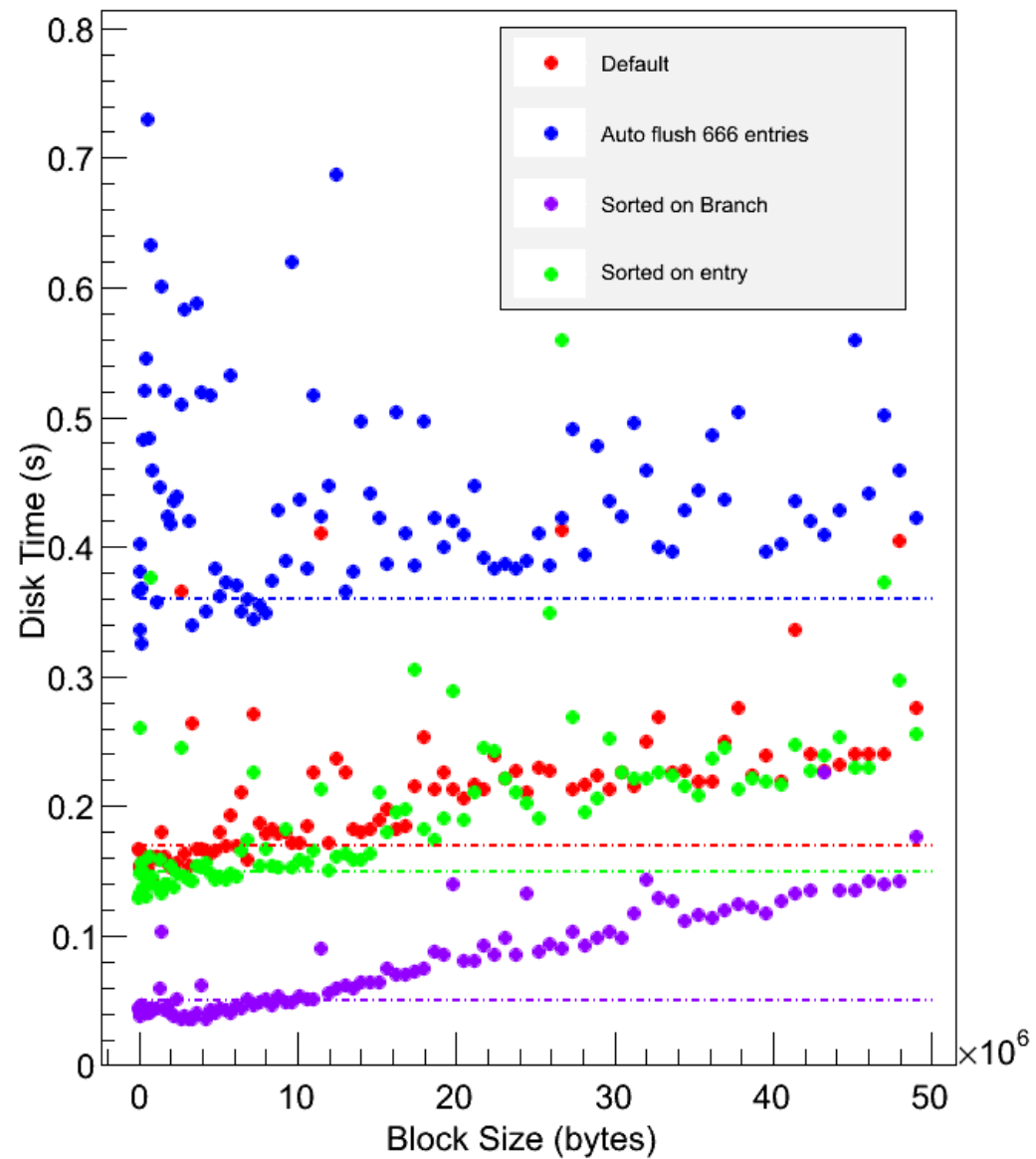
# Large Initial Prefetch (raw)

