

Brief Communication

**nf-LO: A scalable, containerised workflow for genome-to-genome
lift over**

Andrea Talenti¹ and James Prendergast¹

¹ The Roslin Institute, University of Edinburgh, Easter Bush, Midlothian, EH25 9RG, UK

Correspondence to: Andrea Talenti – andrea.talenti@ed.ac.uk

Keywords: liftover, assembly, Nextflow, workflow

13 **Abstract**

14 The increasing availability of new genome assemblies often comes with an impaired
15 amount of associated genomic annotations, limiting the range of studies that can be
16 performed. A common workaround is to lift over annotations from better annotated
17 genomes. However, generating the files required to perform a liftover is
18 computationally and labour intensive and only a limited number are currently publicly
19 available.

20 Here we present nf-LO (nextflow-LiftOver), a containerised and scalable Nextflow
21 pipeline that enables liftovers within and between any species for which assemblies
22 are available. nf-LO will consequently facilitates data interpretation across a broad
23 range of genomic studies.

24

25 **Main body**

26 The advent of third generation sequencing and ultra-fast assemblers (Joseph et al.
27 2018; Ruan and Li 2020) allows for the generation of high quality *de novo* assemblies
28 in a fraction of the previous time. As a result increasingly large numbers of new
29 genomes for several species are being generated (Zoonomia consortium 2020).

30 Despite this increased availability, novel assemblies most often lack the extensive
31 annotation data required to perform downstream analyses. Not only simple
32 annotations such as gene models, but also supplementary resources for researcher
33 to understand the biological significance of their studies. Unfortunately, such
34 resources are generally only available for a small number of model organisms (OMIA;
35 Amberger et al. 2015; Carithers and Moore 2015; Hu et al. 2019).

36 A solution to the problem is to liftover positions and annotation (i.e. cross-mapping of
37 the loci) to the new genome from well-annotated assemblies, using tools such as
38 LiftOver (Navarro Gonzalez et al. 2021) and NCBI Remap (Luu et al. 2020). However,
39 the alignment files required to perform these analyses are only currently publicly
40 available for a small number of pairs of genomes. For all other pairs of genomes
41 researchers have to generate their own liftover files. Only a few algorithms address
42 the problem in an easy to implement and distributable way, e.g. flo for same species
43 liftovers (Pracana et al. 2017) and LiftOff for ultra-fast liftovers (Shumate and Salzberg
44 2020)). In this study we present nf-LO, a scalable workflow to generate liftover files for
45 any pair of genomes based on the UCSC liftover pipeline. Nf-LO can directly pull
46 genomes from public repositories, supports parallelised alignment using a range of

47 alignment tools and can be finely tuned to achieve the desired sensitivity, speed of
48 process and repeatability of analyses.

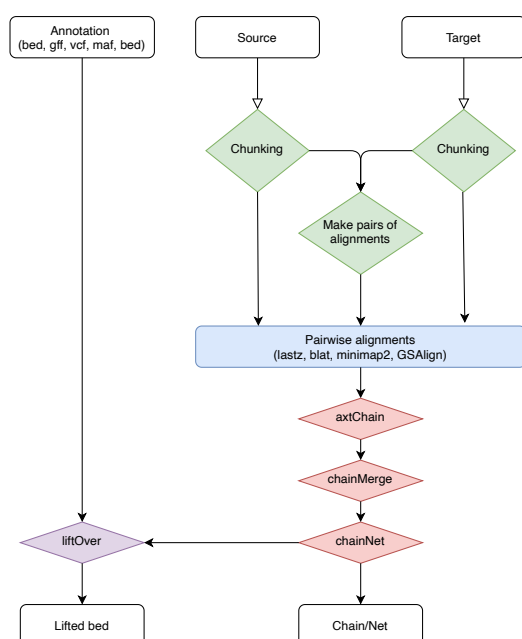
49 nf-LO is a workflow to facilitate the generation of genome alignment chain files
50 compatible with the LiftOver utility. It is written in Nextflow, a domain specific language
51 v2 (DSL) and workflow manager, that allows easy implementation, redistribution and
52 scalability of complex workflows across every Unix-based operating system; ranging
53 from a desktop machine to cloud computing and HPC clusters. The dependencies are
54 shipped alongside the workflow as docker containers or as an anaconda environment,
55 facilitating the diffusion and adoption of the workflow across different systems.

