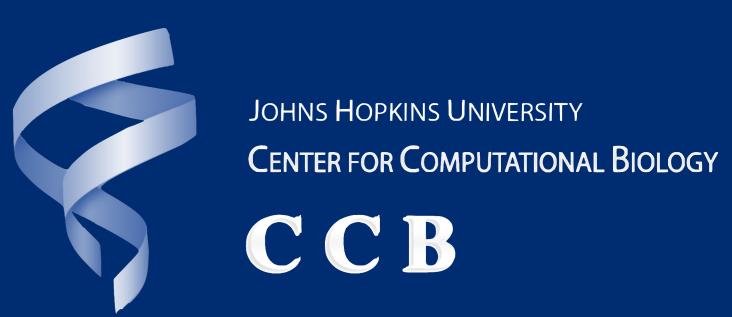


# Combining DNA and protein alignments to improve genome annotation with LiftOn



Kuan-Hao Chao<sup>1,2,\*</sup>, Jakob M. Heinz<sup>5</sup>, Celine Hoh<sup>1,2</sup>, Alan Mao<sup>1,2,3</sup>, Alaina Shumate<sup>2,3</sup>, Steven Salzberg<sup>1,2,3,4,\*</sup>, Mihaela Pertea<sup>1,2,3,\*</sup>

<sup>1</sup>Department of Computer Science,

<sup>2</sup>Center for Computational Biology,

<sup>3</sup>Department of Biomedical Engineering,

<sup>4</sup>Department of Biostatistics, Johns Hopkins University,

<sup>4</sup>Department of Biomedical Informatics, Harvard University

<u>k</u>

kh.chao@cs.jhu.edu



https://ccb.jhu.edu/lifton/

Protein identity



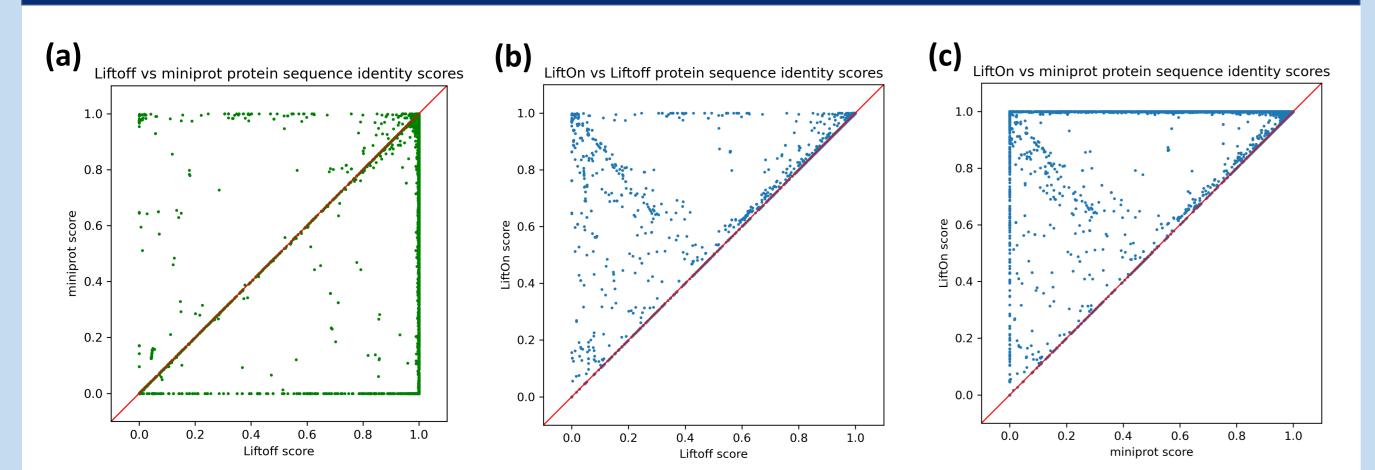
https://github.com/Kuanhao-Chao/LiftOn



### Introduction

- LiftOn uses both DNA-DNA alignments (from Liftoff) & protein-DNA alignments (from miniprot) to map annotations between genome assemblies of the same or different species.
- LiftOn's protein-maximization algorithm improves the annotation of protein-coding genes in the T2T- CHM13 genome.
- LiftOn can map annotation between relatively distant species species, at least as divergent as mouse and rat.

