



**Dovetail**Genomics

Canis Lupus Lupus HiC HiRise Genome  
Assembly Report

Phil Ewels

Swedish Museum Of Natural History

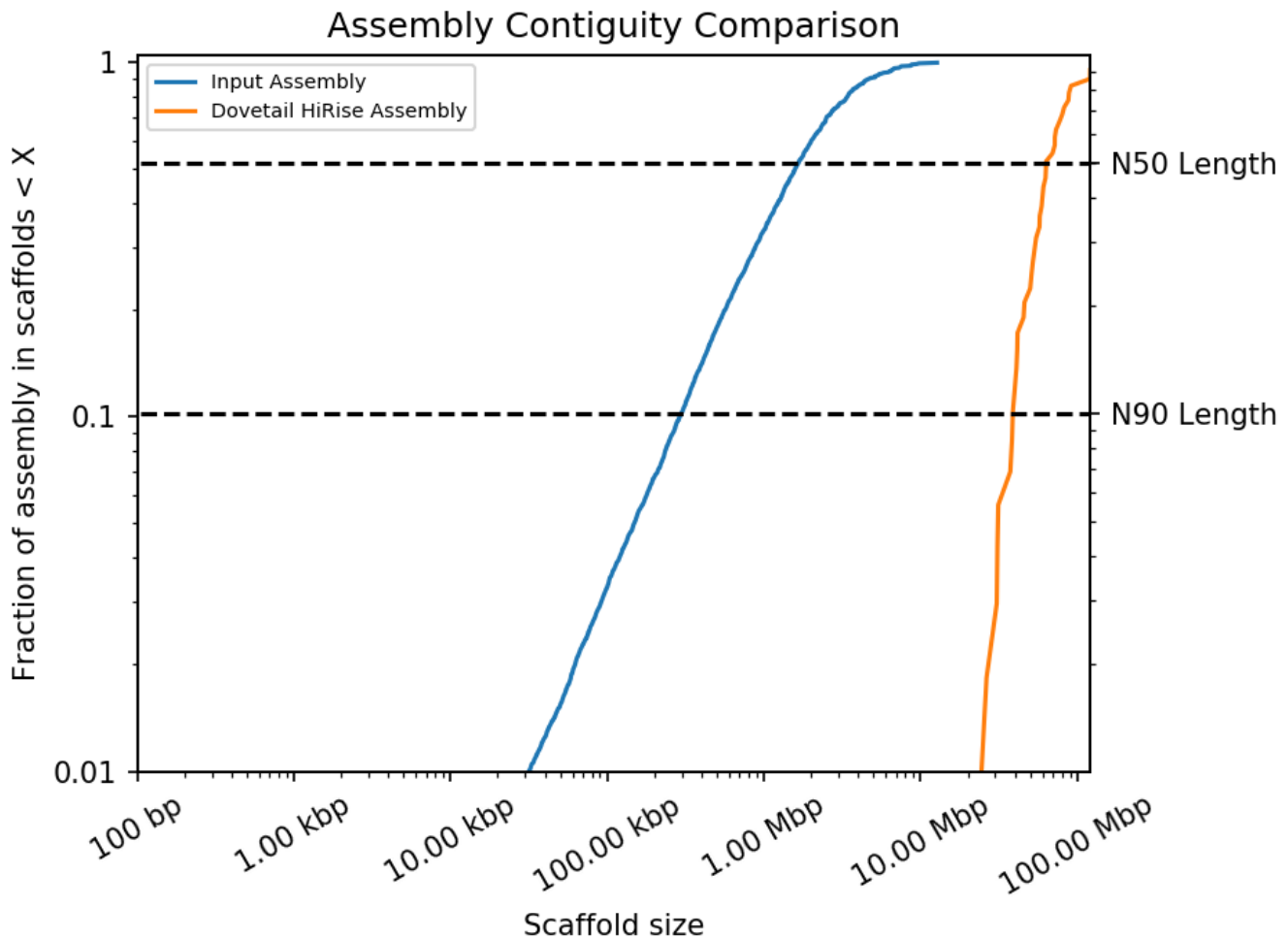
