GPU-accelerated k-mer analysis in MetaHipMer

Prashant Pandey University of Utah September 2022

MetaHipMer2 Metagenome Assembler

- Part of ECP Exabiome project
- Original MetaHipMer (v1) was released in 2017 (written in UPC/MPI)
- Released version 2 of MHM on Sept 30th 2020
 - Entirely rewritten in UPC++
 - Runs up to 10x faster than MHM v1 and uses 2x less memory
- Increasing use of GPUs

Science Analysis with Metagenomics



What happens to microbes after a wildfire?



What are the seasonal fluctuations in a wetland mangrove?



What are the microbial dynamics of soil carbon cycling?



How do microbes affect disease and growth of switchgrass for biofuels



Combine genomics with isotope tracing methods for improved functional understanding



JGI-NERSC-KBase FICUS (Facilities Integrating Collaborations for User Science) call called out MetaHipMer

Metagenome Assembly



microbial genomes (1000s)

- varying abundance (frequency)
- sequence depth (how many times sequenced, e.g. 50x)

reads

- typical length for short reads 150-250
- error prone, e.g 0.24% per base
- number of reads for a genome dependant on sequence depth and abundance

contigs

- contiguous sequences
- the longer the better
- the fewer errors the better

The MHM2 pipeline



Walk contig graph (iterate)

MHM2 Performance on Summit

- Using GPUs for local assembly and CPUs for kmer analysis
- kmer analysis is most of runtime (up to 45%)



Can we speed up k-mer analysis phase using GPUs?

k-mer counting



- The size of the raw sequencing data makes the problem challenging
- *k*-mer counts follow *highly skewed distributions* making the problem computationally intensive

Two optimizations for the GPU k-mer analysis phase

- Reducing communication volume across nodes using supermers to speed up
- Reducing memory usage on GPU nodes using mapping filters

Two optimizations for the GPU k-mer analysis phase

 Reducing communication volume across nodes using supermers to speed up



Reducing memory usage on GPU nodes using mapping filters

GPU hash table

- Use open addressing
 - fixed size, determined at start by using max available memory to minimize load factor/collisions
 - quadratic probing allows high load factors
- Consists of two arrays (need to be separately initialized):
 - keys: array of uint64_t, containing k-mer packed into 4 bases per byte
 - values: array of 8 uint16_t for the counts of each of the 4 extensions on both sides and one uint32_t for the count of the kmer
- Empty slots indicated by last position of KEY_EMPTY = 0xff...ff
 - always use odd values of k, and each nucleotide is packed into 2 bits, so last two bits will never be used
- Locked slots indicated by last position of KEY_TRANSITION = 0xff...fe
 - also never overlaps with packed k-mer

k-mer analysis for creating distributed de Bruijn graph

- Our goal is to build a de Bruijn graph as a distributed hash table
- k-mers come from reads and from contigs in previous rounds
- Don't use singletons (likely errors)
- Each vertex is a compressed k-mer with extensions and k-mer count



k-mer analysis for creating distributed de Bruijn graph



Splitting reads into compressed supermers



concatenate reads into strings separated by underscores _

each thread processes one position and computes target from (revcomp) kmer minimizer

(+ CPU spins on upcxx::progress +)
output is device array of target for
every kmer (-1 for spaces)

each thread processes one position and builds a supermer if the start position target is different from the left target position

array with offset, length, target for each supermer slot obtained by atomic add

input string compressed separately, one thread per 2 positions (compressing two ACGTNacgt chars into one byte (lower case = low quality)

Inserting k-mers from reads

- Reads are split into compressed supermers on the GPU and then passed to the CPU for network transfer
- The network transfer is more complicated than this (hierarchical three-tier aggregating store)
- The local hash table is not the same as the final hash table (which stores only a kmer plus left and right extension characters - A,C,G,T,X,F)



Processing supermers at the receiver



k-mer count: integer left ext counts: array 4 ints right ext counts: array 4 ints concatenate supermers into string separated by underscores _

each thread processes one position and computes (revcomp) compressed kmer plus extensions

low quality extensions are set to 0

each thread inserts k-mer into hash table on GPU

Performance of k-mer analysis

- Dataset is WA0 (71GB) on 32 Summit nodes
- Total speedup from GPUs is 2.9x
- Parse and pack (PnP) gives 29% speedup
- Adding GPU hash tables gives another 56%
- Adding supermers gives another 44%
- Supermers actually *slow down* CPU version by 55% (mostly the cost of extracting supermers in RPCs on receiving processes)



Performance of pipeline

Speedups from GPUs



- k-mer analysis is now 2.7x to 2.9x faster
- k-mer analysis is now less than 25% of runtime (down from max 45%)
- overall speedup is 70 to 86% on up to 256 nodes
- GPU computation is a few percent of k-mer analysis time

Communication





- supermers reduce communication volume by 5.3x and number of messages by 3.7x
- communication reduction results in 44% speedup for k-mer analysis
- most communication is now alignment reduce in future through earlier read shuffling

Two optimizations for the GPU k-mer analysis phase

- Reducing communication volume across nodes using supermers to speed up
- Reducing memory usage on GPU nodes using mapping filters



K-mer analysis phase



K-mer analysis phase



Singleton k-mers are erroneous and need to be cleaned out.

