



PACIFIC  
BIOSCIENCES®

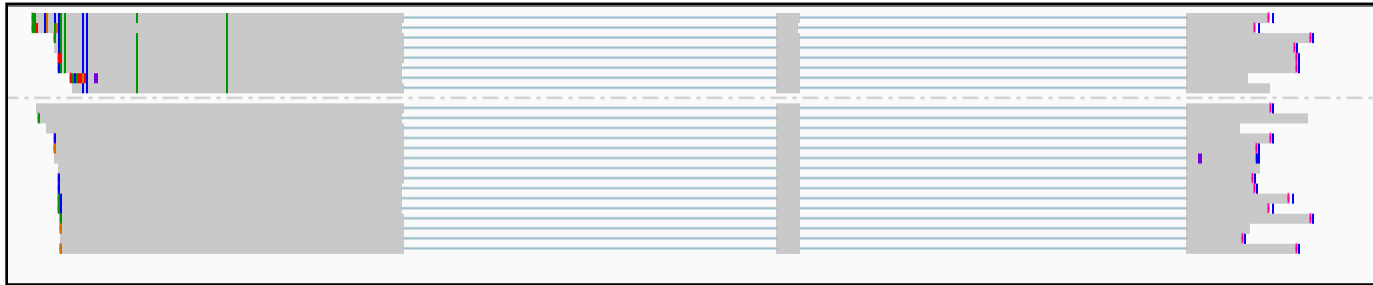
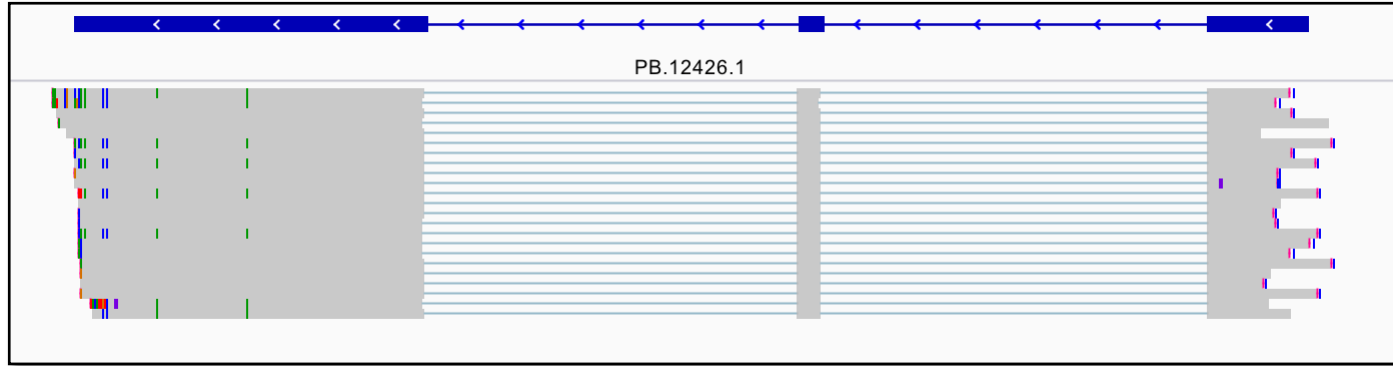


# IsoPhase: Isoform-Level Phasing

Elizabeth Tseng, Principal Scientist, PacBio

 @Magdoll

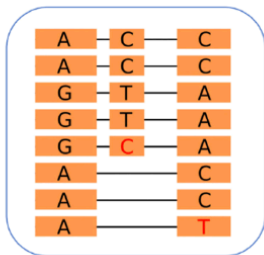
# ISOFORM-LEVEL PHASING USING ISO-SEQ READS



# ISOPHASE: ISOFORM-LEVEL PHASING

a

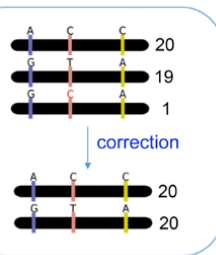
ALIGNMENT



SNP CALLING

Position	SNPs
POS1	A, G
POS2	C, T
POS3	C, A

PHASING



VCF OUTPUT

```
##fileformat=VCFv4.2
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT ISOFORM1 ISOFORM2
chr1 105 . A G . PASS DP=40;AF=0.50 GT:HQ 0|1:20,20 0:15
chr1 190 . C T . PASS DP=40;AF=0.50 GT:HQ 0|1:20,20 0:15
chr1 336 . C A . PASS DP=40;AF=0.50 GT:HQ 0|1:20,20 0:15
```