56 The software accepts any two input genomes in fasta format, or alternatively can
57 download a resource by providing a web address, an iGenome identifier or an NCBI
58 GenBank or RefSeq accession. The workflow is shown in Figure 1, and in brief
59 consists of three core steps, and one optional one: 1) chunking the two genomes, 2)
60 pairwise alignment of the blocks, 3) generating the chain-net file that can be used to
61 perform the liftover and, if a bed/gff/gtf/vcf/bam/maf file is provided, 4) performing the
62 liftover from source to target. The chunking approach dramatically reduces the runtime
63 of the analysis by parallelizing the alignments.

64

65

66 Figure 1



67

68 The alignment phase can be performed in different ways, depending on the type and
69 sensitivity required by the user. For same-species alignments, we provide native
70 support for both blat (Kent 2002), the aligner of choice for same species liftover files
71 from the UCSC genome browser, and GSAIgn (Lin and Hsu 2020), a new, high speed
72 same-species alignment software. For performing different-species liftovers, nf-LO
73 also incorporates lastz (Harris 2007), used by the UCSC genome browser to generate
74 between species liftover files, and minimap2 (Li 2018), one of the fastest genome-to-
75 genome aligners. All these aligners are integrated within the workflow, keeping
76 unchanged the UCSC backbone for downstream stages (UCSC 2018). We provide
77 canned configurations for each aligner based on how distant the two genomes are
78 (e.g. near or far), with the possibility to provide sets of custom parameters to achieve
79 the desired balance between speed and sensitivity (Supplementary table 1). nf-LO
80 achieves similar liftover coverage as liftover files from UCSC with appropriate tuning
81 of the parameters (Supplementary table 2).

82 The third stage processes the alignments analogously to the UCSC processing
83 pipeline, obtaining the chain-net files to perform the actual liftover. Finally, the fourth
84 step supports both the standard bed format with the LiftOver software, or several
85 additional formats using CrossMap (Zhao et al. 2014), including popular formats such
86 as VCF, BAM and GFF.

87 In conclusion, we provide a transposition of the UCSC liftover pipeline within the
88 Nextflow language, together with the necessary containers to run the analyses,
89 allowing an easy, streamlined implementation in any Unix-based system. We believe
90 that this workflow will be of use across genomics studies, facilitating research work
91 and enabling data interpretation.

92

93 **Code availability**

94 The code described in the paper is publicly available on GitHub at the repository
95 <https://github.com/evotools/nf-LO>. The documentation for the software can be
96 accessed in the wiki page of the website (<https://github.com/evotools/nf-LO/wiki>).

97

98 **Authors' contributions**

99 AT and JP conceived the study. AT developed the software. AT and JP tested the
100 code. AT and JP contributed to data interpretation and drafted the manuscript. All
101 authors reviewed and approved the final manuscript.

102

103 **Acknowledgements**

104 This work was supported by BBSRC grants BB/T019468/1 and BBS/E/D/10002070.

105

106 **Captions**

107 Figure 1 - Scheme of the workflow of nf-LO with the chunking (step 1, in green),
108 alignment (step 2, in blue), generation of the liftover files (step 3, in red) and optionally
109 lifting of the variants to the target genome (step 4, in purple).

110 Supplementary Table 1 – Comparison of the run times of different aligners and
111 configurations using the human genome GRCh38 as the source and four other large
112 genomes (>1Gbp) as targets on a Scientific Linux 6.9 system with AMD Opteron 6376
113 2.3GHz 64-cores and 500 GB of RAM.