K-mer analysis phase

- Hash tables store all k-mers but discard new k-mers once they are full

 (load factor = 1)
- GPUs have limited memory (16GBs on V100)
- Lots of wasted space due to singleton k-mers
 - Each k-mer can take up to 18 bytes
- Non-singletons k-mers end up being dropped resulting in poor assembly quality.

A filter is an approximate dictionary



A filter supports *approximate* membership queries on *S*.

A filter guarantees a false-positive rate



False-positive rate enables filters to be compact

space $= \Omega(n \log |U|)$ space $\geq n \log(1/\epsilon)$ Small Large Filter Hash table

False-positive rate enables filters to be compact

space
$$\geq n \log(1/\epsilon)$$

Small
Sm

For most practical purposes: $\varepsilon = 2\%$, a filter requires ≈ 8 bits/item

Using filters to remove singletons



Mapping filter on GPU

K-mer analysis with filter



Mapping filter on GPU

We use a quotient filter^[Pandey et al. SIGMOD 2017] as a mapping filter. It can approximately map a key to a small value. Key is k-mer and value is the extensions.

Two choice filter on the GPU

- The two choice filter is based on **power-of-two-choice** hashing.
- We can perform insert/query/delete operations using only **2 cache line probes** and **1 cache line write.**
- Cooperative groups is used to perform operations inside each block.
- Both point and batch operations supported.
 - Batching operations helps to speed up the throughput



K-mer analysis results using filters



Fig. 6: Memory usage of the GQF and hash table in MHM for WA and Rhizo datasets.

K-mer analysis results using filters

Dataset	Method	Nodes	GQF mem	HT mem mem	GPU insert time	k-mer analysis time	N50 quality
WA	GQF GQF No GQF	32 64 64	96 96 0	522 522 2281	13.5 6.2 4	236.5 125 123	1336 1336
Rhizo	GQF GQF No GQF	14 24 24	35 35 0	131 131 795	13.7 8.4 3.9	151 76 72	901 910 910

TABLE VI: GQF and Hash table memory usage and running time during MetaHipMer runs across two datasets. Note: the N50 quality score for the WA dataset is not available as the assembly process failed at a later stage after the k-mer analysis phase. Memory is in GB and time is in seconds.

Summary

- Accelerating k-mer analysis with GPUs is complicated
- Gives a 2.5x to 2.9x speedup
- Can **reduce the memory usage by 2x** using advanced filter data structures
- Open questions:
 - Is there a better GPU hash table that can offer consistent high performance?
 - Can we use a better partitioning scheme to perform efficient distributed join during contiging phase?
 - Can we rewrite the whole pipeline using batch parallel model?

Hash table insertions

```
// kmer is array of uint64_t of length N
// all key slots initially set to KEY_EMPTY
int slot = hash(kmer) % capacity; // starting slot
for (i = 0; i < MAX_PROBE; i++) { // maximum number of probes</pre>
  do { // acquire "lock" on key
   key = atomicCAS(keys[slot][N - 1], KEY_EMPTY, KEY_TRANSITION)
 } while (key != KEY TRANSITION);
  if (key == KEY_EMPTY) { // lock was acquired on an empty slot, set key
    for (l = 0; l < N - 1; l++) \{ key = atomicExch(keys[slot][l], kmer[l]); \}
    found key = true; break;
 } else if (\text{key} == \text{kmer}[N - 1]) { // this could be the same key
    found key = kmers equal(kmer, keys[slot]); // check for same key
    if (found key) break;
  }
  slot = (start slot + (i + 1) * (i + 1)) % capacity; // quadratic probe
if (found slot)
  atomicADD(values[slot][...]...; // inc k-mer count and extension counts
```

k-mers from contigs

- k-mers from contigs inserted into a separate, smaller *ctg-kmers* GPU hash table
- processed similarly to k-mers from reads, except:
 - count for k-mer is min of all inserts (set with atomicMin, after atomicCAS to check for zero)
 - after all insertions are done, a new GPU kernel purges conflicts (2 or more valid extensions on a side) - marked KEY_EMPTY
- after *ctg-kmers* completed, insert each entry into *read-kmers* hash table:
 - if the k-mer already exists (from a read) but it doesn't have good extensions, add the extension counts from the contig k-mer
 - if the k-mer doesn't exist, add the ctg k-mer

Final stages

- purge invalid entries, one GPU thread per hash table slot:
 - purge (set to KEY_EMPTY) k-mers with a count of 1 or without valid extensions on both sides
- compute final extensions, one GPU thread per hash table slot:
 - \circ $\,$ for each side, convert array of counts for A,C,G,T into a single base extension or X or F $\,$
- each GPU thread inserts (kmer, count, left ext, right ext) into new compact GPU hash table
- copy compact GPU hash table into host memory
- iterate through array on CPU, inserting each (kmer, count, left ext, right ext) element into local CPU hash table for later use in deBruijn graph traverse