

Hubert Pausch

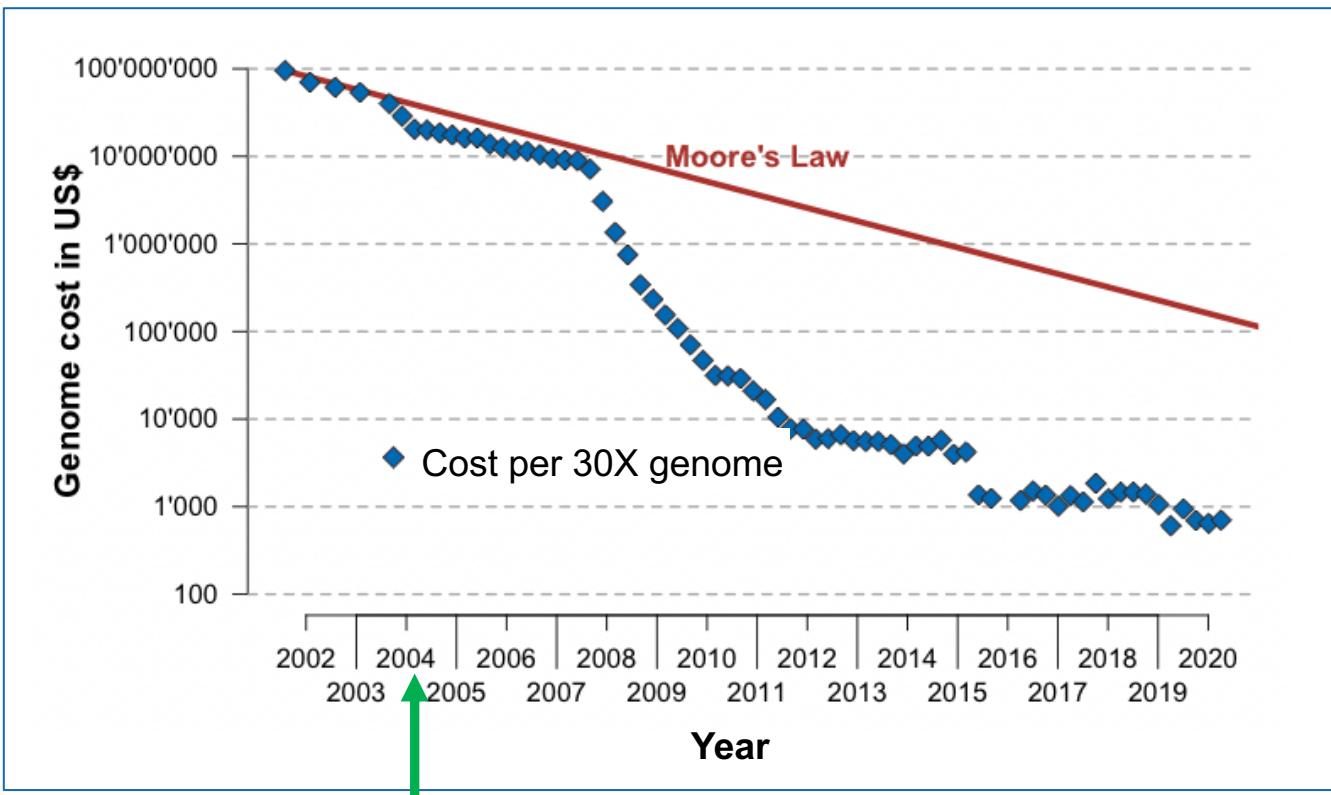
ETH Zürich

Animal Genomics | Institute of Agricultural Sciences

Towards establishing the bovine pangenome from multiple reference-quality assemblies

Zürich / Bern, 16 March 2021

Thousands of cattle genomes have been sequenced ...