May 3, 2019

# Canis Lupus Lupus

## HiC HiRise assembly

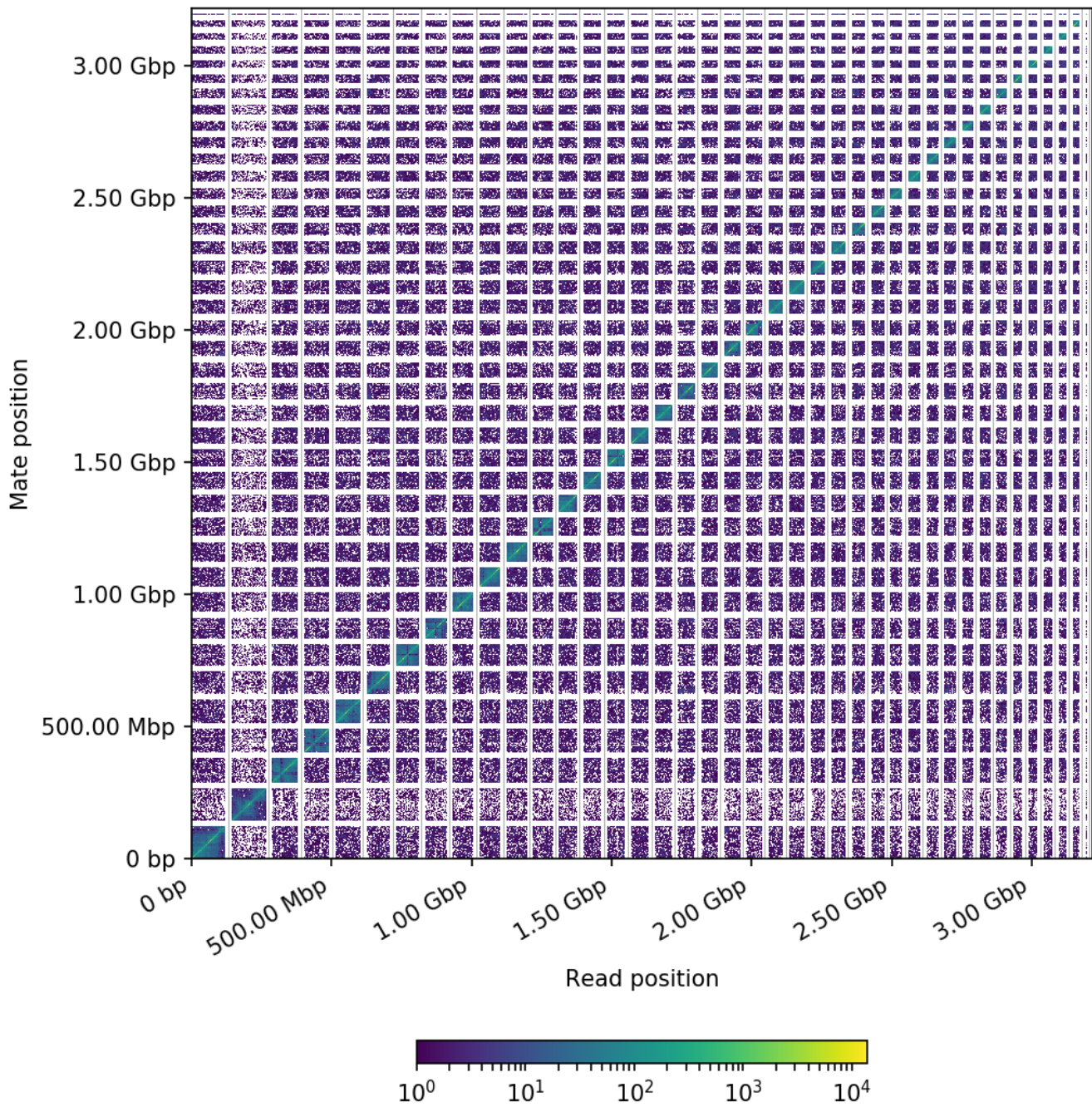
Estimated physical coverage (10-10,000 kb pairs): 4,321.70X

