





Computational methods to improve genome annotation, splice site prediction, and gene expression prediction

Kuan-Hao Chao

2024.08.06







Part I

Genome Annotation

- LiftOn: genome annotation lift-over
- Applications: mapping genes between two genomes



Genome annotation



Genome (FASTA)

chr1 BestRefSeq exon 454672 459678 ID=exon-NM_001005221.2-3;Parent=rna-NM_001005221.2; chr1 BestRefSeq CDS 450740 451678 0 ID=cds-NP_001005221.2-1;Parent=rna-NM_001005221.2; chr1 BestRefSeq CDS 452658 453675 0 ID=cds-NP_001005221.2-2;Parent=rna-NM_001005221.2;	Annotation (GFF / GTF)	chr1 BestRefSeq chr1 BestRefSeq chr1 BestRefSeq chr1 BestRefSeq chr1 BestRefSeq chr1 BestRefSeq chr1 BestRefSeq chr1 BestRefSeq	gene 450740 451678 . mRNA 450740 451678 exon 450740 451678 . exon 452658 453675 . exon 454672 459678 . CDS 450740 451678 . CDS 452658 453675 .	 	 ID=gene-OR4F29; ID=rna-NM_001005221.2;Parent=gene-OR4F29; ID=exon-NM_001005221.2-1;Parent=rna-NM_001005221.2; ID=exon-NM_001005221.2-2;Parent=rna-NM_001005221.2; ID=cds-NP_001005221.2-1;Parent=rna-NM_001005221.2; ID=cds-NP_001005221.2-2;Parent=rna-NM_001005221.2;
--	---------------------------	--	--	--------------------------	--

Lift-over Problem Definition:





Lift-over Problem Definition:

Reference genome

R



https://www.sanger.ac.uk/data/genome-reference-consortium/

Application: GRCh38 to T2T-CHM13 lift-over

"Finished" the human genome project



nature

Explore content ~ Publish with us About the iournal 🗸

Original Article | Published: 01 February 2001

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium

Science

Current Issue First release papers

SCIENCE > VOL. 291, NO. 5507 > THE SEQUENCE OF THE HUMAN GET

局 SPECIAL REVIEWS

The Sequence of the Human Genome

J. CRAIG VENTER, MARK D. ADAMS, EUGENE W. MYERS, PETER W. LJ. [...]. AND XIAOHONG ZHU 🤇 +269 authors 🌖 Authors Info & Affiliation

Cost per Human Genome



https://youtu.be/MbYvTyIMc84?si=YjijRyLwq31MY5HM

Application: GRCh38 to T2T-CHM13 lift-over

- 238 Mbp added and corrected
- 180 Mbp of centromeric satellites
- 68 Mbp of segmental duplications

- 10 Mbp of rDNAs
- 182 Mbp of entirely novel sequence
- 1956 novel genes including 99 protein-coding



Introduction

Current Method Overview

If you were to use a CHM13 annotation ... Which lift-over tool to use?



Giulio Formenti 3:44 PM

if I was to use an annotation for CHM13, which would it be?

(gene annotation)



Arang Rhie 4:11 PM

https://s3-us-west-2.amazonaws.com/human-pangenomics/T2T/CHM13/assemblies/annotation.chm13v2.0_RefSeq_Liftoff_v5.1.gff3.gz or https://s3-us-west-2.amazonaws.com/human-pangenomics/T2T/CHM13/assemblies/annotation/chm13v2.0_RefSeq_Liftoff_v5.1.bb



Introduction



Telomere-to-Telomere (T2T) consortium slack channel

Lift-over problem, what methods are available?









- How can we do better?
 - Combining DNA- and protein-based alignments!