B73



male B73  
X  
female Ki11



Ki11



male Ki11  
X  
female B73

b

Pooled Reads

B73 FL reads  
Ki11 FL reads  
B73 x Ki11 FL reads  
Ki11 x B73 FL reads

IsoPhase

Allele and Per-sample Read Counts

	B73	Ki11	B73xKi11	Ki11xB73
	30	0	15	15
	0	20	15	3

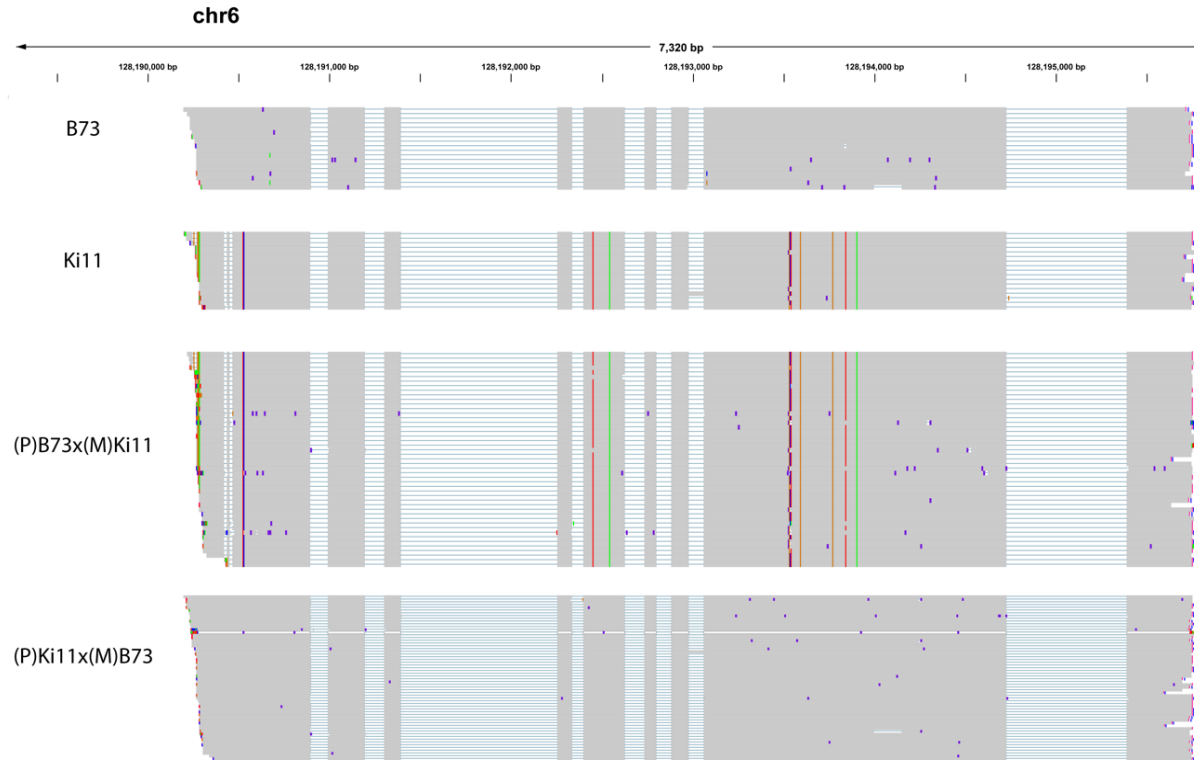
Final Results

Phase 0 (B73) : A-C-C  
Phase 1 (Ki11) : G-T-A  
B73 x Ki11 : 15, 15  
Ki11 x B73 : 15, 3

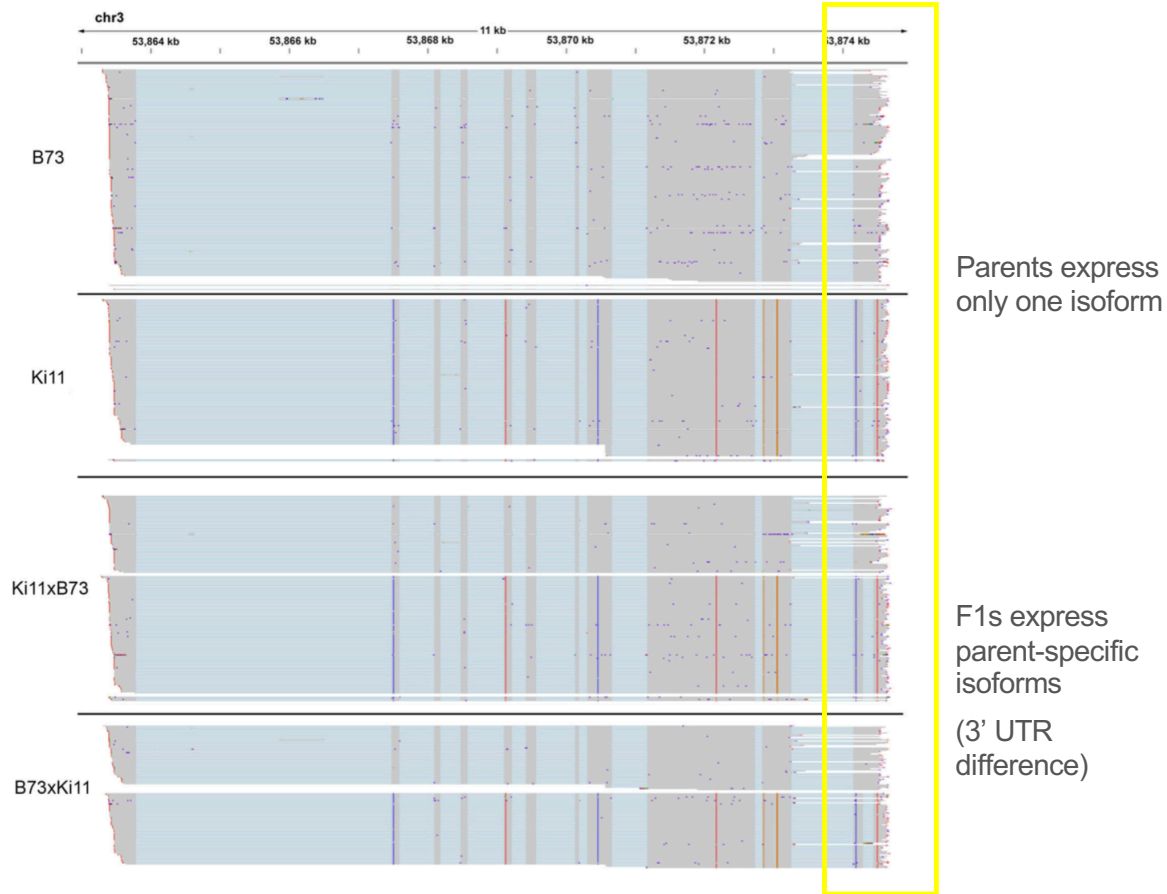
Identify B73-only allele as phase0  
Identify Ki11-only allele as phase1