20 Branches, 50MB Cache



# Large Learn Prefetch (raw)

20 branches, 50MB Cache

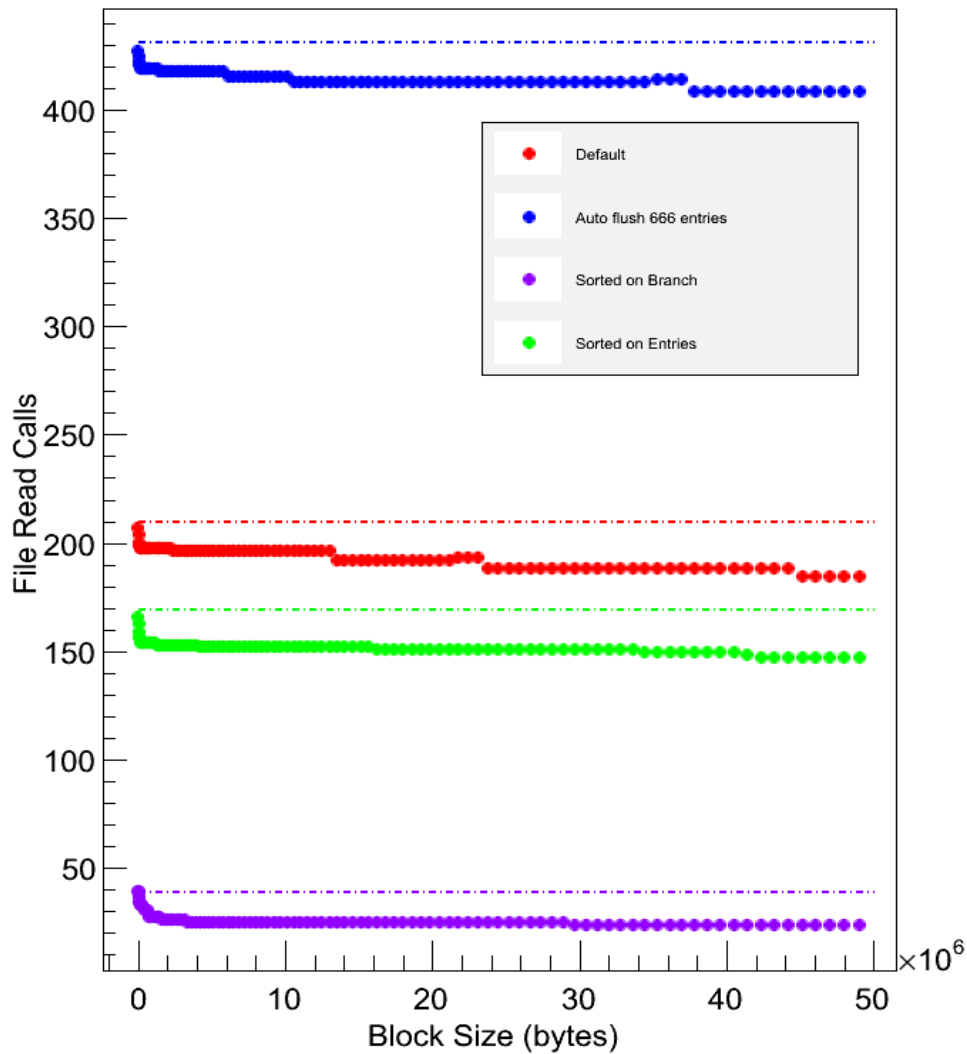


# Neighbor Data Prefetch as a Raw Block

- During learn phase, before cache miss, grab sequential block
- Exploit physical locality of related baskets
- Similar to TFile readahead
  - Don't know next read, no gap to fill
- Smaller blocks are sufficient to reduce reads
- Read overhead increases with branches used

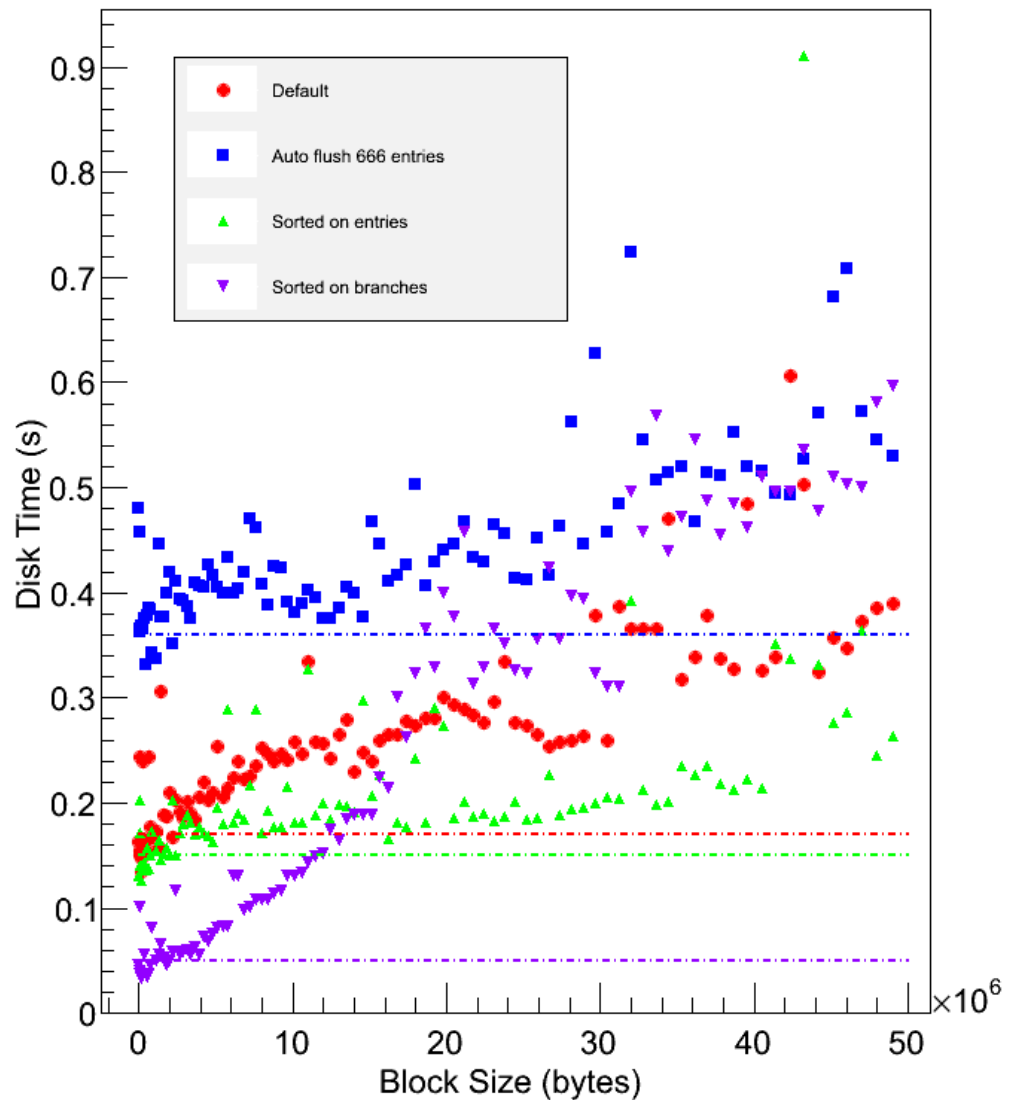
# Neighbor Data Prefetch (raw)