Results



Results

• outperforms state-of-the-art DNA- and protein-based liftover methods

• improves the annotation of protein-coding genes in T2T-CHM13 genome

• Improves annotation lift-over between distant species, such as mouse and rat

Introduction

Current Method Overview

Results

Evaluation Metrics

"Sequence pairwise alignment"





 $\frac{\#Matched_nucleotide}{\#alignment\ column} = \frac{17}{26} = 65.4\%$



"Do not penalize longer proteins"

Introduction

Current Method Overview

Results

Methods

10

e i

Map RefSeq v220 from GRCh38 -> CHM13V2.0

Compressed-gap protein sequence identity





Map RefSeq v220 from GRCh38 -> CHM13V2.0

Compressed-gap protein sequence identity







TDRKH (NM_001083965.2) chr1:150896981 - 150913985



Map RefSeq v220 from GRCh38 -> CHM13V2.0

Compressed-gap protein sequence identity







SLC22A31 (NM_001384763.1) chr16:95276205 - 95280662



e i

Map RefSeq v220 from GRCh38 -> CHM13V2.0

Compressed-gap protein sequence identity







WASHC1 (XM 011517662.4) chr19:6990 - 22049 Gene on reverse strand (2) **Protein identity** q13.11 q13.12 q13.2 q13.31 q13.32 q13.33 q13.41 q13.42 q13.43 a12 38.92% Liftoff Liftoff 99.14% miniprot XM 011517662.4 miniprot 99.35% LiftOn LiftOn ma-XM 011517662 (1) (2) 14 Introduction **Current Method Overview** Methods Results

Result 2: improve CHM13 protein annotations

Protein sequence identity score frequency histogram 10⁵ JHU RefSeqv110 + Liftoff v5.1 LiftOn v1.0.0 10^{4} Frequency 10³ 10⁵ 10¹ 10⁰ 0.2 0.0 1.0 0.4 0.6 0.8 Protein sequence identity score



Result 3: improve distant species lift-over

human to chimp



mouse to rat



Results

Drosophila m. to Drosophila e.



Methods

Introduction

Current Method Overview

ent of Computer Science, Johns Hopkins University, Baltimore, Maryland 21218-2683, USA

Daniel N. Baker and Ben Langmead

Result 3: improve distant species lift-over





LiftOn

NINIPOT

Liftoff

Liftoff

Liftoff

າ 8 1.0

18 1.0 0.7

0.0

17

Result 3: improve distant species lift-over









Thale cress (Arabidopsis thaliana)



Rice (Oryza sativa)

House mouse (*Mus musculus*)

Yeast (Saccharomyces cerevisiae)



(Saccharomyces cerevi

New Results

Follow this preprint

Combining DNA and protein alignments to improve genome annotation with LiftOn

Kuan-Hao Chao, Jakob M. Heinz, Celine Hoh,
Alan Mao, Alaina Shumate, Aliana Pertea,
Steven L Salzberg
https://doi.org/10.1101/2024.05.16.593026



Honey bee

(Apis mellifera)

fruit fly (Drosophila melanogaster)

Introduction

Results

Methods

18



Methods in Details

Protein-maximization algorithm

• Step 1: chaining CDSs

• Step 2: ORF search



Α

B Step 1: Align Liftoff & miniprot proteins to reference protein



C Step 2: Mapped CDS boundaries onto Liftoff & miniprot protein alignments



D Step 3: group CDSs by "accumulated AA in the reference protein"



LiftOn CDS chaining



D Step 3: group CDSs by "accumulated AA in the reference protein"



D Step 3: group CDSs by "accumulated AA in the reference protein"



Introduction

Results

22

D Step 3: group CDSs by "accumulated AA in the reference protein"



Introduction

Results

22



Results

Current Method Overview

Introduction



Google search: **"LiftOn genome"**

LiftOn: Accurate annotation mapping for GFF/GTF across assemblies

♂ ccb.jhu.edu/lifton

- GPL-3.0 license
- ☆ 52 stars 😵 2 forks ⊙ 1 watching
- 🗜 1 Branch 🛭 🕤 Tags 🔸 Activity

- LiftOn is a promising new tool to study comparative genomics
- 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, at least as divergent as mouse and rat.