Btau2.0

ARTICLES

nature
genetics

Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle

Hans D Daetwyler¹⁻³, Aurélien Capitan^{4,5}, Hubert Pausch⁶, Paul Stothard⁷, Rianne van Binsbergen⁸, Rasmus F Brøndum⁹, Xiaoping Liao⁷, Anis Djari¹⁰, Sabrina C Rodriguez⁴, Cécile Grohs⁴, Diane Esquerre¹¹, Olivier Bouchez¹¹, Marie-Noëlle Rossignol¹², Christophe Klopp¹⁰, Dominique Rocha⁴, Sébastien Fritz⁵, André Eggen⁴, Phil J Bowman^{1,3}, David Coote^{1,3}, Amanda J Chamberlain^{1,3}, Charlotte Anderson¹, Curt P VanTassell¹³, Ina Hulsegege⁸, Mike E Goddard^{1,3,14}, Bernt Guldbrandtsen⁹, Mogens S Lund⁹, Roel F Veerkamp⁸, Didier A Boichard⁴, Ruedi Fries⁶ & Ben J Hayes¹⁻³

LETTERS

<https://doi.org/10.1038/s41588-018-0056-5>

nature
genetics

Meta-analysis of genome-wide association studies for cattle stature identifies common genes that regulate body size in mammals

Aniek C. Bouwman¹, Hans D. Daetwyler^{2,3}, Amanda J. Chamberlain^{1,2}, Carla Hurtado Ponce^{2,4}, Mehdi Sargolzaei^{5,6}, Flavio S. Schenkel^{1,5}, Goutam Sahana^{1,7}, Armelle Govignon-Gion⁸, Simon Boitard⁹, Marlies Dolezal¹⁰, Hubert Pausch^{2,11,12}, Rasmus F. Brøndum⁷, Phil J. Bowman², Bo Thomsen⁹, Bernt Guldbrandtsen^{1,7}, Mogens S. Lund⁷, Bertrand Servin^{1,13}, Dorian J. Garrick^{1,14}, James Reecy¹⁴, Johanna Vilkki¹⁵, Alessandro Bagnato¹⁶, Min Wang^{1,2,3}, Jesse L. Hoff¹⁷, Robert D. Schnabel¹⁷, Jeremy F. Taylor¹⁷, Anna A. E. Vinkhuyzen^{18,19}, Frank Panitz¹⁹, Christian Bendixen⁹, Lars-Erik Holm¹⁹, Birgit Gredler²⁰, Chris Hozé^{8,21}, Mekki Boussaha⁸, Marie-Pierre Sanchez⁸, Dominique Rocha⁸, Aurelien Capitan^{8,21}, Thierry Tribout⁸, Anne Barbat⁸, Pascal Croiseau⁸, Cord Drögemüller^{1,22}, Vidhya Jagannathan²², Christy Vander Jagt², John J. Crowley²³, Anna Bieber²⁴, Deirdre C. Purfield²⁵, Donagh P. Berry²⁵, Reiner Emmerling²⁶, Kay-Uwe Götz²⁶, Mirjam Frischknecht²⁰, Ingolf Russ²⁷, Johann Sölkner²⁸, Curtis P. Van Tassell²⁹, Ruedi Fries¹¹, Paul Stothard³⁰, Roel F. Veerkamp¹, Didier Boichard^{1,8}, Mike E. Goddard^{2,4} and Ben J. Hayes^{1,2,3*}

... and aligned to ARS-UCD1.2

▪ First assembly of the bovine **reference sequence** (2004)

- Hereford cow «Dominette»
- Cost over 50 million US\$



The Genome Sequence of Taurine Cattle: A Window to Ruminant Biology and Evolution

The Bovine Genome Sequencing and Analysis Consortium,* Christine G. Elsik,¹
Ross L. Tellam,² Kim C. Worley¹

To understand the biology and evolution of ruminants, the cattle genome was sequenced to about sevenfold coverage. The cattle genome contains a minimum of 22,000 genes, with a core set of 14,345 orthologs shared among seven mammalian species of which 1217 are absent or undetected in nonruminant (marsupial or monotreme) genomes. Cattle-specific evolutionary breakpoint regions in chromosomes have a higher density of segmental duplications, enrichment of repeats, segments, and species-specific variants associated with evolution and ruminant uniqueness. Genes involved in metabolism are generally highly conserved, although five metabolic genes are deleted or extensively diverged from their human orthologs. The cattle genome sequence thus provides a resource for understanding mammalian evolution and accelerating livestock genetic improvement for milk and meat production.