IsoPhase is run individually for each gene

# MATERNAL IMPRINTING IN MAIZE



# PARENTAL-SPECIFIC ISOFORM EXPRESSION IN MAIZE



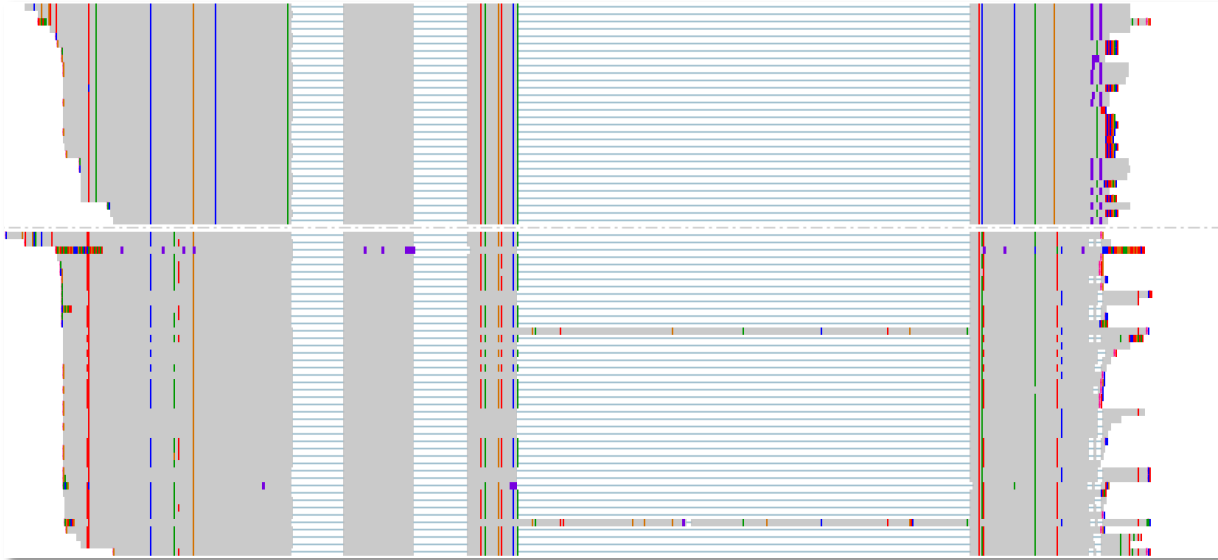
# PHASING BEYOND DIPLOID – A CHALLENGE

## Tetraploid potato poses phasing challenges

Each variant position has two SNPs

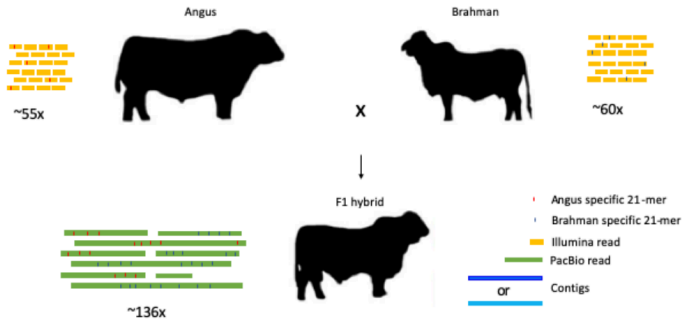
IsoPhase estimated three alleles, not four

Marko Petek  
& Kristina Gruden  
(NIB, Slovenia)





(this figure did not make it to the publication)

# TISSUE-SPECIFIC ALLELIC EXPRESSION IN CATTLE



Article | [Open Access](#) | Published: 29 April 2020

## Haplotype-resolved genomes provide insights into structural variation and gene content in Angus and Brahman cattle

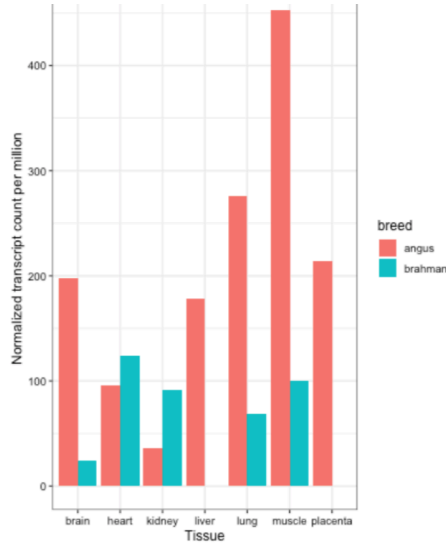
Wai Yee Low, Rick Tearle, Ruijie Liu, Sergey Koren, Arang Rhie, Derek M. Bickhart, Benjamin D. Rosen, Zev N. Kronenberg, Sarah B. Kingan, Elizabeth Tseng, Françoise Thibaud-Nissen, Fergal J. Martin, Konstantinos Billis, Jay Ghurye, Alex R. Hastie, Joyce Lee, Andy W. C. Pang, Michael P. Heaton, Adam M. Phillippy, Stefan Hiendleder , Timothy P. L. Smith  & John L. Williams 

*Nature Communications* 11, Article number: 2071 (2020) | [Cite this article](#)