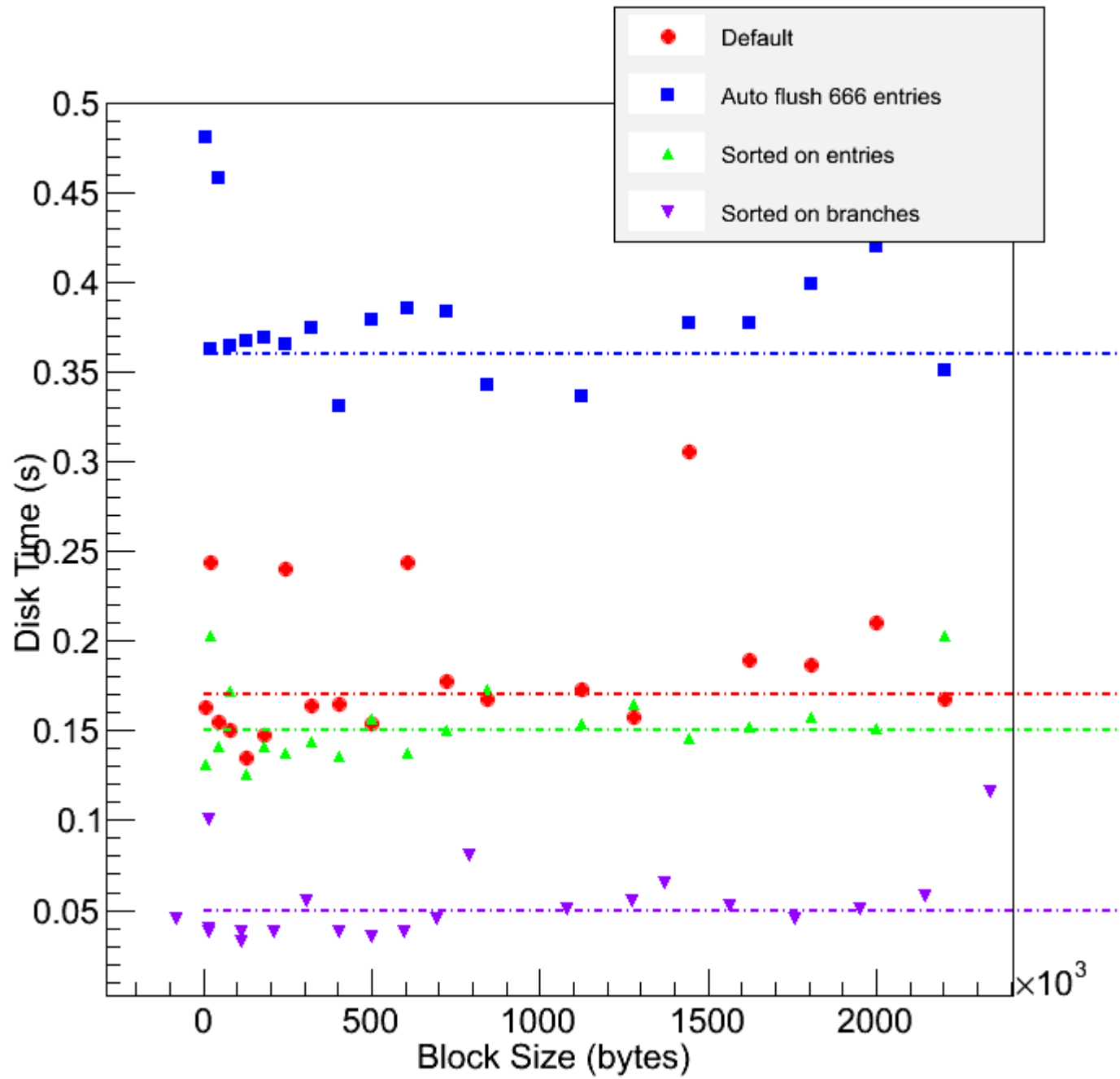
20 Branches, 50MB Cache



# Neighbor Data Prefetch (raw)

20 Branches, 50MB Cache



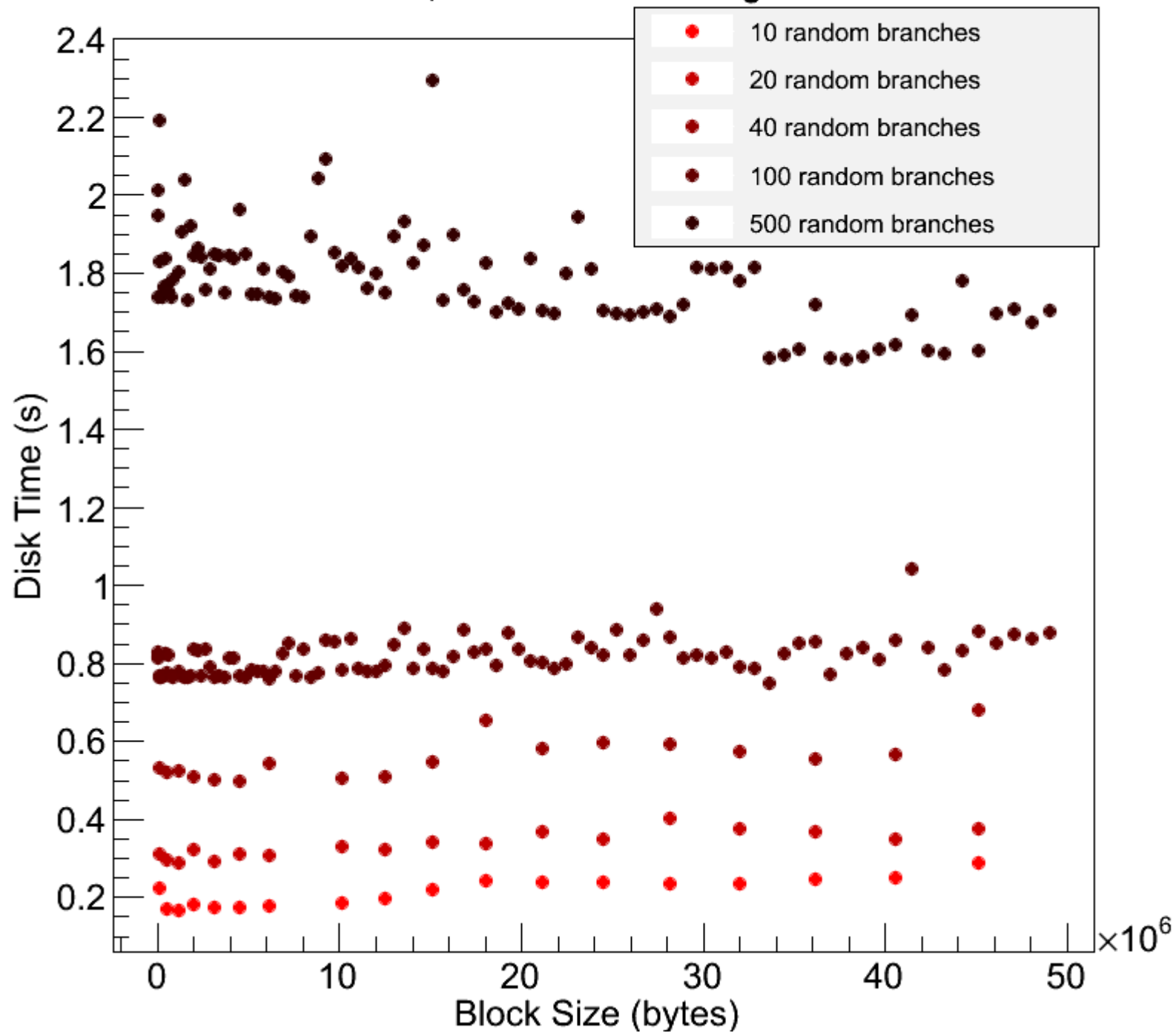


# With More/Different Branches

- Greater number of random branches
  - Read baskets get closer
  - File read calls decrease more sharply
- Neighbor data prefetch makes more overhead reads