## LifOn vs Liftoff vs miniprot

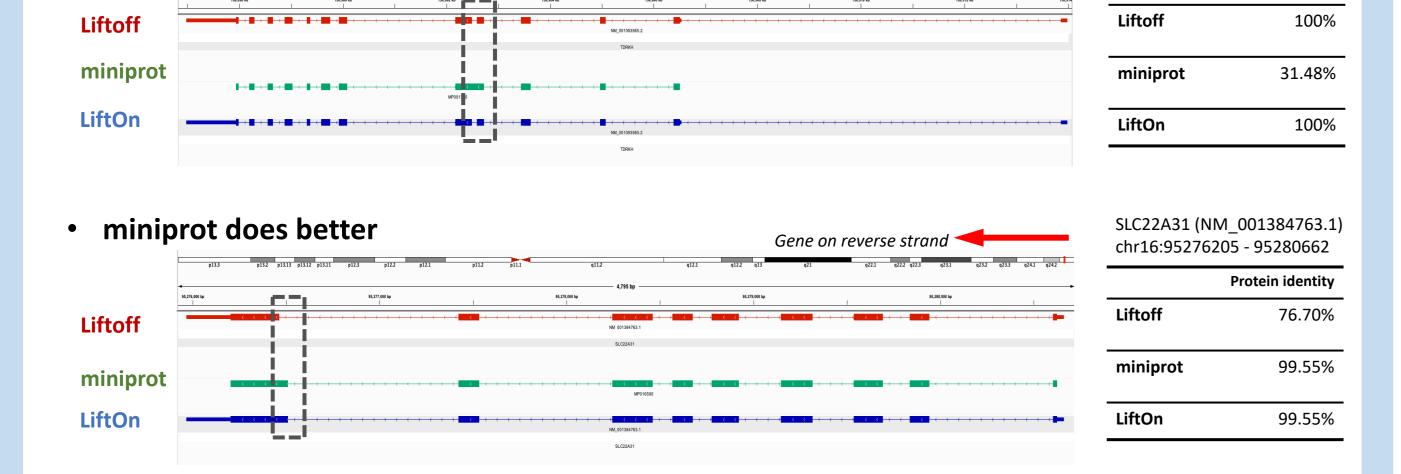


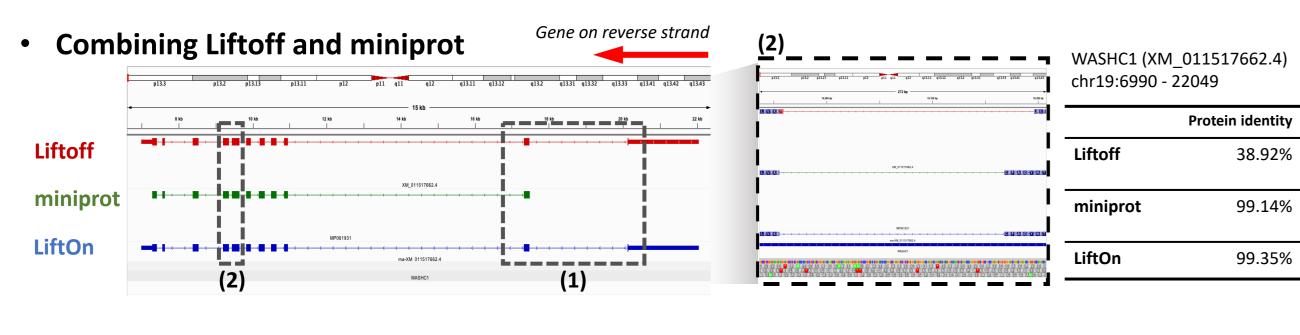
The scatter plot of protein sequence identity comparing between (a) miniprot vs Liftoff,