# TISSUE-SPECIFIC ALLELIC EXPRESSION IN CATTLE

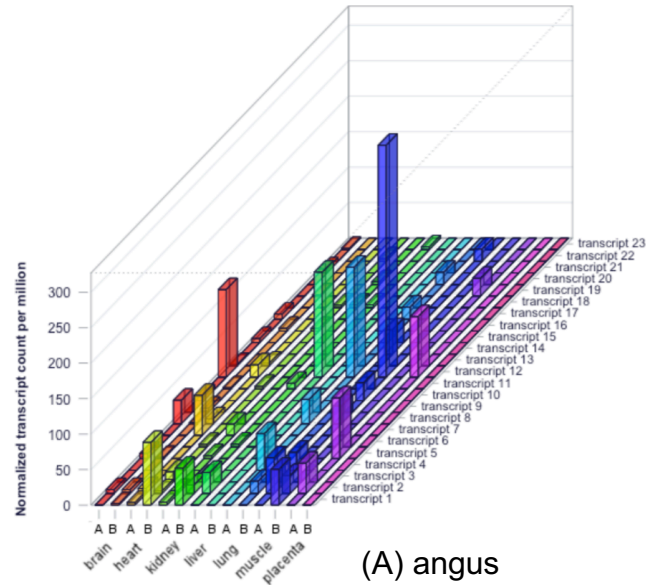
**ARIH2**

**gene level expression**



**ARIH2**

**isoform expression by tissue / allele**



(A) angus

(B) brahman



# ISOPHASE ON GITHUB (PART OF CUPCAKE)

## IsoPhase: Haplotyping using Iso Seq data

Elizabeth Tseng edited this page yesterday · 8 revisions

---

Last Updated: 09/26/2020

---

1. [Prerequisite](#)
  2. [Install IsoPhase](#)
  3. [What you will need for phasing](#)
  4. [Running IsoPhase](#)
  5. [Summarizing IsoPhase output](#)
- 

[https://github.com/Magdoll/cDNA\\_Cupcake/wiki/IsoPhase:-Haplotyping-using-Iso-Seq-data](https://github.com/Magdoll/cDNA_Cupcake/wiki/IsoPhase:-Haplotyping-using-Iso-Seq-data)



PACIFIC  
BIOSCIENCES®



# Cogent: Coding Genome Reconstruction Without A Reference Genome

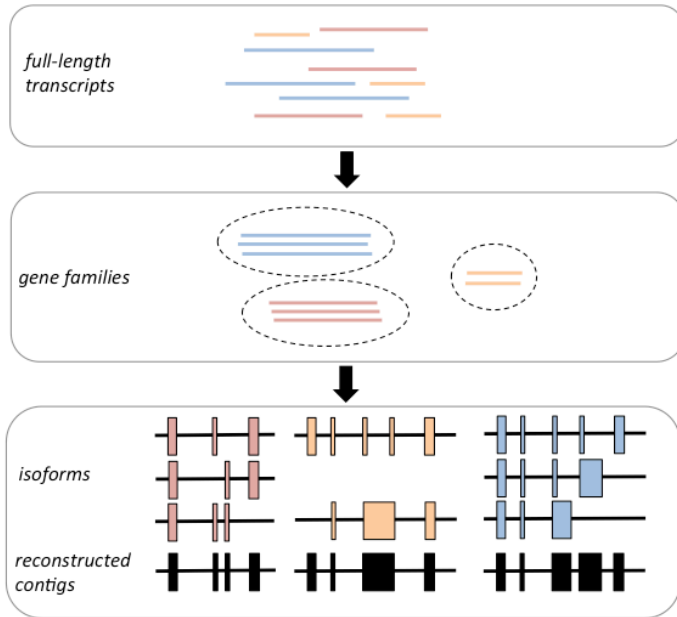
Elizabeth Tseng, Principal Scientist, PacBio

 @Magdoll

# NO GENOME? NO PROBLEM

## COGENT workflow

Using only Iso-Seq data to find gene families and reconstruct a fake “genome”

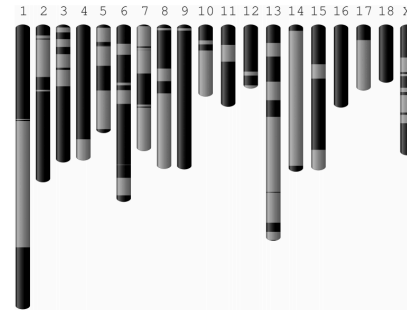


<https://github.com/Magdoll/Cogent>

## Use COGENT results to...

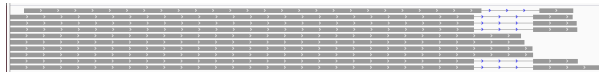
### Evaluate genome assemblies

Pig Iso-Seq Cogent rescued 5 missing genes for the new pig assembly



### Visualize alternative splicing

You can still see skipped exons!



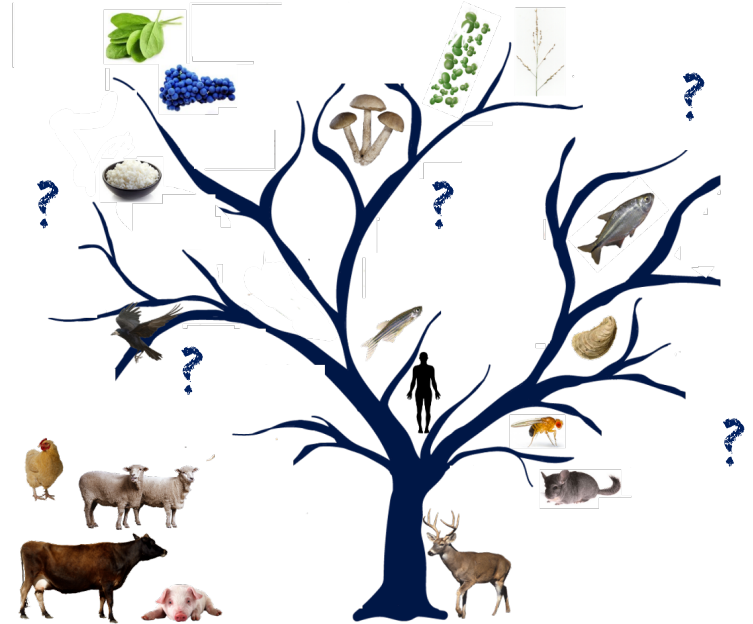
# COGENT: WHY?