114 Supplementary Table 2 – Coverage for the liftover chain files both generated by us
115 and those available from the UCSC genome database, calculated by converting the
116 chain files to maf (chainToAxt > axtToMaf) and then using mafCoverage (Earl et al.
117 2014).

118

119

120 References

- 121 Amberger JS, Bocchini CA, Schiettecatte F, Scott AF, Hamosh A. 2015. OMIM.org:
122 Online Mendelian Inheritance in Man (OMIM®), an Online catalog of human
123 genes and genetic disorders. *Nucleic Acids Res.* 43:D789–D798.
- 124 Carithers LJ, Moore HM. 2015. The Genotype-Tissue Expression (GTEx) Project.
125 *Biopreserv. Biobank.*
- 126 Earl D, Nguyen N, Hickey G, Harris RS, Fitzgerald S, Beal K, Seledtsov I, Molodtsov
127 V, Raney BJ, Clawson H, et al. 2014. Alignathon: A competitive assessment of
128 whole-genome alignment methods. *Genome Res.* 24:2077–2089.
- 129 Harris RS. 2007. Improved pairwise alignment of genomic DNA. Available from:
130 http://www.bx.psu.edu/~rsharris/rsharris_phd_thesis_2007.pdf
- 131 Hu ZL, Park CA, Reecy JM. 2019. Building a livestock genetic and genomic
132 information knowledgebase through integrative developments of Animal QTLdb
133 and CorrDB. *Nucleic Acids Res.* 47:D701–D710.
- 134 Joseph S, O'Connor RE, Al Mutery AF, Watson M, Larkin DM, Griffin DK. 2018.
135 Chromosome level genome assembly and comparative genomics between three
136 falcon species reveals an unusual pattern of genome organisation. *Diversity* 10.
- 137 Kent WJ. 2002. BLAT---The BLAST-Like Alignment Tool. *Genome Res.* 12:656–664.
- 138 Li H. 2018. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics*
139 34:3094–3100.
- 140 Lin HN, Hsu WL. 2020. GSAIalign: An efficient sequence alignment tool for intra-
141 species genomes. *BMC Genomics* 21.
- 142 Luu P-L, Ong P-T, Dinh T-P, Clark SJ. 2020. Benchmark study comparing liftover
143 tools for genome conversion of epigenome sequencing data. *NAR Genomics*
144 *Bioinforma.* 2.
- 145 Navarro Gonzalez J, Zweig AS, Speir ML, Schmelter D, Rosenbloom KR, Raney BJ,
146 Powell CC, Nassar LR, Maulding ND, Lee CM, et al. 2021. The UCSC genome
147 browser database: 2021 update. *Nucleic Acids Res.* 49.
- 148 OMIA. Online Mendelian Inheritance in Animals. *Sydney Sch. Vet. Sci.* [Internet].
149 Available from: <https://omia.org/>
- 150 Pracana R, Priyam A, Levantis I, Nichols RA, Wurm Y. 2017. The fire ant social
151 chromosome supergene variant Sb shows low diversity but high divergence
152 from SB. *Mol. Ecol.* 26:2864–2879.
- 153 Ruan J, Li H. 2020. Fast and accurate long-read assembly with wtdbg2. *Nat.*

- 154 *Methods* [Internet] 17:155–158. Available from:
155 <http://dx.doi.org/10.1038/s41592-019-0669-3>
- 156 Shumate A, Salzberg SL. 2020. Liftoff: accurate mapping of gene annotations.
157 *Bioinformatics*.
- 158 UCSC. 2018. Minimal steps for liftover. Available from:
159 http://genomewiki.ucsc.edu/index.php/Minimal_Steps_For_LiftOver
- 160 Zhao H, Sun Z, Wang J, Huang H, Kocher JP, Wang L. 2014. CrossMap: A versatile
161 tool for coordinate conversion between genome assemblies. *Bioinformatics* 30.
- 162 Zoonomia consortium. 2020. A comparative genomics multitool for scientific
163 discovery and conservation. *Nature* [Internet] 587:240–245. Available from:
164 <http://www.nature.com/articles/s41586-020-2876-6>
165