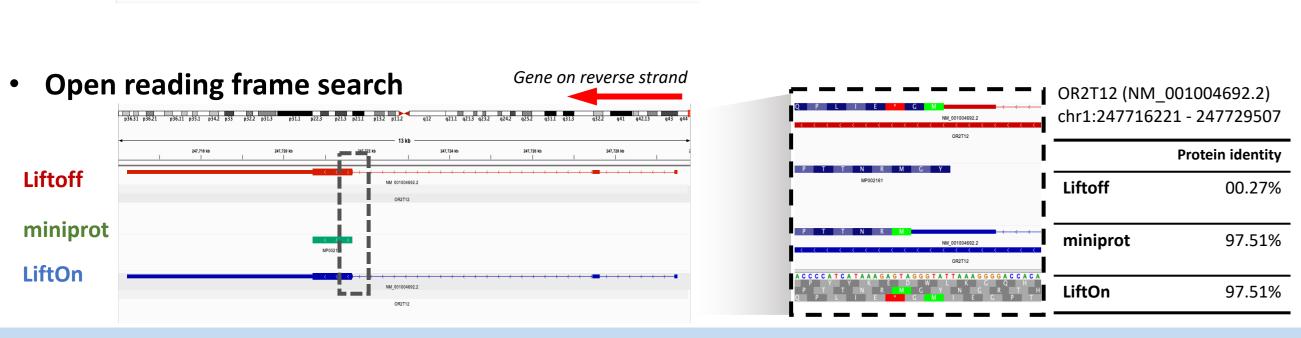
(b) LiftOn vs Liftoff, and (c) LiftOn vs miniprot. Each dot represent a protein-coding transcript.