Domesticated cattle (*Bos taurus* and *Bos tauri indicus*) provide a significant source of nutrition and livelihood to nearly 6 billion humans. Cattle belong to a clade phylogenetically distant from humans and rodents, the Cetartiodactyl order of eutherian mammals, which

2009 VOL 324 SCIENCE www.sciencemag.org



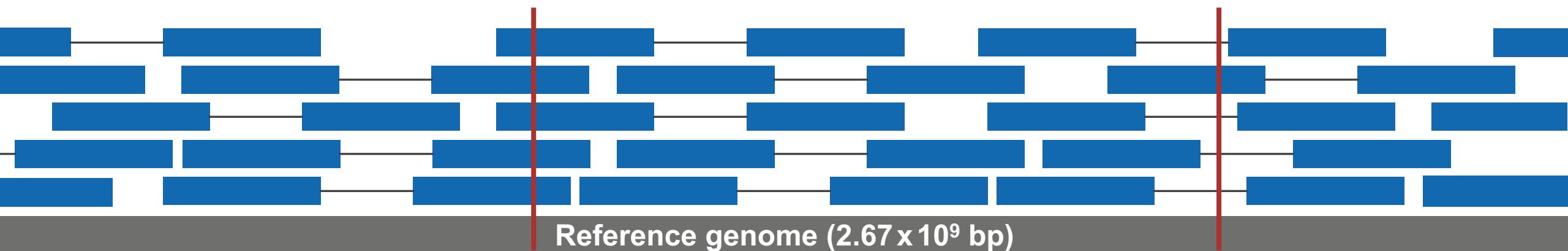
Source: Michael MacNeil, USDA

Typically, more than 100 million sequence read pairs (2 x 150 bases) for each individual corresponding to >10-fold genome coverage

```
CCCCCATGTAGCACAGCTCTGGAGGCTCTCACCACGTAGGAAATGTAGCGTCTCTGTAGAGAACGCTGCCGGATGGTCGGGGGGAGAAAGCGGATC  
GTGAAGAACCTGCCTCTGATGAAACTGCCTGTACGGAGTCTCGAAGTTCAAAGGCGATGGGAGGTTGAGACCATCAAACCGAGGCCTCTGAGAGCCTG  
TCACTGATGACCTTGTTGAGCCATCGTCGCCCTCGATGCCACGCTGTCCAGGATCTCTTGATGTCCTTGGCGCTAGGAGAAGTATTGCCCGAG  
. . .  
CTGAAAGTTCCACATGTGGAATATAGATACAACATTGAAACAAAATATGTGGCCTCCCATGTACATTGGTTACCTATGTACAAGTATCCTATACACCAAGTA  
CTTGGAAAATCCATGTTCTCCACCCGTGTTCTGTCTGCCCTGTCTGTAGCACAAACTGTGACTCTGAGAACGCAGACTGCCGTGGCCCTGTGCCGTCAGG  
GATCGGAAGAGCACACGTCTGAACCTCCAGTCACCCGTCGGATCTCGTATGCCGTCTCTGTTGAAAAAACAAAAAGAGAGTCTTGCAGGCATGCCCTG
```

5-fold coverage

1-fold coverage



Genetic variation in molecules and matrices

Thousands of animals ↓

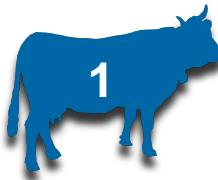
Millions of SNPs →

	SNP1 # of T 's	SNP2 # of C 's	SNP3 # of T 's
1	2	1	2
2	1	2	1
3	0	0	0

SNP1 **SNP2** **SNP3**

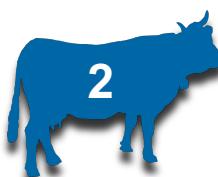
...TAT**T**CACTTTTAAT**G**TTTCATTAAAGT**T**AT...

...TAT**T**CACTTTTAAT**C**TTTCATTAAAGT**T**AT...



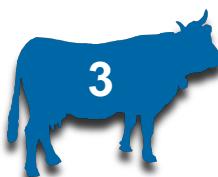
...TAT**T**CACTTTTAAT**C**TTTCATTAAAGT**G**AT...