	Input Assembly	Dovetail HiRise Assembly
Total Length	2,312.94 Mb	2,313.38 Mb
L50/N50	404 scaffolds; 1.584 Mb	14 scaffolds; 63.093 Mb
L90/N90	1,682 scaffolds; 0.290 Mb	33 scaffolds; 38.334 Mb



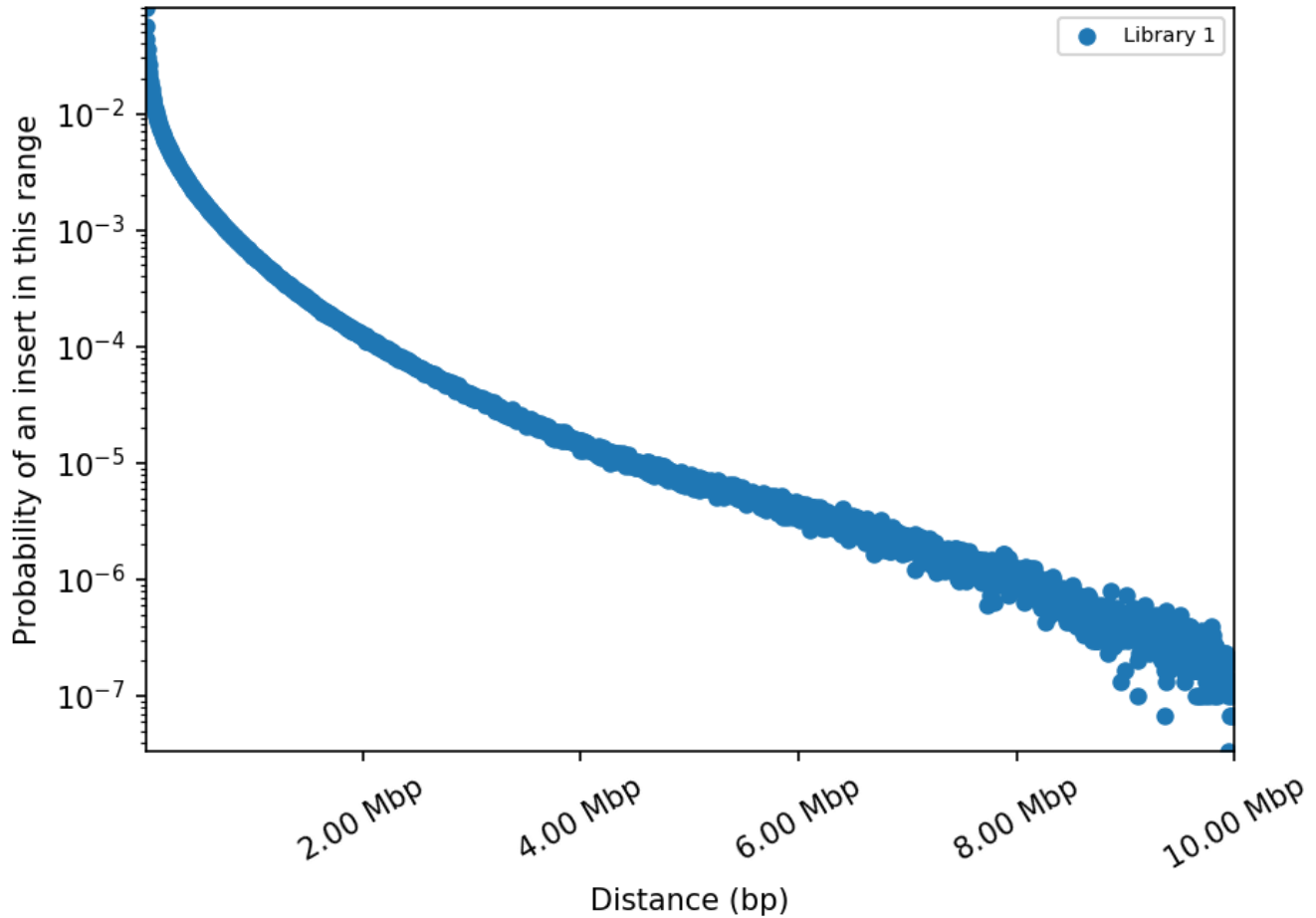
A comparison of the contiguity of the input assembly and the final HiRise scaffolds. Each curve shows the fraction of the total length of the assembly present in scaffolds of a given length or smaller. The fraction of the assembly is indicated on the Y-axis and the scaffold length in basepairs is given on the X-axis. The two dashed lines mark the N50 and N90 lengths of each assembly. Scaffolds less than 1 kb are excluded.