# Large Initial Prefetch (raw)

50MB Cache, default basket arrangement





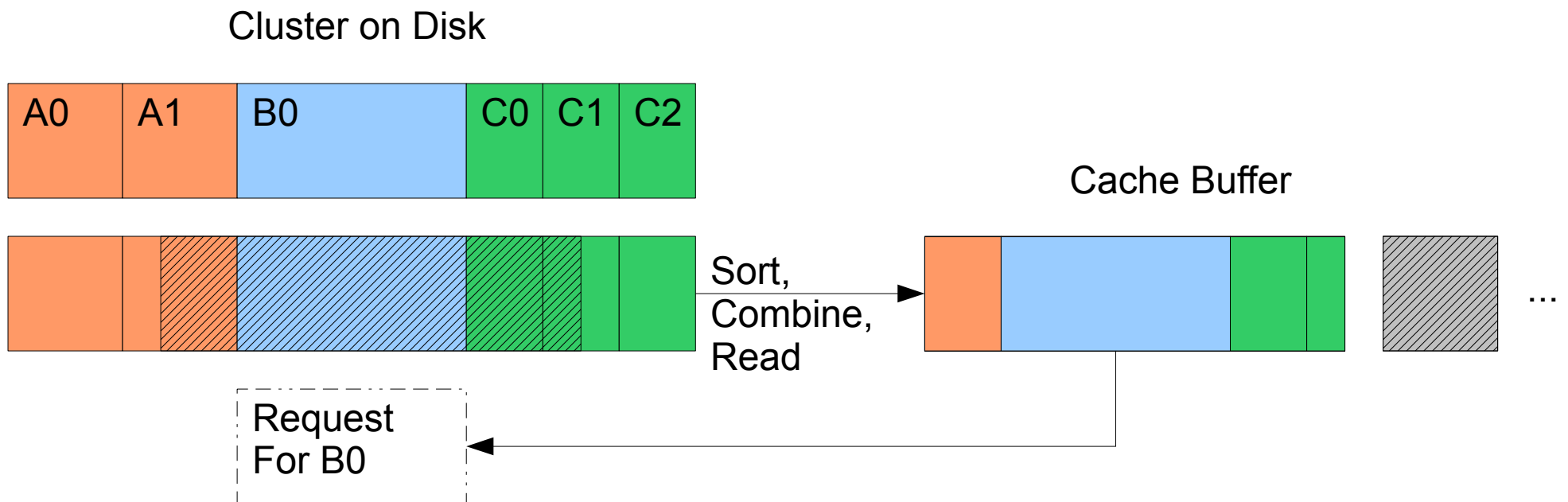
# Conclusions

- Neighbor Data Prefetch works well for small block sizes
  - Sharp decrease in read calls with block size
- Large Initial Prefetch works well for “large” blocks compared to cluster size
  - Constant overhead disk time for fixed block sizes
  - Slower decrease in read calls
- Most cases, trade read calls for disk time



# ReadBuffer Overload

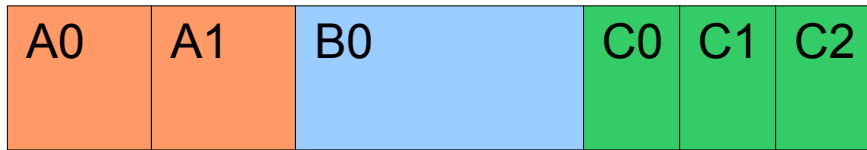
- TTreeCache::ReadBufferExtNormal
  - Overloads TFileCacheRead::ReadBufferExtNormal
  - Extends functionality



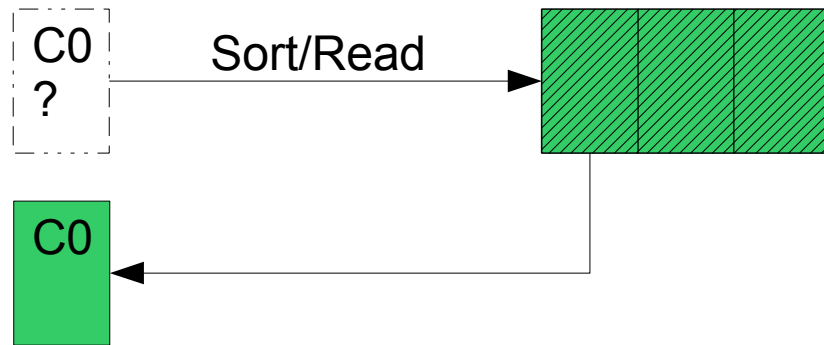
# Afterthought

- It would be nice to be able to read data into the cache without clearing the cache
  - Recycle reads
  - Would work well with neighboring data prefetch
  - Could mix large initial prefetch with neighboring data prefetch

