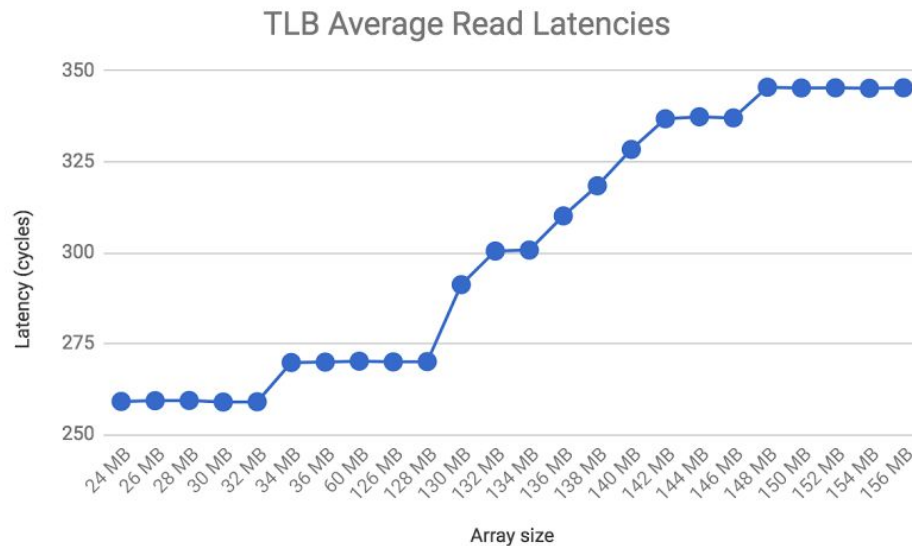


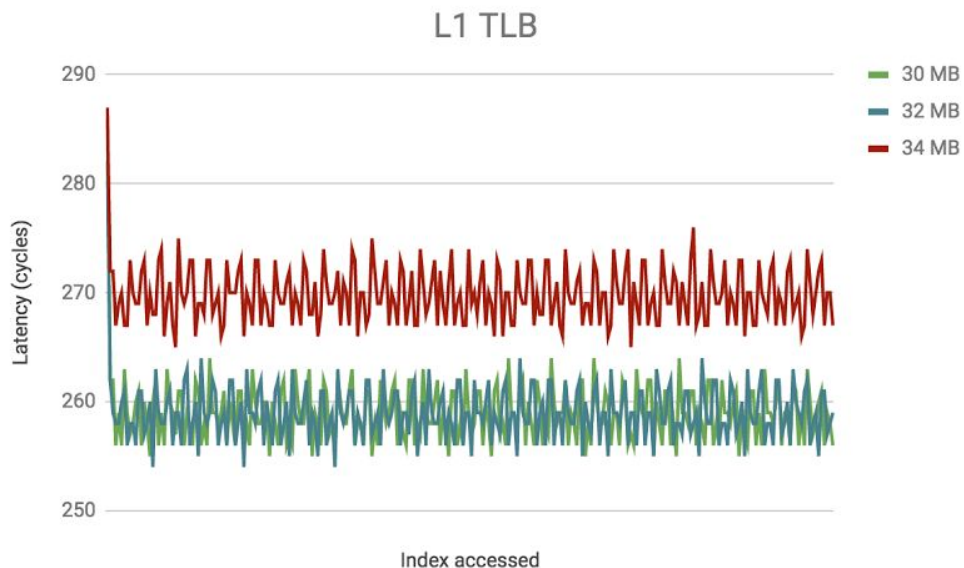
TLB Levels

The TLB could have 2 levels (same as previous generations' TLBs). From the graph below, the plateau from 24 MB to 32 MB is due to L1 TLB hits, the plateau from 34 MB to 128 MB is due to L2 TLB hits, and the plateau from 148 MB to 156 MB is due to L2 TLB misses. (Another hypothesis is that the TLB has 4 levels, which I explain later.)