Link density histogram



In this figure, the x and y axes give the mapping positions of the first and second read in the read pair respectively, grouped into bins. The color of each square gives the number of read pairs within that bin. White vertical and black horizontal lines have been added to show the borders between scaffolds. Scaffolds less than 1 Mb are excluded.

## Library insert size distribution



This figure shows the distribution of insert sizes in the Dovetail library. The distance between the forward and reverse reads is given on the X-axis in basepairs, and the probability of observing a read pair with a given insert size is shown on the Y-axis.

<b>Comparative Assembly Statistics</b>		
	<b>Input Assembly</b>	<b>Dovetail HiRise Assembly</b>
Longest Scaffold	12,722,267 bp	121,009,298 bp
Number of scaffolds	8,747	4,416
Number of scaffolds > 1kb	8,643	4,312
Contig N50	83.96 kb	83.96 kb
Number of gaps	54,447	58,783
Percent of genome in gaps	1.88%	1.90%

\* Note: Every join made by HiRise creates a gap.

<b>Other Statistics</b>	
Number of breaks made to input assembly by HiRise	2
Number of joins made by HiRise	4,334
Library 1 stats	452M read pairs; 2x151 bp

BUSCO Stats					
	Single copy	Duplicated	Fragmented	Missing	Total
Input Assembly	251	13	11	28	303
Dovetail HiRise Assembly	250	13	15	25	303

Number of BUSCO (Benchmarking Universal Single-Copy Ortholog) genes found in the assembly before and after HiRise using the eukaryota odb9 dataset. Genes are split into four categories: complete and single-copy, complete and duplicated, fragmented, and missing.

## Glossary

**Sequence Coverage** - For a given position in the genome, the sequence coverage is the number of times this basepair is directly observed in the sequencing data. Typically given as an average over the whole genome, or estimated by the total length of reads divided by the genome size.

**Physical Coverage** - For a given position in the genome, the physical coverage is the number of read pairs that span this position. Typically given as an average over the whole genome, or estimated by the area under the insert distribution divided by the genome size.

**Contig** - A contiguous genomic sequence without any gaps in an assembly.

**Scaffold** - A genomic sequence consisting of contigs that have been ordered and oriented relative to each other. Contigs within scaffolds are separated by gaps (indicated by stretches of Ns).

**N50** - The scaffold length such that the sum of the lengths of all scaffolds of this size or larger is equal to 50% of the total assembly length.

**N90** - The scaffold length such that the sum of the lengths of all scaffolds of this size or larger is equal to 90% of the total assembly length.

**L50** - The smallest number of scaffolds that make up 50% of the total assembly length.

**L90** - The smallest number of scaffolds that make up 90% of the total assembly length.