...TAA**C**ACTTTTAAT**C**TTTCATTAAAGT**T**AT...



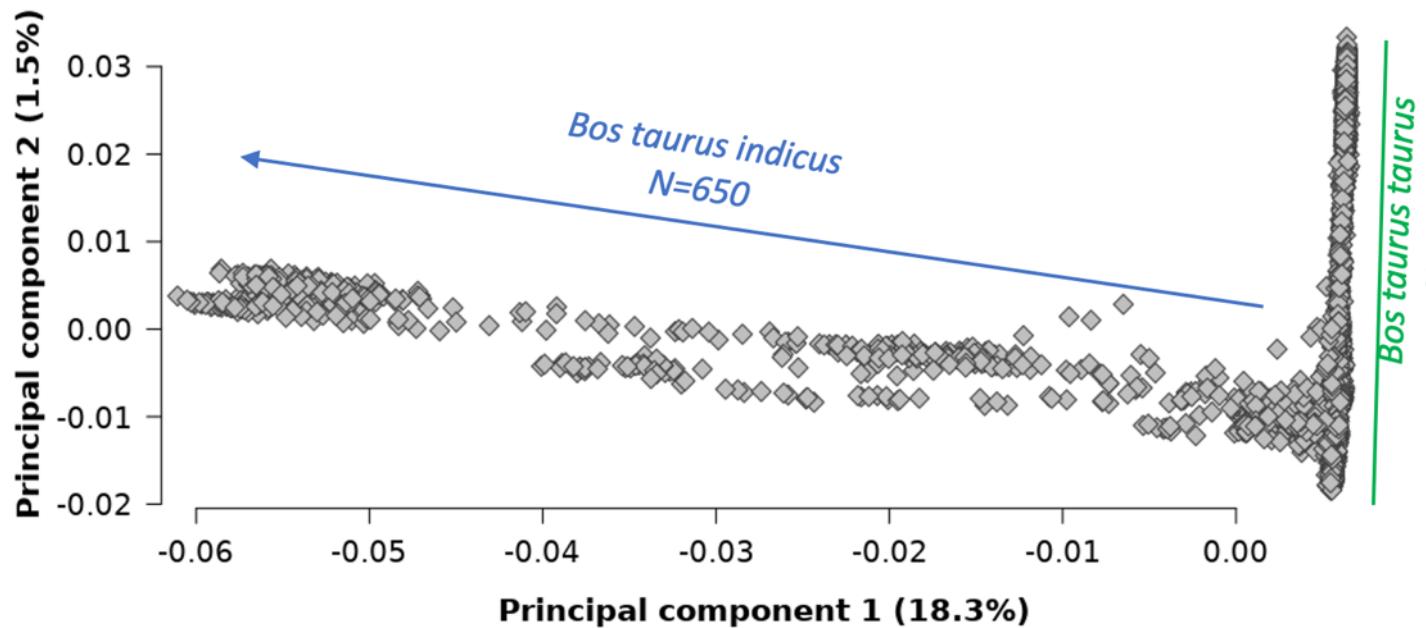
...TAA**C**ACTTTTAAT**G**TTTCATTAAAGT**G**AT...

...TAA**C**ACTTTTAAT**G**TTTCATTAAAGT**G**AT...