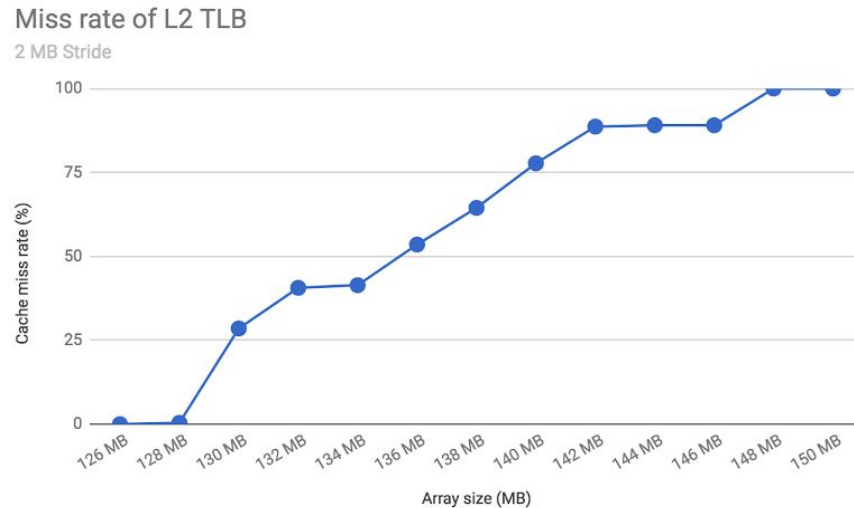
TLB Level 1 Size

The L1 TLB could be 32 MB. From the graph below, array sizes of 32 MB and less have an average read latency of ~259 cycles. An array size of 34 MB triggers L1 TLB misses and L2 TLB hits, with an average read latency of ~269 cycles. Previous generations' L1 TLBs are fully associative, and this one seems to be too because of the relatively smooth and flat graph of latencies.



TLB Level 2 Size

The L2 TLB could be 128 MB. From the graph below, an array size of 130 MB immediately triggers cache misses.



TLB Level 2 Structure

The older L2 TLBs are set associative. However, my graph of the L2 TLB miss rate is not as smooth as the microbenchmarking paper's graph of the L2 TLB miss rate, so the structure of the newer L2 TLB is probably more complicated. This is the paper's original graph:

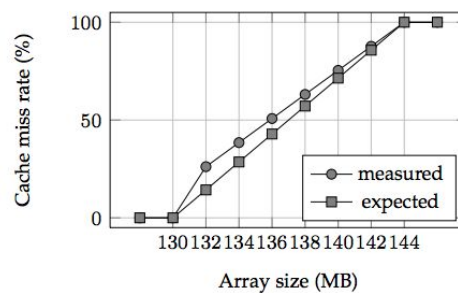


Fig. 8. Miss rate of L2 TLB (2 MB stride).

From the "measured line" above, they reasoned that there were 7 sets in the L2 TLB, and the size of the 1st set was larger than the other 6 sets.

From the graph of the L2 TLB miss rate, I think there are 10 sets.

There are 6 sets that seem to be the same size:

- The cache miss rates for array sizes of 134 MB through 142 MB increases smoothly/linearly, so those 4 sets must be the same size.