## Not every species has a high-quality reference genome

If genome is poor, genome annotation (ab initio prediction & mapping) will suffer

## Iso-Seq bioinformatics does not require a genome

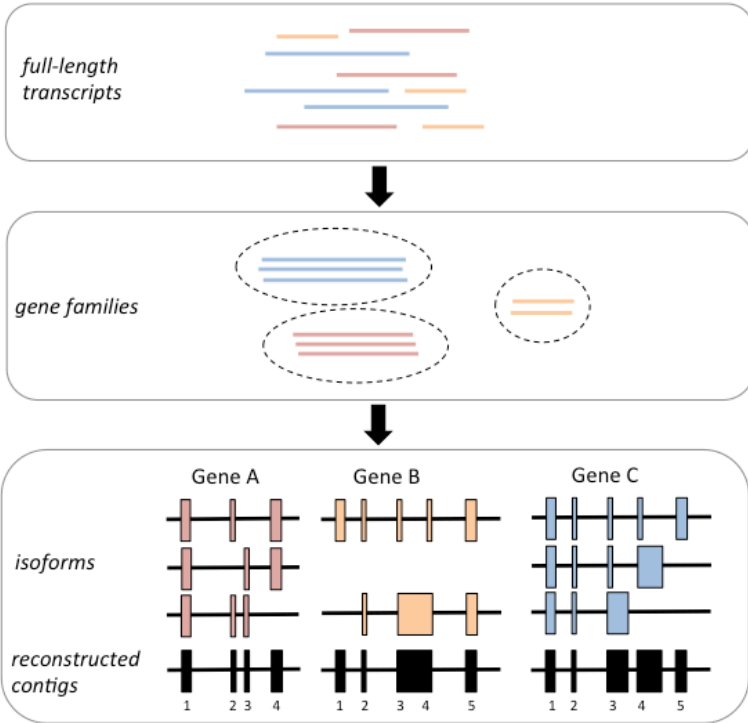
There is often enough information in Iso-Seq transcripts itself to identify gene families and the “coding” regions of the genome



# COGENT: HOW

## Cogent

COding GEName reconstruction Tool



Iso-Seq analysis generates full-length, high-quality ( $\geq 99\%$  accuracy) transcript sequences.

Transcripts are partitioned into gene families based their k-mer similarity.

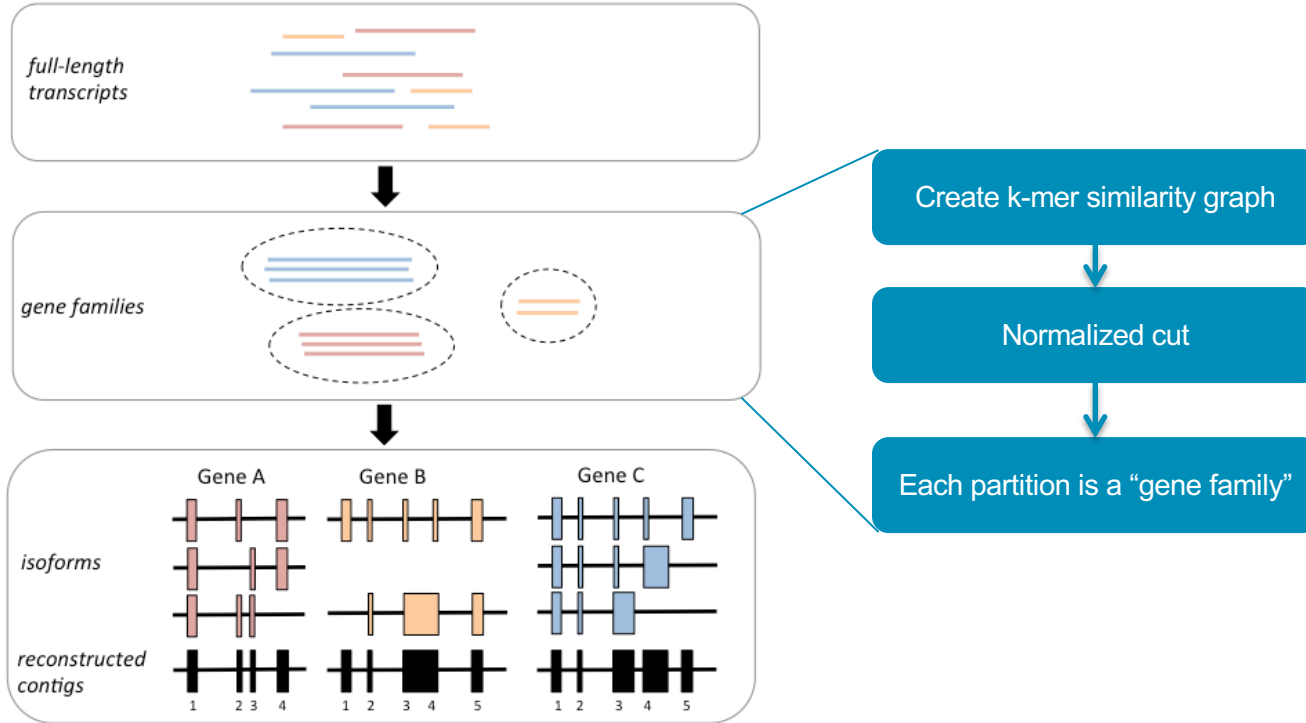
Each reconstructed contig represents the “union” of all coding bases in a particular gene locus.