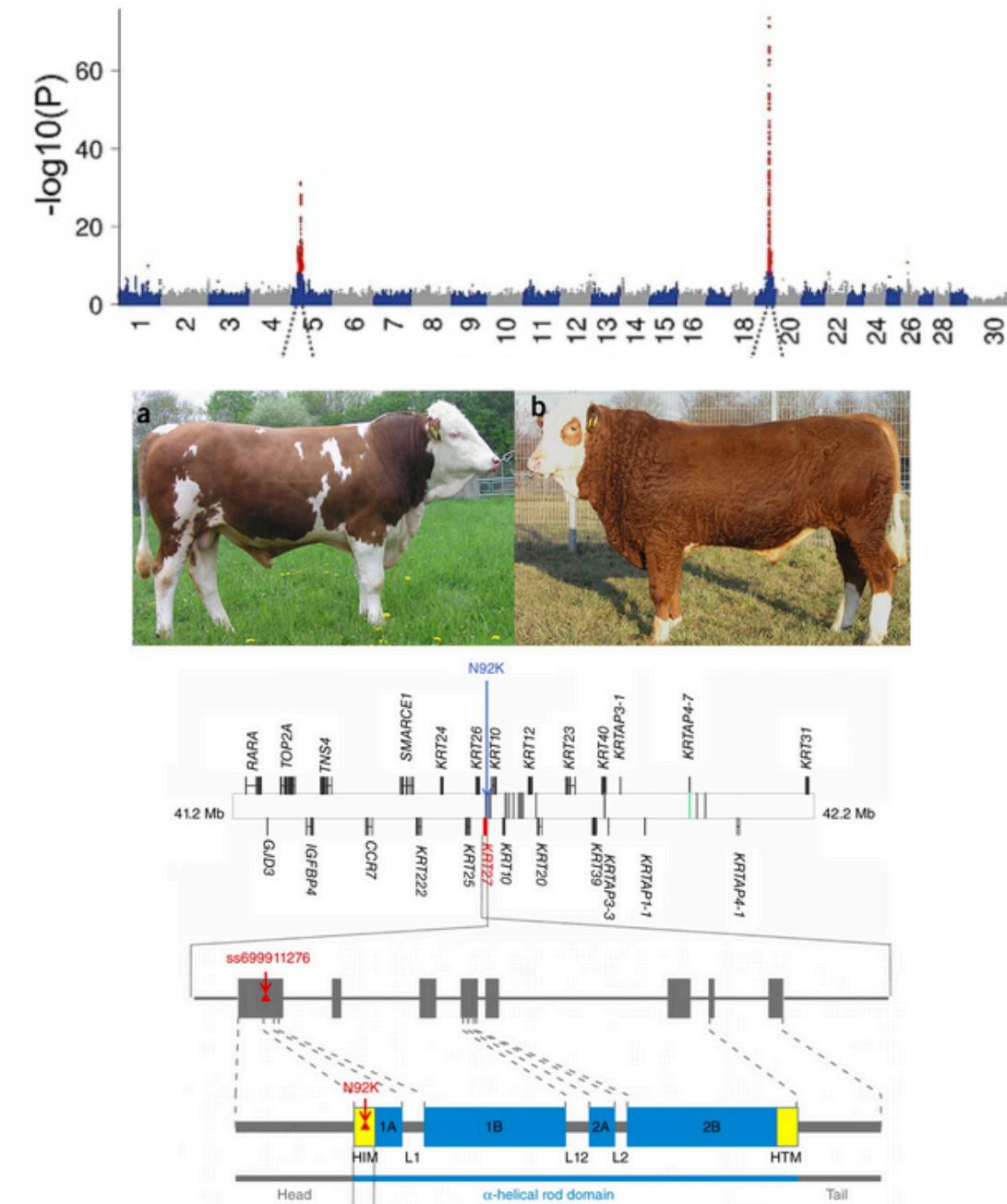
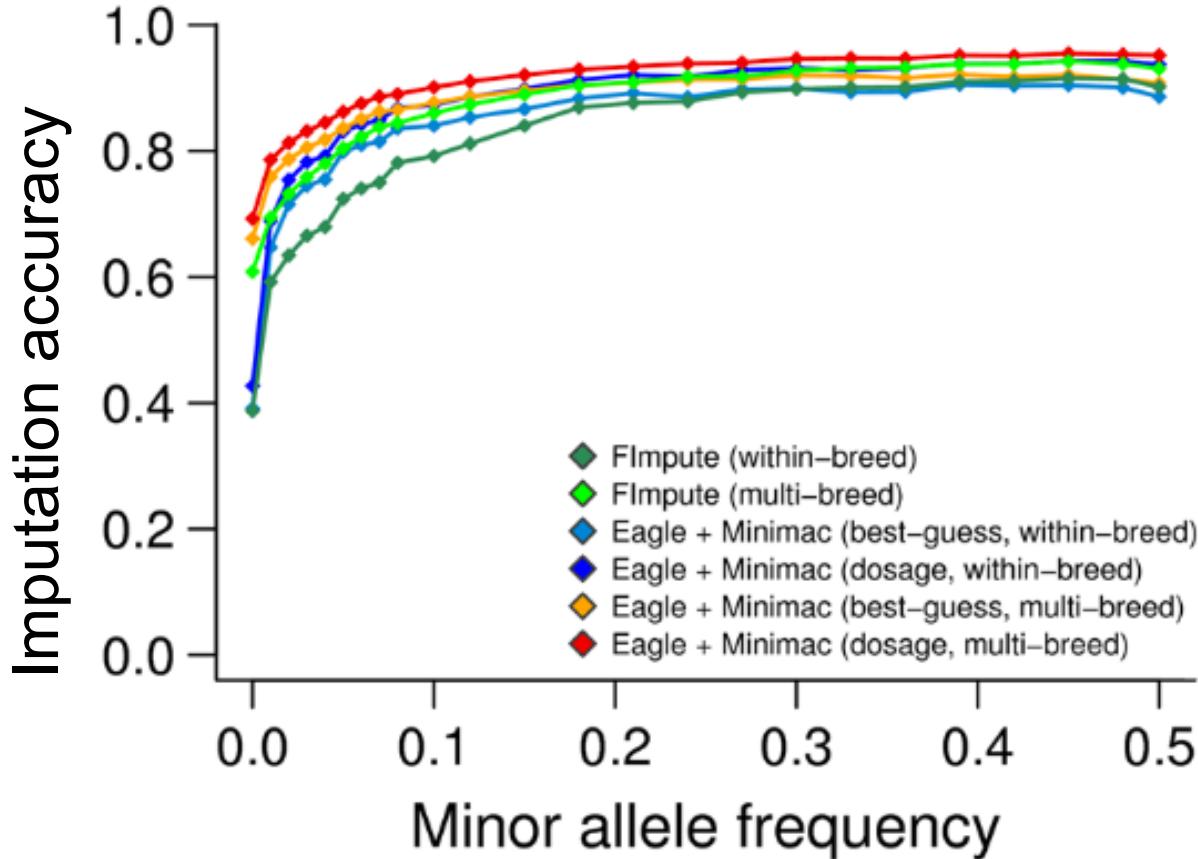