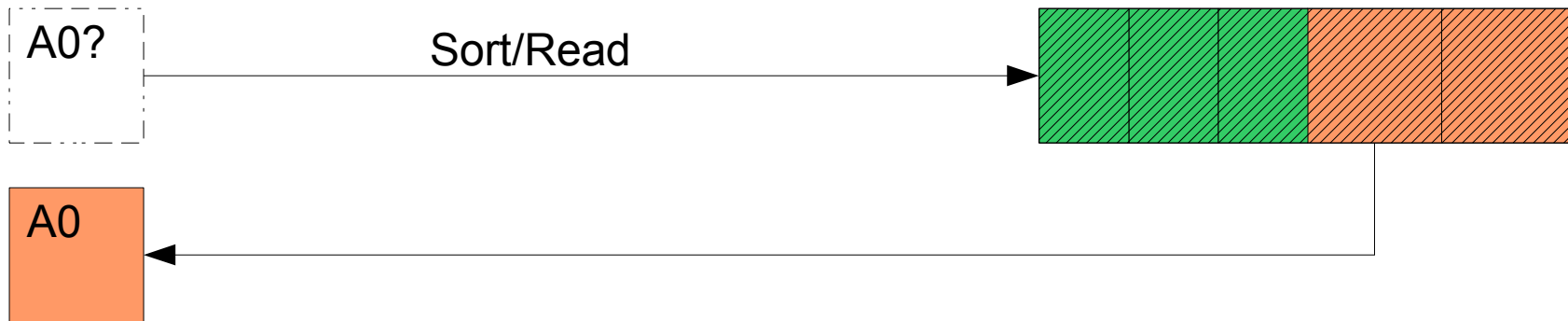
# Cluster on Disk



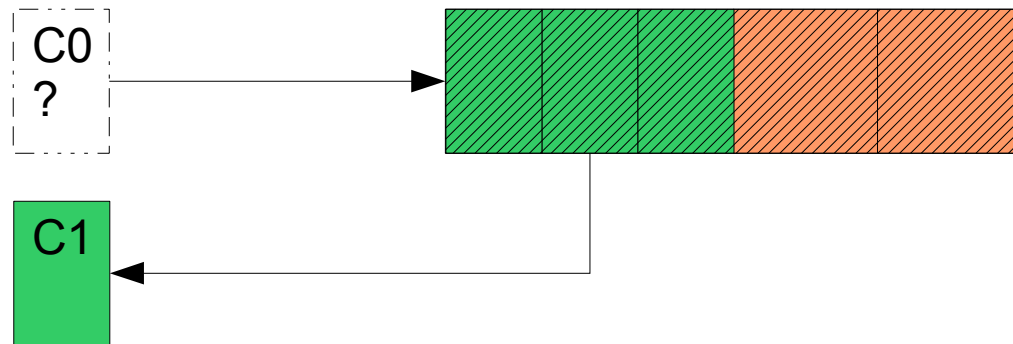
# Cache Buffer



= 1 read total



= 2 reads total



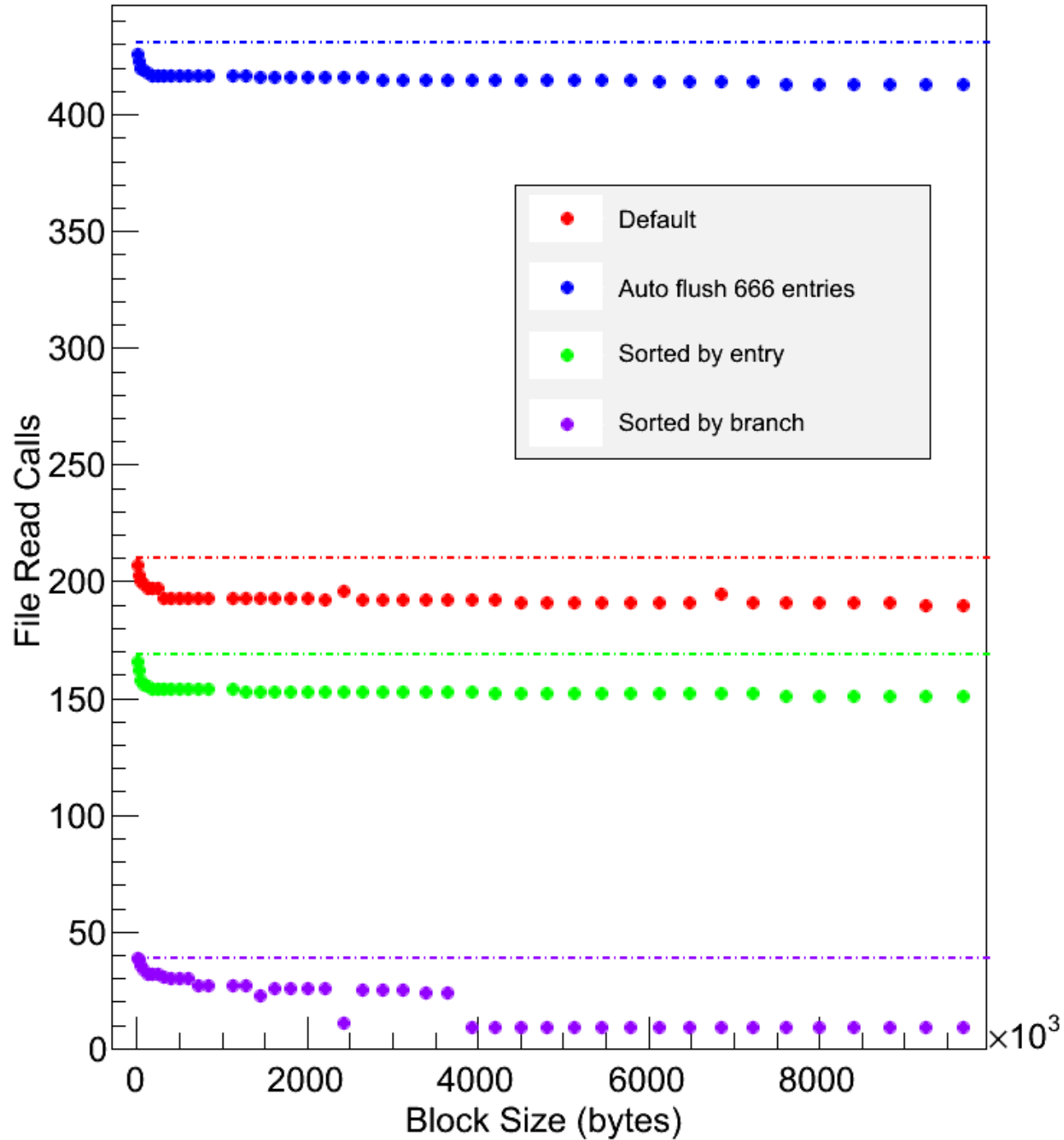
= 2 reads total

# Neighbor Data Prefetch with Cache Modifications

- Don't clear cache (until after learn phase, before cache fill)
- Don't throw away learn phase reads
- Overhead in bytes read is never more than the cache size
- Larger decrease in disk reads
- Slight decrease in overall disk time for small block sizes