Chao, K. H., Heinz, J. M., Hoh, C., Mao, A., Shumate, A., Pertea, M., & Salzberg, S. L. (2024). Combining DNA and protein alignments to improve genome annotation with LiftOn. **bioRxiv**.

ccb.jhu.edu/lifton

github.com/Kuanhao-Chao/LiftOn ²⁴

Part II

Deep learning splice site prediction

- OpenSpliceAI: Pytorch reimplementation of SpliceAI
- Splam: splice junction recognizer



Introduction



Sequence models map a sequence to a sequence

				2
SpliceAl-t	oolkit	Sp	olam	Future work

Reference: https://www.youtube.com/watch?v=luCBXCErkCs&t=197s



Neural ODEs



CNNs





Future work



SpliceAl-toolkit

Splam

Future work
Input Prompt: Recite the first law of robotics

Output:



https://jalammar.github.io/how-gpt3-worksvisualizations-animations/

Hidden Layer	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Hidden Layer	\bigcirc	\bigcirc	\bigcirc	\bigcirc	\bigcirc	0	\bigcirc	\bigcirc	\bigcirc	\bigcirc	0	\bigcirc	\bigcirc	\bigcirc	\bigcirc	
Hidden Layer	\bigcirc															

Output 😑 😑

https://deepmind.google/discover/blog/wavenet-agenerative-model-for-raw-audio/





Spectrum of Sequential Data



Future work





SpliceAl-toolkit

29













Calico

DeepMind







Calico

DeepMind





Troyanskaya Laboratory



DeepMind

Why Convolutional Neural Network to DNA?





Troyanskaya

Laboratory



DeepMind

Why Convolutional Neural Network to DNA?







Mihaela Pertea

Steven Salzberg

Anqi Liu

OpenSpiceAl : splice site predictor



Chao, K. H., Mao, A., Liu, Anqi, Salzberg, S. L., & Pertea, M. (**2024**). OpenSpliceAI. Manuscript in preparation. **[https://ccb.jhu.edu/spliceai-toolkit/**

Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** [https://ccb.jhu.edu/splam/

30

Can we predict splice sites using only DNA? Yes!

Х

Y



AGACTCAGCCCCCGGAGACTTAGTTAGAGGAAGAAAAAGGTAGGACAGAAGAAAAGGCAGGACATACAAGGTGCTGGCCCAGGGCGG





-**A**-

SpliceAI: splite site predictor









33









Splam

33 Future work

$$W = 5000$$
 $F = 10,000$



©penSp∛iceAl : better than SpliceAl!



OpenSpiceAI : retrain on different species В

Donor Precision 0.9 0.8 0.7 0.6 ⁸ ₀.5 - SpliceAl-Pytorch (mouse) 0.4 SpliceAl-Keras(Human) 0.3 400

С

1.0

0.9

0.8 o 0.7

0.6

5 0.5

å o.4

0.3

0.2

1.0 0.9

0.8

0.7

al 0.6

ō 0.5

0.4 0.3

0.2









Donor Precision Donor Recall

Splice site prediction metrics for arabadop

Splice site prediction metrics for mouse



Flanking Size







Splice site prediction metrics for zebrafish





Splice site prediction metrics for bee





36

Introduction

SpliceAl-toolkit

Splam

Future work

10000

OpenSpiceAl : Calibration



OpenSpiceAl : Calibration



0	0.2	0.4	0.6	0.8 Model predicted	1.0 probability ³⁸
Introduction		SpliceAl-toolkit	Splam		Future work

OpenSpiceAl : Calibration





0	0.2	0.4	0.6	0.8 Model predicted	1.0 I probability ³⁸
Introduction	>	SpliceAl-toolkit	Splam		Future work