(common introns will be invisible)

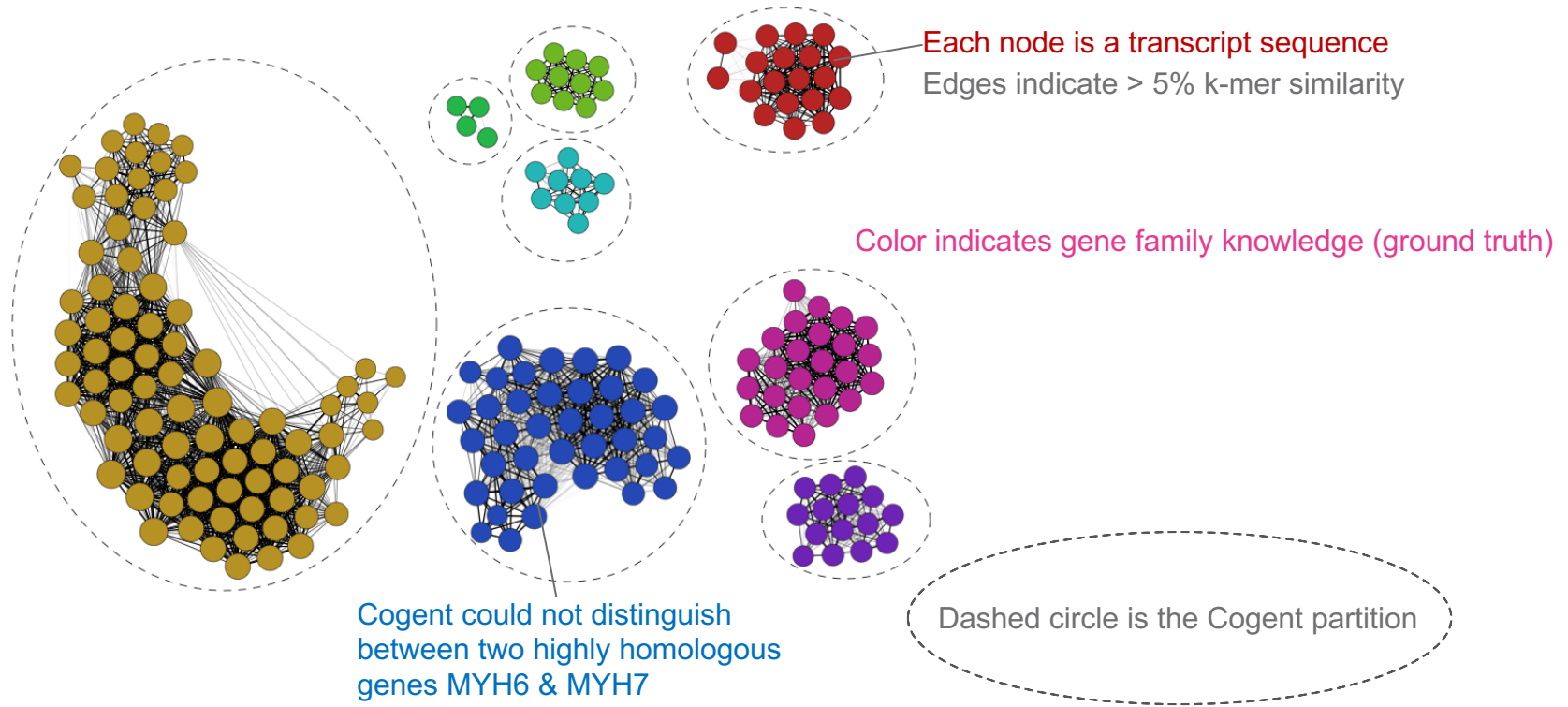
# COGENT: GENE FAMILY FINDING

## Cogent

COding GENome reconstruction Tool



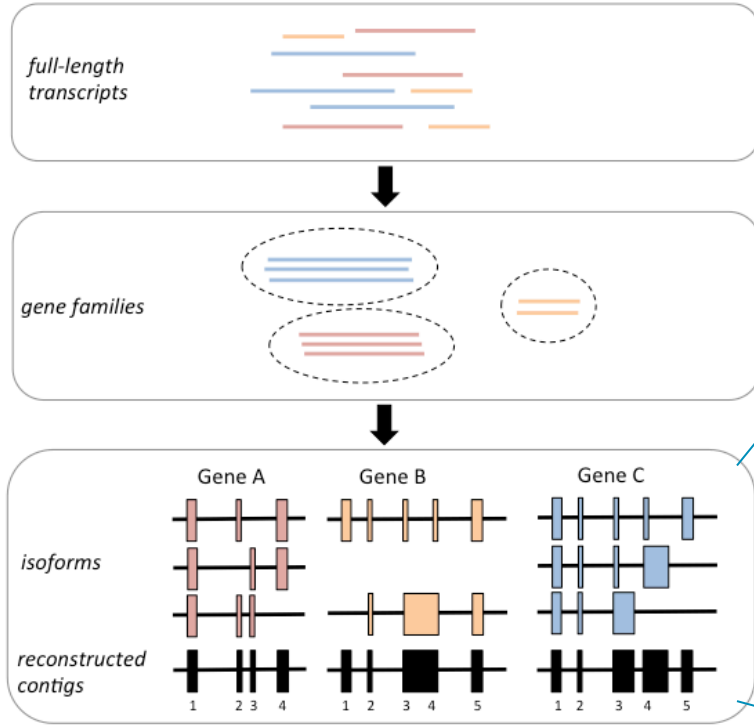
# GENE FAMILY PARTITIONING OF 9 HUMAN GENES



# COGENT: GENOME RECONSTRUCTION

## Cogent

COding GENome reconstruction Tool



for each gene family (partition):

Create de Bruijn graph

Collapse unipaths  
Resolve bubbles  
Resolve ends

Output reconstructed contig(s)