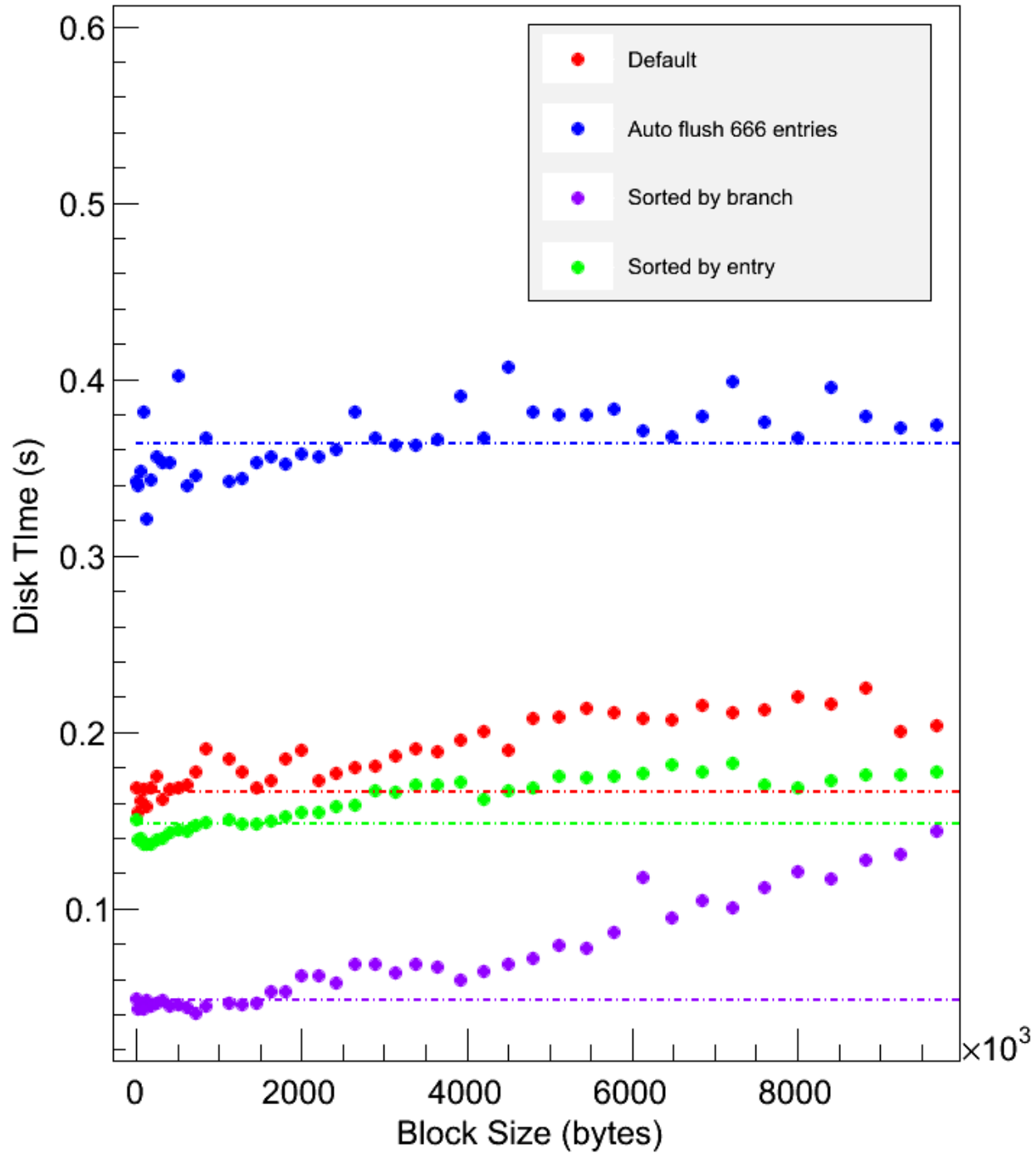
# Neighbor Data Prefetch (raw)

20 branches, 50MB Cache, with Cache Modifications



# Neighbor Data Prefetch (raw)

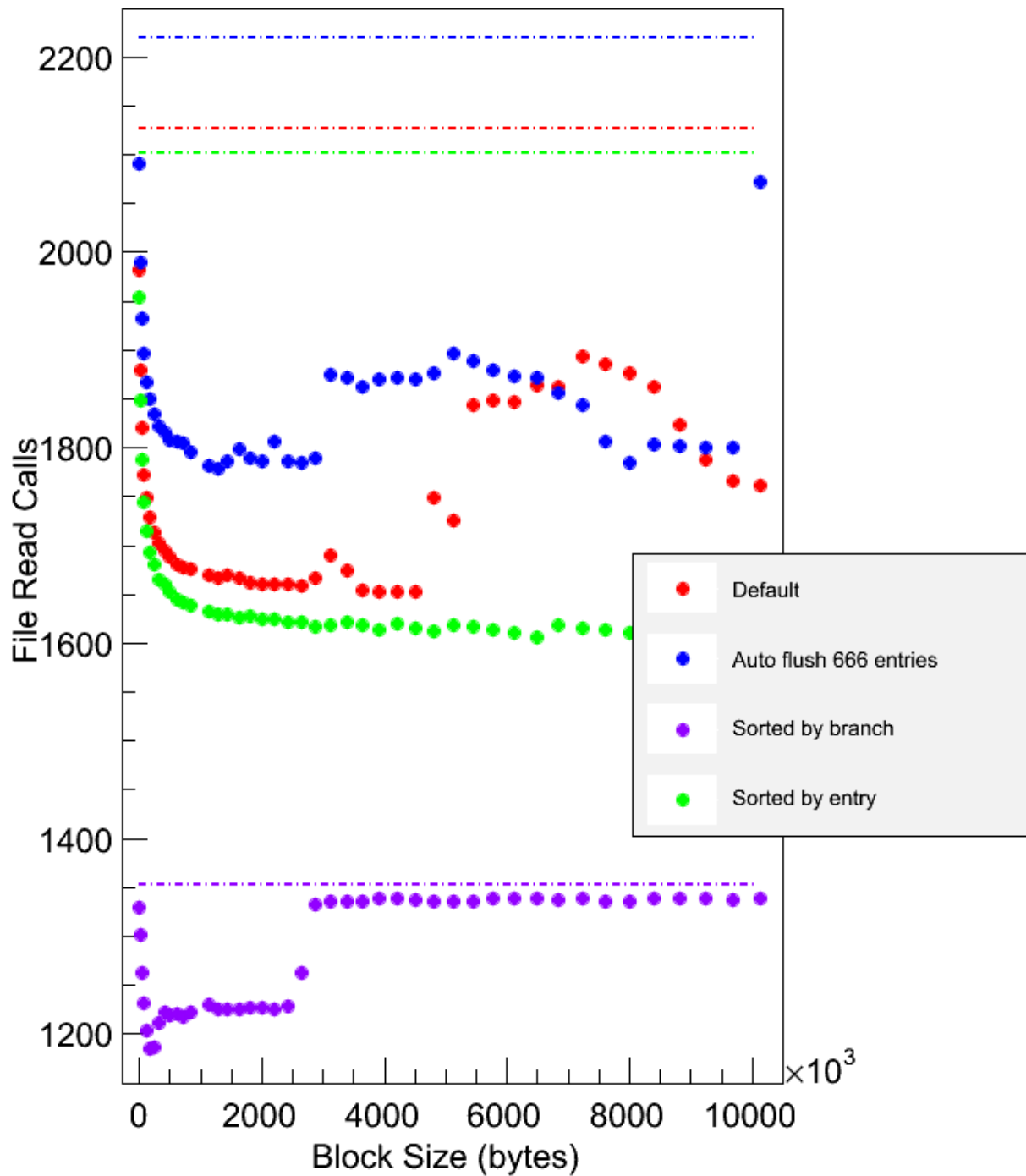
20 branches, 50MB Cache, with Cache Modifications