OpenSp iceAl : new concept – Calibration



	0 0.2	0.4	0.6	0.8 Model predict	1.0 ed probability ¹⁵		
Introduc	tion	SpliceAl-toolkit	Splam		Future work		

OpenSpector in the second se





OpenSpiceAl : new concept – Calibration





OpenSpiceAl : Calibration 1.0 (fraction of chihuahua) Empirical probability 0.8 0.6 0.5 0.4 0.2 0.0 1.0 0.2 0.4 0.6 0.8 0 38 Model predicted probability SpliceAl-toolkit Introduction Splam Future work

OpenSpiceAl : Calibration 1.0 (fraction of chihuahua) 0.8 0.66 0.6 0.5 0.4 0.2 0.0 1.0 0.2 0.4 0.6 0.8 0 38 Model predicted probability SpliceAl-toolkit Introduction Future work Splam

Empirical probability











OpenSpiceAl : Calibration



Introduction

SpliceAl-toolkit





Alan Mao

OpenSpiceAI : benchmark

Introduction



Splam

Future work

43

©penSp∛iceAl : Summary

- 1. Data preprocessing: sliding window chunking
- 2. Easy-to-run framework to train your own SpliceAl
- 3. Pretrained SpliceAl-MANE
- 4. Pretrained SpliceAl on different species
- 5. Predict genetic variants' effect on splice sites



6. Model calibration: temperature scaling

Chao, K. H., Mao, A., Liu, Anqi, Salzberg, S. L., & Pertea, M. (<u>2024</u>). OpenSpliceAI. Manuscript in preparation. <u>https://ccb.jhu.edu/spliceai-toolkit/</u> (in preparation)

Splam

Future work

SpliceAl-toolkit

Is canonical labelling approach correct?






Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome Biology* in press. https://ccb.jhu.edu/splam/

Introduction

SpliceAl-toolkit

Splam

Future work





Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome Biology* in press. https://ccb.jhu.edu/splam/

Introduction

SpliceAl-toolkit

Splam

Future work

SPLMM : deep-learning splice site predictor



Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome Biology* in press. https://ccb.jhu.edu/splam/

Introduction

SpliceAl-toolkit

Splam

Future work

SPLMM : deep-learning splice site predictor

Interpretability: ablation study

Interpretability: input sequence



Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome Biology* in press. https://ccb.jhu.edu/splam/

Introduction

SpliceAl-toolkit

Splam

Future work



- Residual connection is powerful
- Grouped convolution helps (cardinality)



Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** [https://ccb.jhu.edu/splam/

Introduction

SpliceAl-toolkit

Splam

Future work

Future sequence models in genomics?

CNN or ?

Future sequence models in genomics?

CNN or/and Transformer?





sales*f*orce

Future? – Protein transformer-based models

Meta

DeepMind



DeepMind Calico **>**InstaDeep™ Arc Institute Future? – DNA transformer-based models



SpliceAl-toolkit

Application? – Genome annotation

Home / A-Z Publications / Annual Review of Genomics and Human Genetics / Early Publication / Review in Advance

ANNUAL REVIEW OF GENOMICS AND HUMAN GENETICS Deep Learning Sequence Models for Transcriptional Regulation

Ksenia Sokolova¹, Kathleen M. Chen¹, Yun Hao², Jian Zhou³, and Olga G. Troyanskaya^{1,2,4}

SegmentNT: annotating the genome at single-nucleotide resolution with DNA foundation models

 Bernardo P. de Almeida, Hugo Dalla-Torre, Guillaume Richard, Christopher Blum, Lorenz Hexemer, Maxence Gélard, Priyanka Pandey, Stefan Laurent, Alexandre Laterre, Maren Lang, Uğur Şahin, Karim Beguir, D Thomas Pierrot
doi: https://doi.org/10.1101/2024.03.14.584712

Part III

Fungi Language Model + Gene Expression Prediction



Foundation model

 Stanford researchers called transformers "foundation models" in an August 2021 paper because they see them driving a paradigm shift in Al.



