

Summary statistics and quality of the *Dictyostelium discoideum* reference genome assembly

Jason Moggridge

29/01/2021

Introduction

In this work, I examine the reference genome of the social amoeba *Dictyostelium discoideum* (Eichinger *et al.*, 2005). *D. discoideum* is a interesting model organism for many areas of research in molecular biology, particularly because of their ability to cooperatively form a multicellular super-organism for the purpose of sporulation under starvation conditions (Williams, 2010). The *de-novo* genome assembly was created by whole chromosome shotgun sequencing using the Sanger method. For this, chromosomes were separated using PFGE. Isolated chromosomes were fragmented and then cloned using plasmids and yeast artificial chromosomes (YACs). As this was among the first protozoan genome projects, there was no draft genome to guide contig assembly. As such, data from HAPPY mapping, previously mapped genes, and YACs were used in creating the assembly. HAPPY mapping is akin to linkage mapping but uses random DNA fragmentation and PCR to get distance information for markers (sequence-tagged sites in this case) instead of cloning (Dear and Cook, 1993). This information was used to assign reads from each library to chromosome-specific bins through BLAST or **Atlas Overlapper**. Binned reads were then joined into contigs by either GAP4 or PHRED/PHRAM/CONSED software. Then read-pair data and BLAST searches of all sequence data were used to extend contigs and create scaffolds. Gap were closed by various strategies, depending on whether sequencing was challenged by repetitive or A+T rich regions. The authors presented an analysis of the genome's nucleotide composition, repeats, transposons, tRNAs, telomeres, and centromere regions. In addition, gene prediction and translation was performed to create a phylogeny to place *D. discoideum* among other eukaryotic phyla based on the predicted proteome.

Links - ([Browse genome record on NCBI](#)) - ([Zipped FASTA file through ftp](#))

Analysis of the *D. discoideum* genome

Lines starting with '\$' are unix commands and '#' are comments, otherwise lines are R code

The *D. discoideum* assembly has a size of 34.2 Mbp.

```
# remove lines with headers (>), strip newline characters (\n), count characters
$ grep -v '>' dicty_genome.fna | tr -d '\n' | wc -c
```

There are 41 scaffolds in the assembly.

These include 6 genomic chromosomes, the mitochondrial chromosome and a plasmid chromosome, as well as 33 unplaced scaffolds.

```
# count scaffold headings
$ grep '>' dicty_genome.fna | wc -l
```

The mean scaffold length is 834.3 kbp.

Scaffold lengths are plotted in fig. 1.

```
# create a table with headings: Scaffold, Length
$ touch seqLens.tsv
$ echo -e 'Scaffold\tLength' > seqLens.tsv
# loop over headings, get length of following sequence for each
$ awk '/^>/ {if (seqLen){print seqLen}; printf $0"\t";seqLen=0;next; }
{seqLen += length($0)}END{print seqLen}' dicty_genome.fna >> seqLens.tsv
```

N50 & L50

I was able to recreate the N50 and L50 statistics reported on NCBI (5,450,249 bp and 3) as follows:

```
N50 <- function(scaffold_lengths) {
  lengths <- sort(scaffold_lengths, decreasing = T)
  total <- sum(lengths)
  cumsum <- 0
  for (i in seq_along(lengths)) {
    cumsum <- cumsum + lengths[i]
    if (cumsum >= total/2)
      return(lengths[i])
  }
}
seqLen <- read_delim("data/seqLens.tsv", "\t")
dicty.N50 <- N50(seqLen$Length)
dicty.L50 <- which(sort(seqLen$Length, decreasing = T) == dicty.N50)
```

Ns are rare in the assembly (0.07 %).

```
# N-content; total bp
$ grep -v '^>' dicty_genome.fna | tr -cd 'Nn' | wc -c
$ grep -v '^>' dicty_genome.fna | tr -d '\n' | wc -c
```

```
N_percent <- 23142 * 100/34204973
```

The ratio of repetitive to unique sequence is 1.12.

```
# repetitive; unique
$ grep -v '^>' dicty_genome.fna | tr -d '\n' | tr -cd acgt | wc -c
$ grep -v '^>' dicty_genome.fna | tr -d '\n' | tr -cd ACGT | wc -c
# 18093064, 16088767
```

Table 1: Less than 1% of the assembly is in unplaced scaffolds

```
seqLen <- seqLen %>%
  mutate(Type = case_when(
    str_detect(Scaffold, 'chrUn_') ~ 'Unplaced',
    str_detect(Scaffold, 'plasmid') ~ 'Plasmid',
```

```

str_detect(Scaffold, 'mitochond') ~ 'Mitochondrial',
TRUE ~ 'Nuclear chr.'))
seqlen %>%
  group_by(Type) %>%
  summarize('Percent of Genome' = round(sum.Length)/sum(seqlen$.Length)*100, 2))