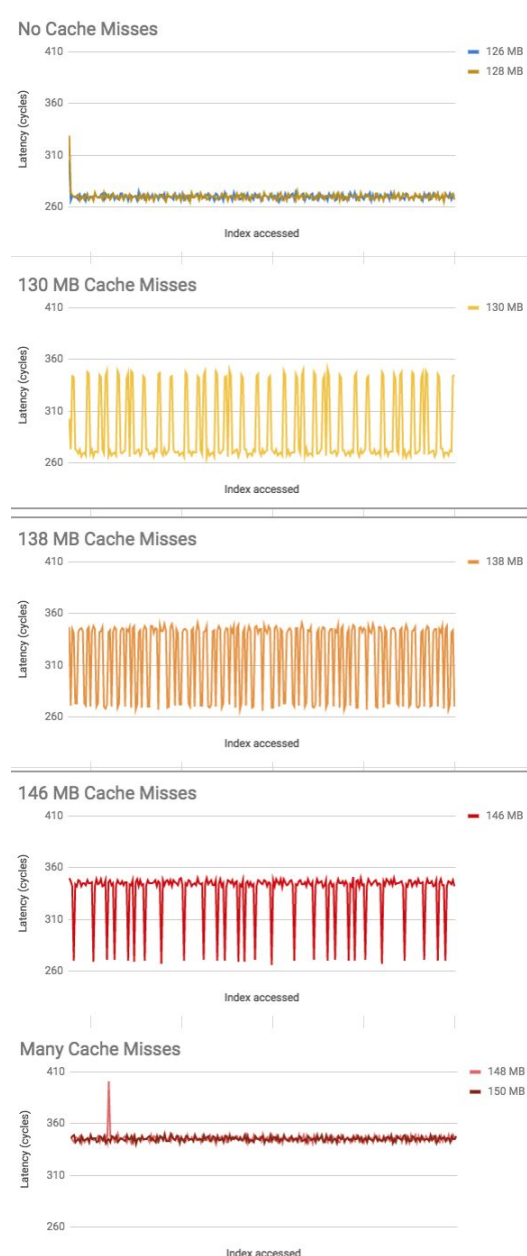
- The cache miss rates from 130 MB to 132 MB and from 146 MB to 148 MB increase at the same rate, so those sets must be the same size as well.

Finally, the cache miss rate from 128 MB to 130 MB is a larger difference, so the first set must be larger than the other sets.

Something I'm not sure about is the plateaus - the cache miss rate differences from 132 MB to 134 MB, and from 142 MB to 146 MB, are flat. There are 2 explanations I can think of:

- The TLB actually has 4 levels. This means that the 4 levels respectively have 1 set (fully associative), 2 unequal sets, 4 equal sets, and 1 set (fully associative).
- It has something to do with the size of the sets in the L2 TLB, but I'm not sure how it works.

I also graphed the read access latencies over each index so we can see the hit and miss patterns:

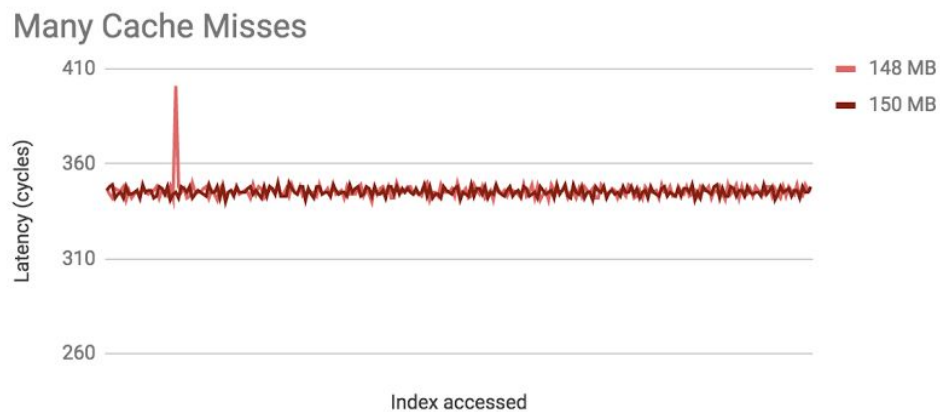


All together, the access patterns look like:



Random Latency Spikes

There are sometimes random spikes of higher latencies. I ignored these as just random anomalies for now. For example, in the graph below, the read latencies for arrays of size 148 MB and 150 MB are mostly smooth, but there's a random spike for 148 MB.



Data + More Graphs

All my data tables and graphed results are [here](#). (In the Google Sheets, there are 3 sheets: L2, L1, and Levels.)