• GPT-3, GPT-4 by

opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258.

Foundation model

- Versatility: wide range of downstream tasks
- **Transfer learning**: learn general representation of data. Task-specific is limited
- Efficiency: computational efficiency of fine-tuning models
- Generalization: "zero-shot" or "few-shot"
- Emergent abilities:
 - basic arithmetic
 - simple programming tasks
 - summarization, translation, or question-answering.

Goals

- Interpretable fungi LLM:
 - Strains / Species diversity
 - coding vs noncoding regions
 - repeat regions
- Predicting ChIP-exo, histone marks, and RNA-Seq
- Does fine-tuning a pretrained language model outperform training a new model from scratch?

Rossi, M. J., Kuntala, P. K., Lai, W. K., Yamada, N., Badjatia, N., Mittal, C., ... & Pugh, B. F. (2021). A high-resolution protein architecture of the budding yeast genome. *Nature*, *592*(7853), 309-314.





Fungi Language Model

• Q1: To what evolutionary distance should we include in our LM?

• Q2: What is the quality of the annotation? Coding vs non-coding regions

• Q3: How repetitive are the genomes?



Genome evaluation – Distance between species/strains

R64 Reference Yeast

80 strains of yeasts

165 Saccharomycetales







Genome evaluation – coding / noncoding regions





R64 Reference Yeast

Saccharomyces cerevisiae: 72.46%

Benegas, G., Batra, S. S., & Song, Y. S. (2023). DNA language models are powerful predictors of genome-wide variant effects. Proceedings of the National Academy of Sciences, 120(44), e2311219120.

Zhai, J., Gokaslan, A., Schiff, Y., Berthel, A., Liu, Z. Y., Miller, Z. R., ... & Kuleshov, V. (2024). Cross-species modeling of plant genomes at single nucleotide resolution using a pre-trained DNA language model. bioRxiv, 2024-06.

Genome evaluation – repeat regions





R64 Reference Yeast

Saccharomyces cerevisiae: 7.39%

Benegas, G., Batra, S. S., & Song, Y. S. (2023). DNA language models are powerful predictors of genome-wide variant effects. Proceedings of the National Academy of Sciences, 120(44), e2311219120.

Zhai, J., Gokaslan, A., Schiff, Y., Berthel, A., Liu, Z. Y., Miller, Z. R., ... & Kuleshov, V. (2024). Cross-species modeling of plant genomes at single nucleotide resolution using a pre-trained DNA language model. bioRxiv, 2024-06.





Data preprocessing



Coding regions



~ 7 genes per window



Data preprocessing



Coding regions





~ 7 genes per window



Data preprocessing



Coding regions







Language model performance – different models





Linder, J., Srivastava, D., Yuan, H., Agarwal, V., & Kelley, D. R. (2023). Predicting RNA-seq coverage from DNA sequence as a unifying model of gene regulation. Biorxiv, 2023-08.

Supervised ChiP-exo, histone marks, RNA-Seq prediction





Fine-tuning vs Training from Scratch (16 bp resolution)

- Divide genome into 8 folds.
- Train 8 models with distinct validation and test folds.





Fine-tuning vs Training from Scratch (4 bp resolution)



Fungi LM: Summary

- 1. Fungi language model: The Saccharomycetales order is a good evolutionary distance, offering good species diversity.
- 2. Repeat & coding down-weight masking are important
- 3. Under the exact model architecture, pretrained LM weights & finetuning can outperform training a model from scratch.





Acknowledgement



Steven Salzberg



Mihaela Pertea



Alaina Shumate



Alan Mao



Jakob Heinz



Kuanhao-Chao

Celine Hoh

Salzberg Lab

















Pertea Lab









Acknowledgement









Kuanhao-Chao

Johannes Linder Majed Mohamed Magzoub David Kelley

Sean Hackett

Kelley Lab & Calico Computing Team











Great mentors, collaborators and good friends!