```

| Type | Percent of Genome |
|---------------|-------------------|
| Mitochondrial | 0.16 |
| Nuclear chr. | 99.23 |
| Plasmid | 0.04 |
| Unplaced | 0.56 |

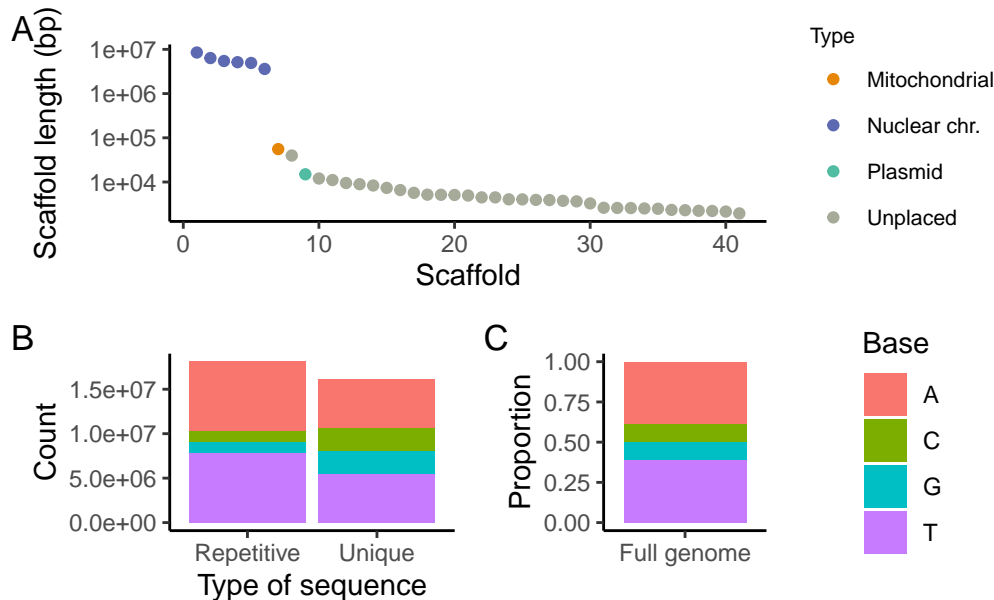
The *D. discoideum* genome is very A+T-rich. All bases were counted and these are presented in fig. 1.

```

$ touch base_count.csv; echo "Base, Count" > base_count.csv
$ for base in A C G T a c g t
$ do
$   count=$(cat dicty_genome.fna | grep -v '>' | tr -cd $base | wc -c)
$   echo "$base, $count" >> base_count.csv
$ done

```

Figure 1: *D. discoideum* scaffold lengths and nucleotide composition



*Code omitted for brevity

Discussion

The *D. discoideum* reference genome is relatively small (34.2 Mbp) and the assembly is of remarkably high quality and completeness. The assembly has nearly complete scaffolds for the six nuclear chromosomes, with

very few unplaced scaffolds (33; fig. 1). The percentage of the assembly in these unplaced scaffolds is only 0.56 %. Correspondingly, the assembly has a very large N50 (5.45 Mbp) and a very small L50 (3 scaffolds). The sequence data is of high quality, with only 0.07 % ambiguous bases (N). The genome contains a relatively even split between repetitive sequence and unique sequence (ratio = 1.12).

The quality of the genome is surprisingly good, given that the effort was reliant on cloning and Sanger sequencing technology of the time, and dealt with complex repetitive regions that are difficult to assemble, as well as A+T-rich tracts (fig. 2) that can be difficult to clone and sequence (Eichinger *et al.*, 2011). To create this near-perfect *de-novo* genome assembly, the authors applied data from contiguity methods that were commonly used at the time: genetic mapping (HAPPY maps), physical mapping (with YACs, *etc.*), and paired reads. Overall, this reference genome should prove to be an invaluable resource for studying this intriguing model amoeba.

References

- Dear PH, and PR Cook. 1993. Happy mapping: linkage mapping using a physical analogue of meiosis. *Nucleic Acids Res.* 11; 21(1): 13–20.
- Eichinger L, *et al.*. 2005. The genome of the social amoeba *Dictyostelium discoideum*. *Nature*. May 5; 435(7038): 43–57.
- Williams JG. 2010. *Dictyostelium* finds new roles to model. *Genetics*, 185(3):717–726.