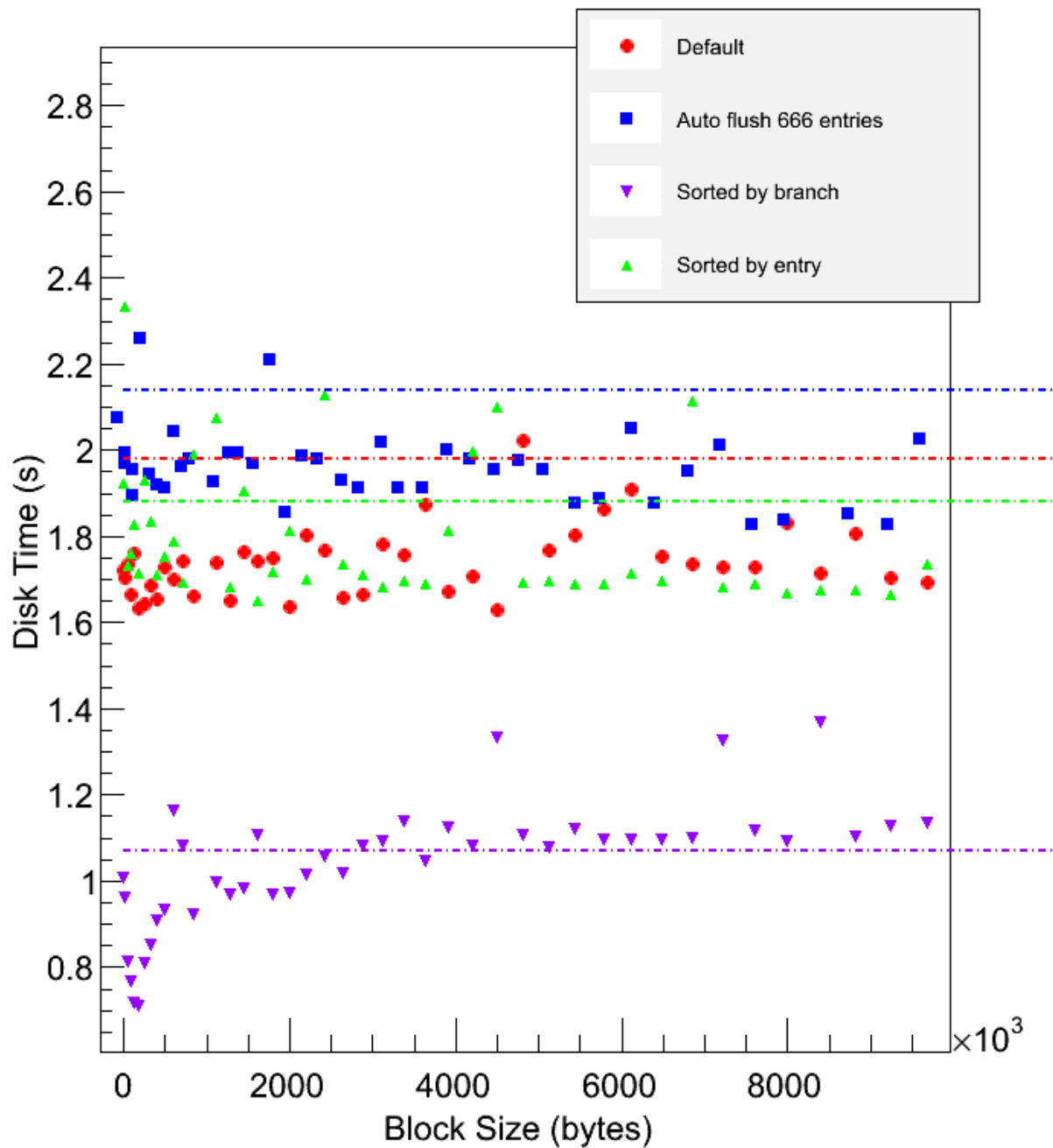
# Neighbor Data Prefetch (raw)