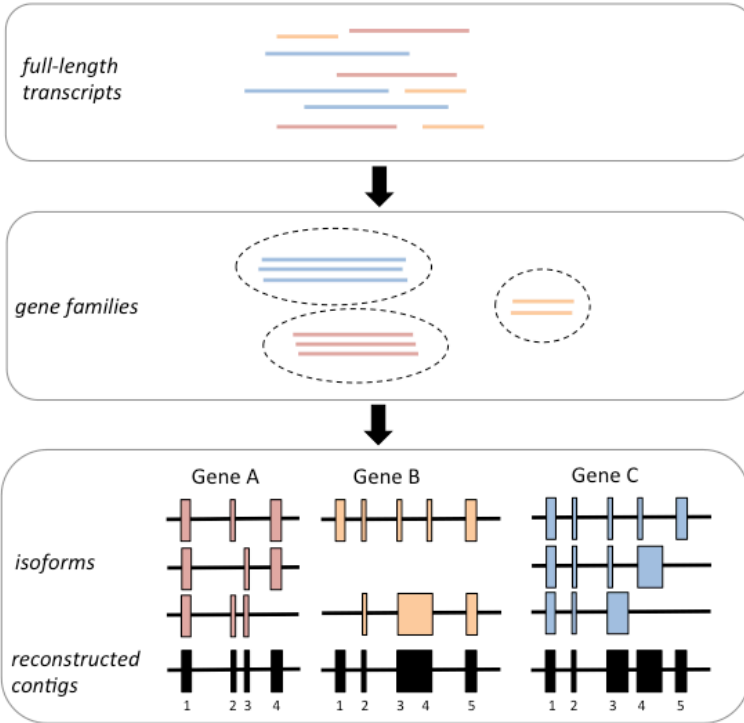
Full Cogent algorithm at <https://github.com/Magdoll/Cogent>



# COGENT: GENOME RECONSTRUCTION

## Cogent

COding GENome reconstruction Tool

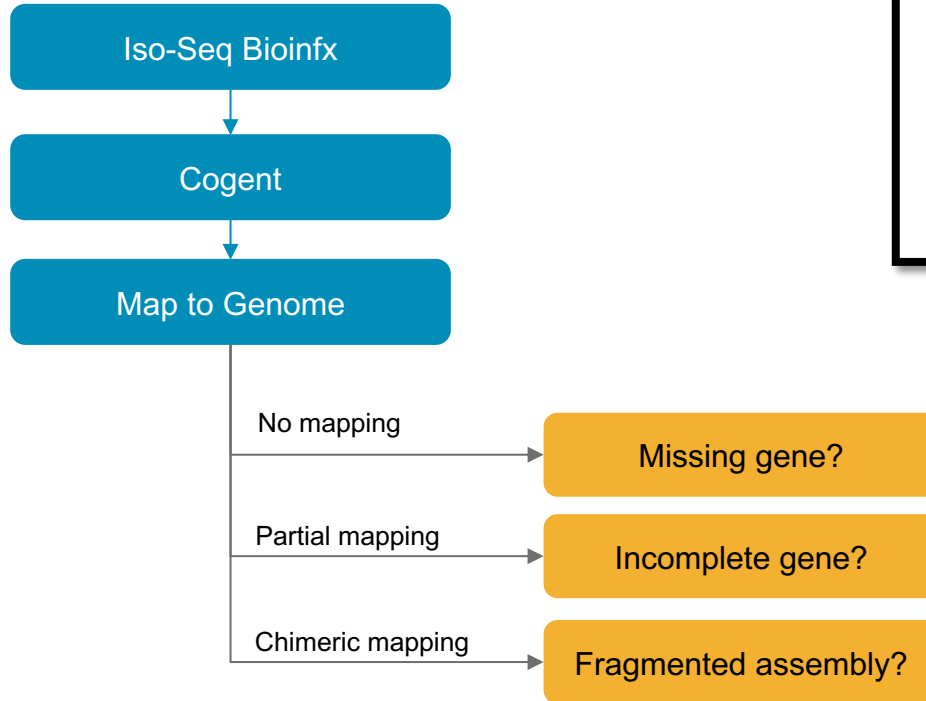


### Shortfalls:

May not separate paralogs

May reconstruct multiple contigs (case of gene C)

# COGENT CASE STUDY: PIG GENOME EVALUATION



## An improved pig reference genome sequence to enable pig genetics and genomics research

Amanda Warr, Nabeel Affara, Bronwen Aken, Hamid Beiki, Derek M Bickhart, Konstantinos Billis, William Chow, Lel Eory, Heather A Finlayson, Paul Flicek ... Show more

*GigaScience*, Volume 9, Issue 6, June 2020, giaa051,  
<https://doi.org/10.1093/gigascience/giaa051>

**Published:** 16 June 2020 **Article history** ▼

# COGENT CASE STUDY: PIG GENOME EVALUATION

Cogent Family	# of Iso-Seq Transcripts	Comment
6667_0	11	Missing CHAMP1 gene
16614_0	2	rotovirus
8757_0	30	Missing ERLIN1 gene
15567_0	8	Missing IRLIN1 gene
17496_0	17	Missing MB gene
17631_0	2	Missing PSD4 gene
17631_1	3	Missing PSD4 gene

## An improved pig reference genome sequence to enable pig genetics and genomics research

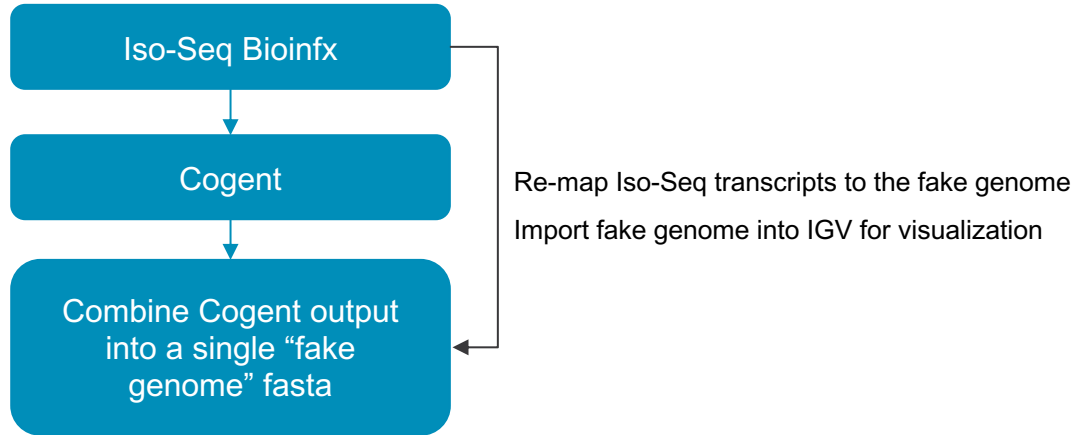
Amanda Warr, Nabeel Affara, Bronwen Aken, Hamid Beiki, Derek M Bickhart, Konstantinos Billis, William Chow, Lel Eory, Heather A Finlayson, Paul Flicek ... Show more