#### IGV screenshots of gene loci examples

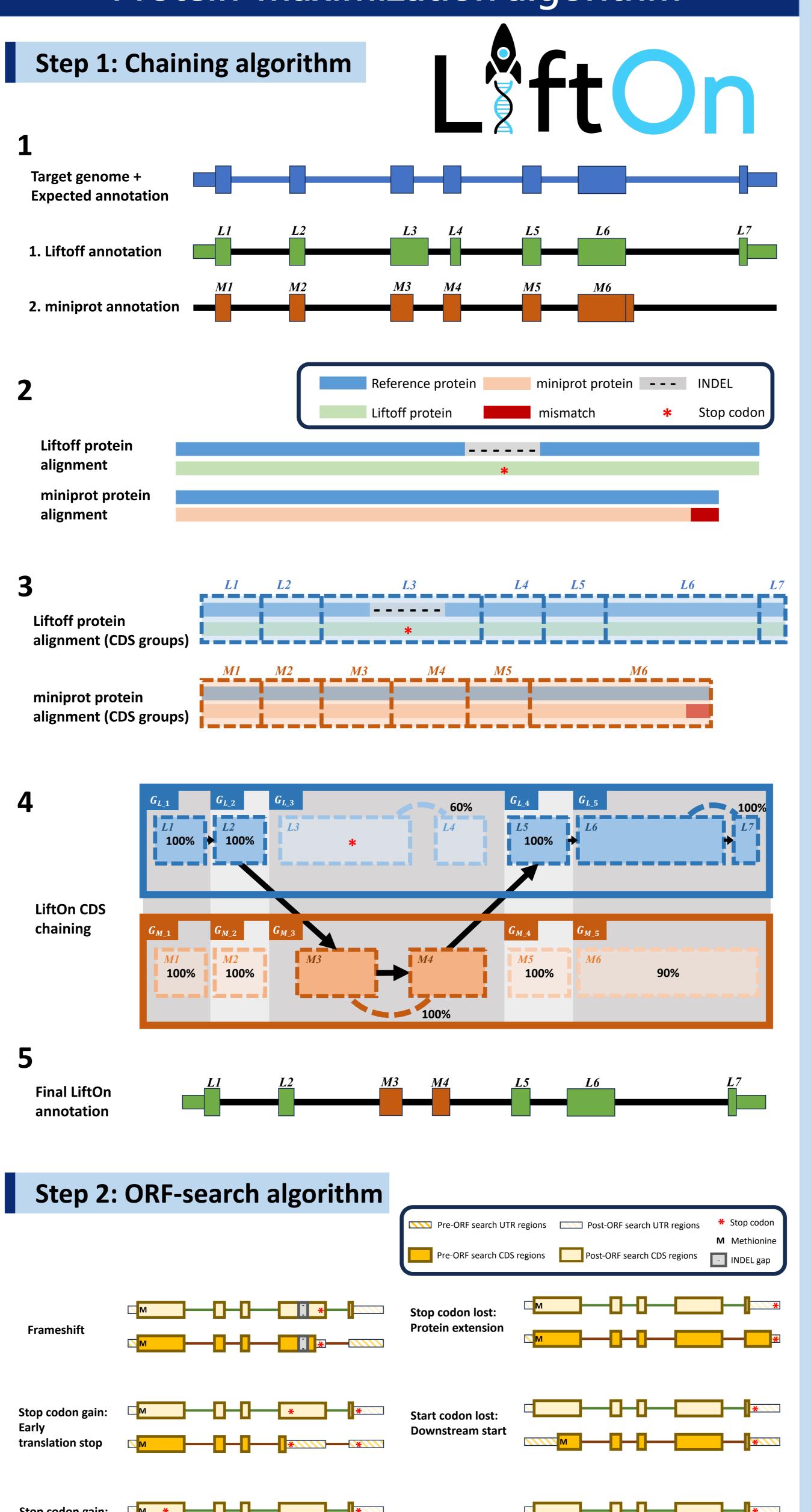
Liftoff does better





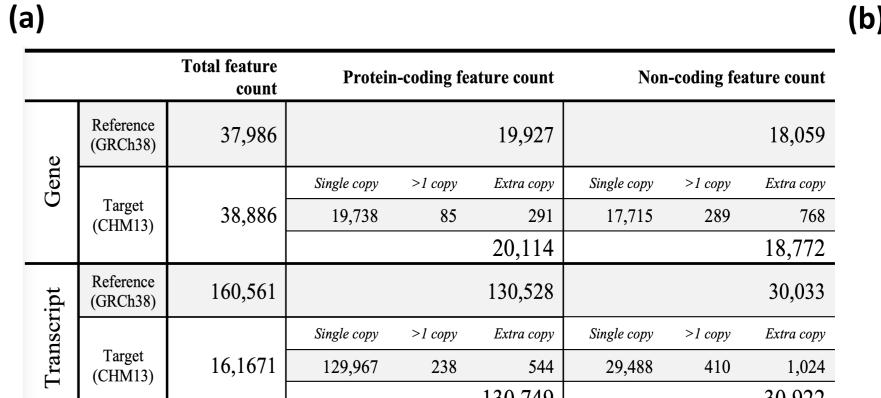


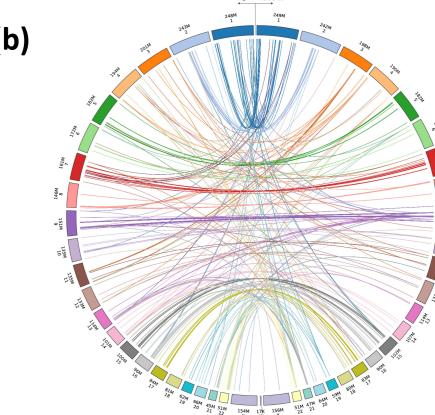
# Protein-maximization algorithm

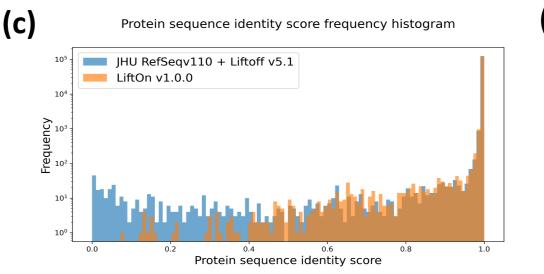


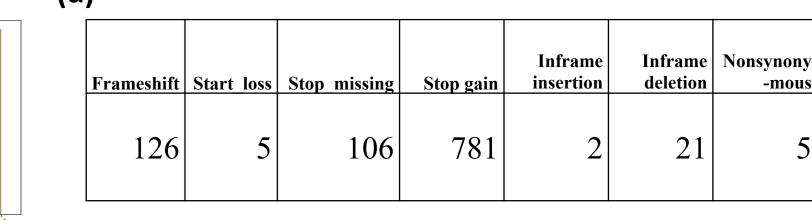
#### Results

#### Mapping RefSeq v220 from GRCh38 to CHM13



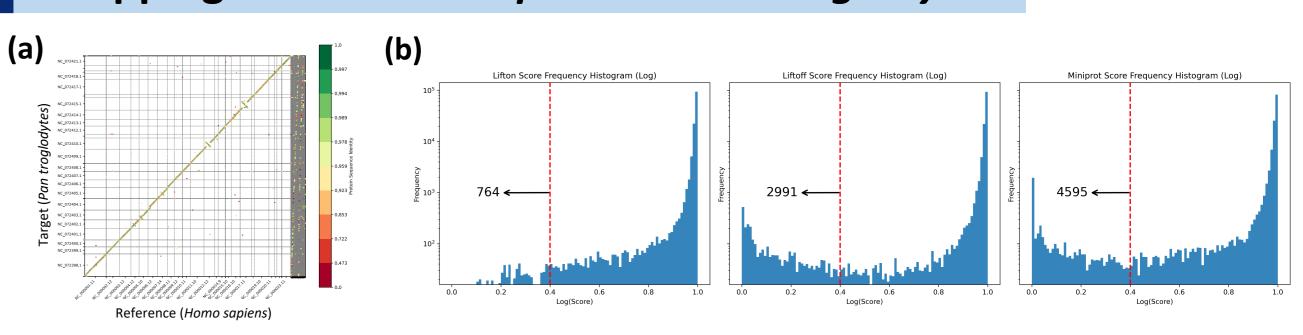




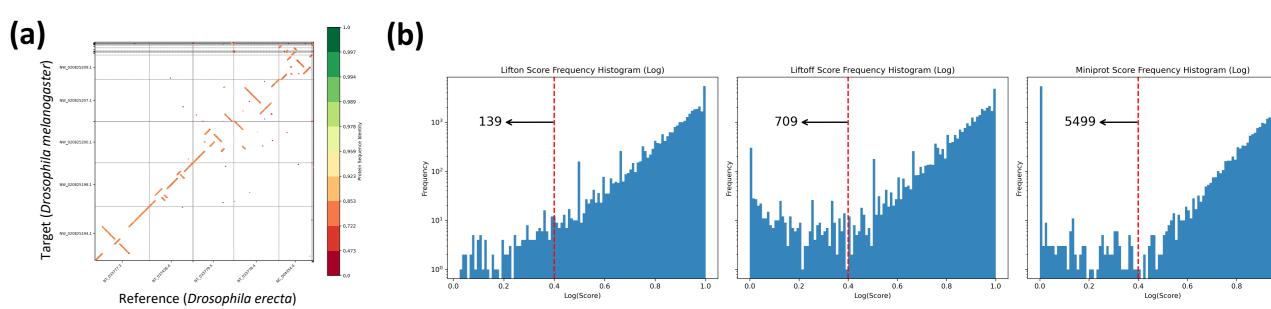


- (a) LiftOn GRCh38-to-CHM13 mapping summary of protein-coding and non-coding genes and transcripts.
- (b) Circos plot shows relative positions of extra gene copies between target (left) and reference (right)
- (c) Protein-coding sequence identity frequency plots comparing LiftOn and JHU RefSeqv110+Liftoff v5.1