500 random branches, 50MB Cache, with Cache Modifications



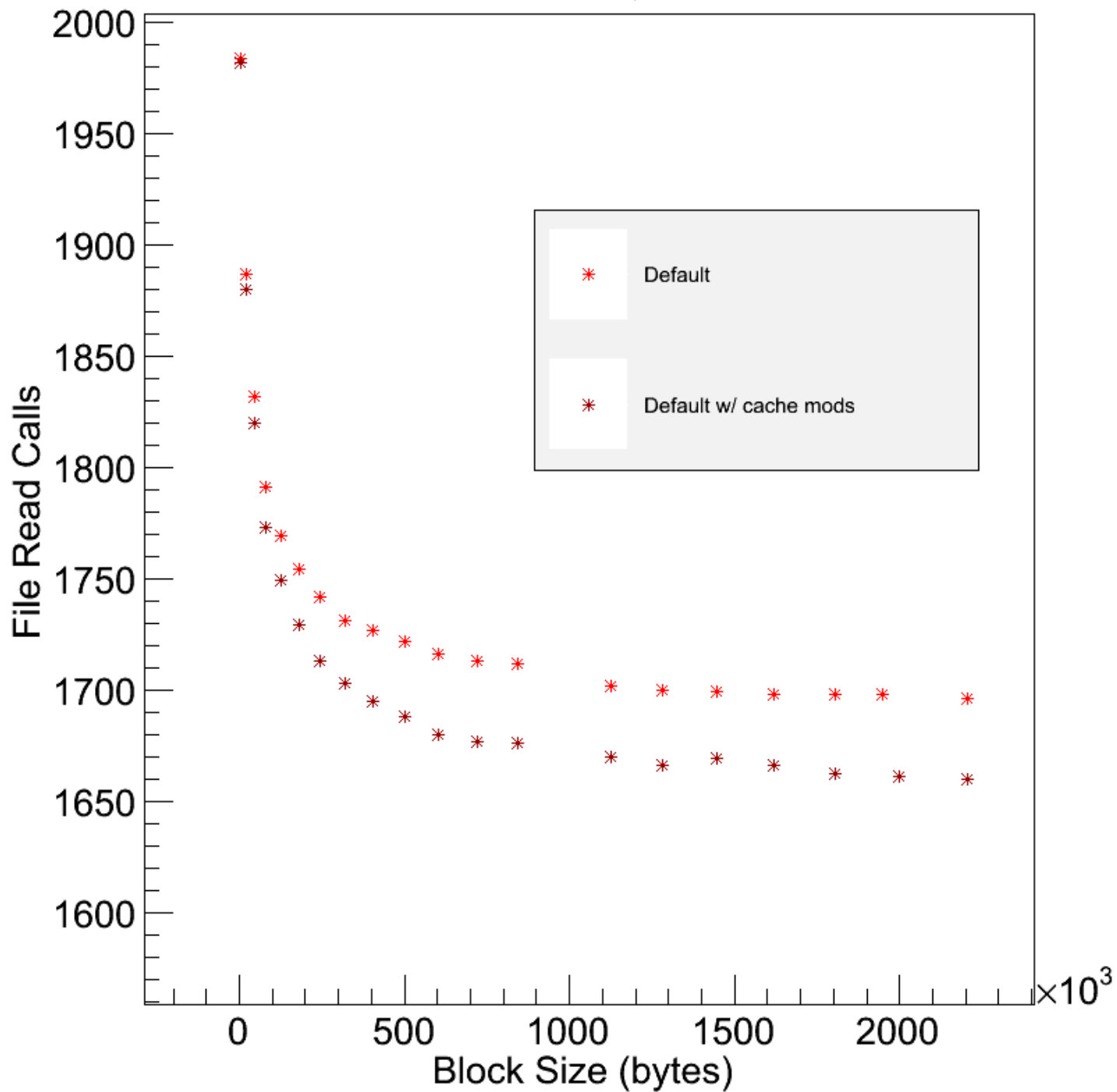
# Neighbor Data Prefetch (raw)

500 branches, 50MB Cache, with Cache Modifications



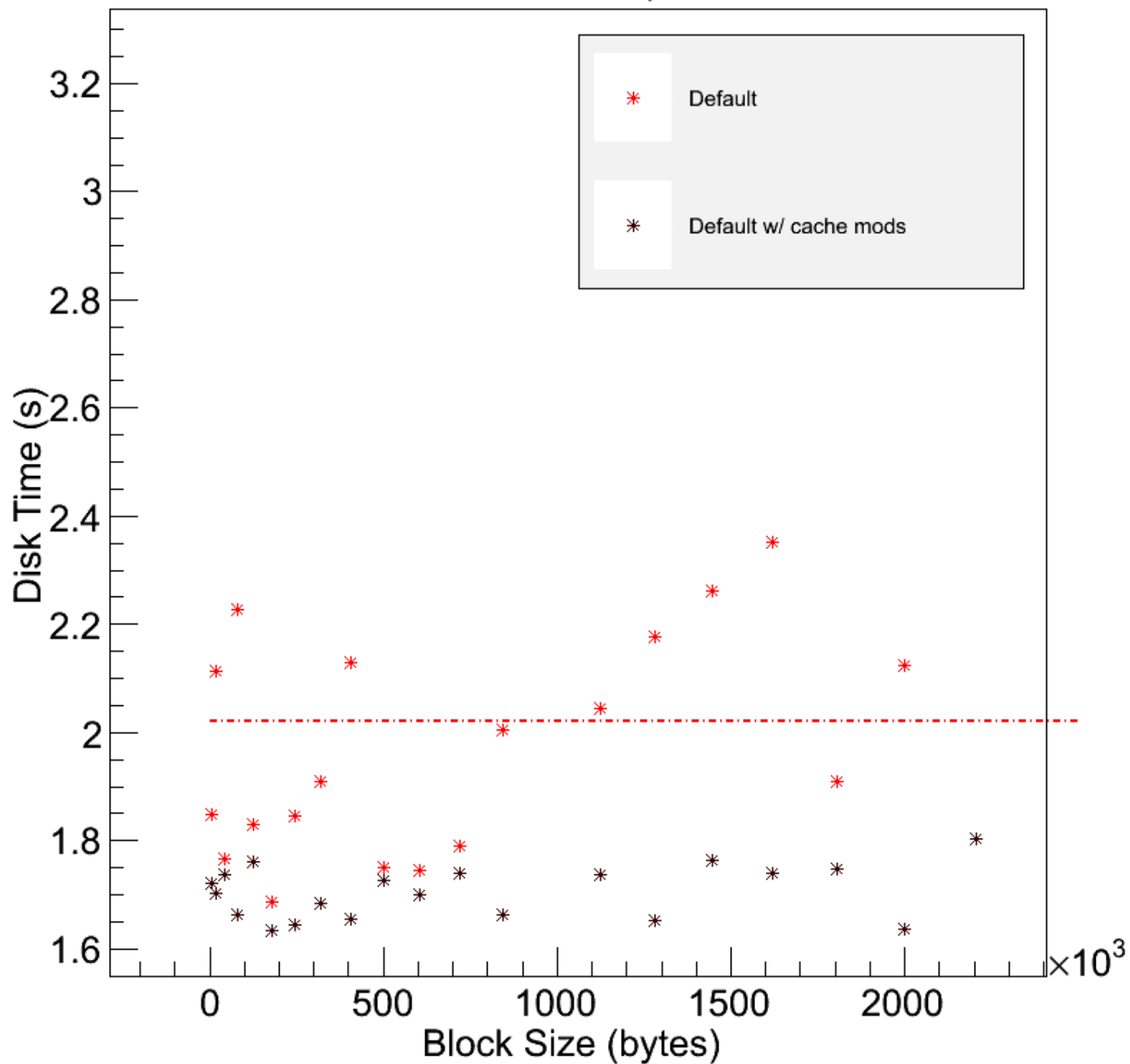
# Neighbor Data Prefetch (raw)

500 random branches, 50MB Cache



# Neighbor Data Prefetch (raw)

500 random branches, 50MB Cache



# Neighbor Data Prefetch (raw)

500 random branches, 50MB Cache