*GigaScience*, Volume 9, Issue 6, June 2020, giaa051,  
<https://doi.org/10.1093/gigascience/giaa051>

**Published:** 16 June 2020 **Article history** ▼

Five missing genes in the assembly were manually put back after Cogent evaluation

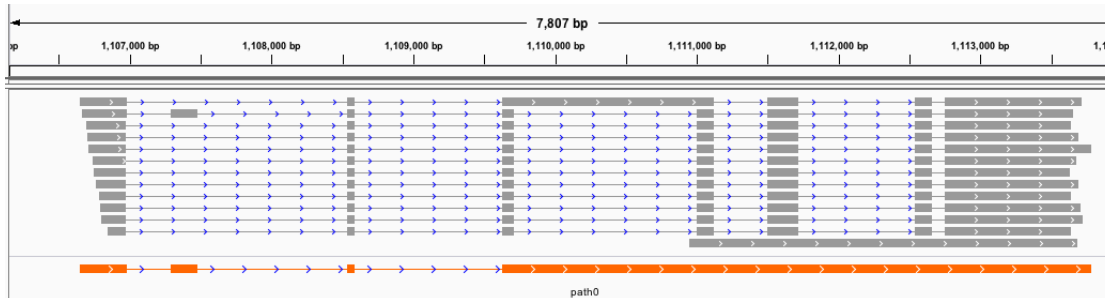
# COGENT CASE STUDY: MAKING A “FAKE CODING” GENOME



# COGENT CASE STUDY: MAKING A “FAKE CODING” GENOME

## Genome-based View

All introns are visible

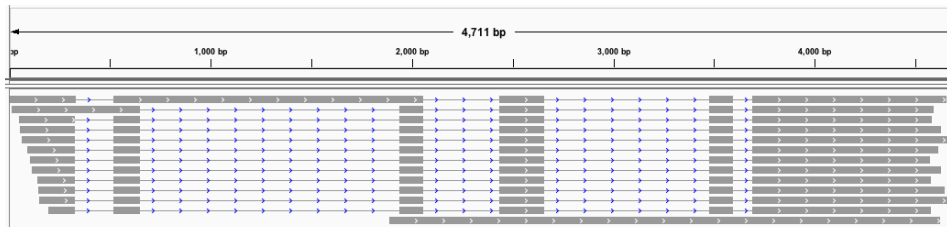


Iso-Seq Transcripts

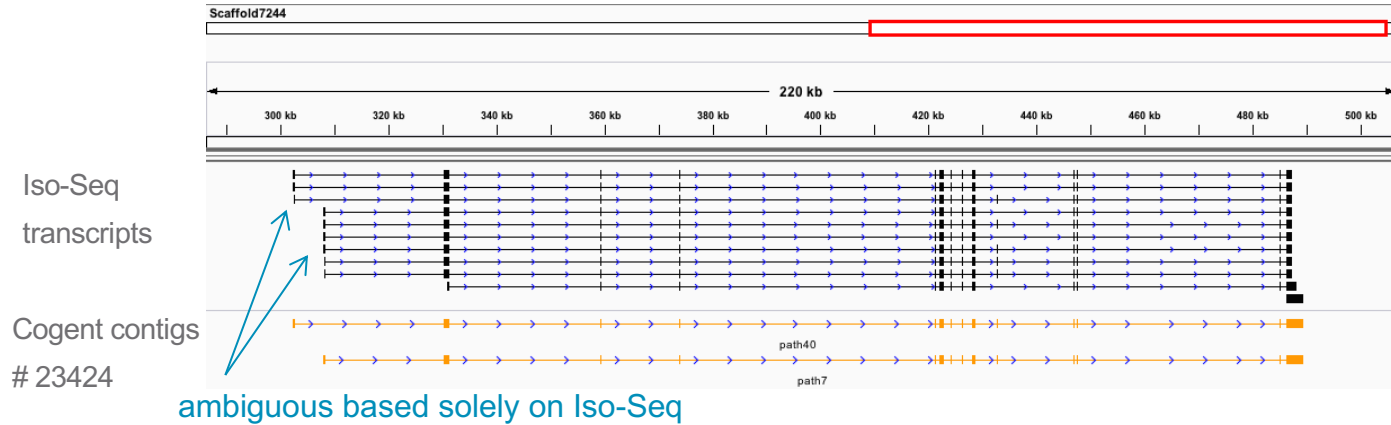
Cogent Output

## Cogent-based View

Introns that are never transcribed are not visible



# WHY COGENT MAY SOMETIMES OUTPUT >1 CONTIGS



- Lack of connectivity information between exon 1 and 2 based solely on transcripts
- Cogent outputs two contigs, one with exon 1 - 3, one with exon 2 – 3
- Mapping back to genome shows that the reconstruction is correct
- This is a case where genome information can be used to order exon 1 and 2

Coming soon: Special Cogent parameter to make “best guess” for ambiguity



PACBIO®

[www.pacb.com](http://www.pacb.com)

For Research Use Only. Not for use in diagnostic procedures. © Copyright 2020 by Pacific Biosciences of California, Inc. All rights reserved. Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT, SMRTbell, Iso-Seq, and Sequel are trademarks of Pacific Biosciences. Pacific Biosciences does not sell a kit for carrying out the overall No-Amp Targeted Sequencing method. Use of these No-Amp methods may require rights to third-party owned intellectual property. FEMTO Pulse and Fragment Analyzer are trademarks of Agilent Technologies Inc.

All other trademarks are the sole property of their respective owners.