- (d) Mutation reports of LiftOn-generated human annotation.

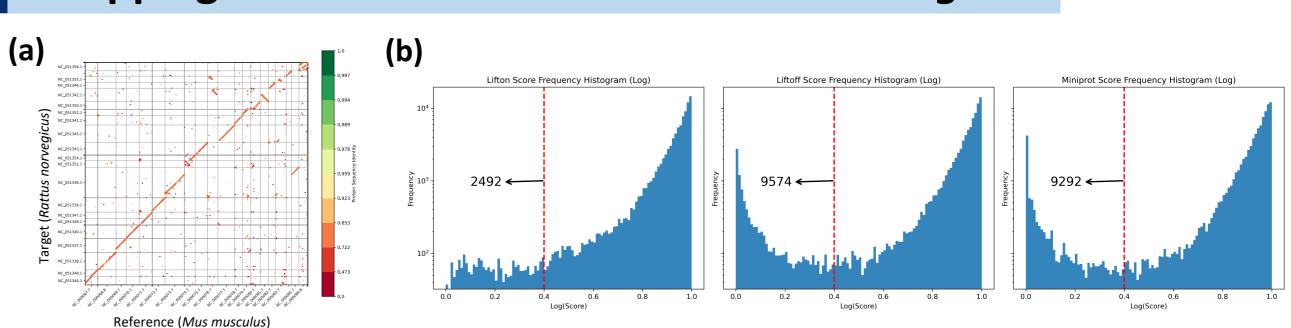
#### Mapping from Homo sapiens to Pan troglodytes



## Mapping from *Drosophila melanogaster* to *Drosophila erecta*



## Mapping from Mus musculus to Rattus norvegicus



- (a)Protein-coding gene order plot with x-axis for reference genome and y-axis for target genome, using a logarithmic color scale from green (identical) to red for protein sequence identities.
- (b) Logarithmic frequency plots of protein sequence identity for LiftOn, Liftoff, and miniprot.

Acknowledgements: supported by NIH under grants R01-HG006677, R35-GM130151, and T32HG002295 and by NSF under grants DBI-1759518.