...TAT**T**CACTTTTAAT**C**TTTCATTAAAGT**T**AT... Reference
(32 out of 2,671,695,104 nucleotides of the bovine sequence)

Reference-guided variant discovery revealed >140 million polymorphic sites in >4700 cattle (1000 Bull Genomes Project)



This catalogue contains causal variants



Reference-guided variant discovery is biased

- DNA sequencing reads that contain only alleles that match corresponding reference nucleotides are more likely to align correctly than DNA fragments that also contain non-reference alleles
- Reads originating from DNA fragments that are highly diverged from corresponding reference nucleotides will either obtain low alignment scores, or align at incorrect locations, or remain un-mapped
- Reads originating from DNA fragments that are missing in the reference genome remain unmapped. Reference-guided variant discovery is blind to these sites.

Reference allele bias

Paten et al. *Genome Res* 2017;27:665–76

Pritt et al. *Genome Biol* 2018;19:220

Sherman et al. *Nat Genet* 2019;51:30–35

Strategies to mitigate reference allele bias

- *Personalized reference genomes*
 - Easy to implement
 - Standard genomic analysis tools work well
 - Coordinates remain unchanged
- *Graph-based references*
 - Relatively easy to implement
 - Standard genomic analysis tools work with workarounds
 - Variant prioritization is crucial
 - Coordinates remain unchanged

A *Bos taurus* reference graph



- ARS-UCD1.2 used as backbone
- Augmented with variants (SNPs and Indels) that were prioritized based on allele frequency
 - Breed-specific genome graphs
 - Pan-genome graphs
- Read mapping and variant genotyping accuracy compared between graphs and the linear reference sequence

Crysanto and Pausch *Genome Biology* (2020) 21:184
<https://doi.org/10.1186/s13059-020-02105-0> Genome Biology

RESEARCH Open Access

Bovine breed-specific augmented reference graphs facilitate accurate sequence read mapping and unbiased variant discovery

Danang Crysanto * and Hubert Pausch

Check for updates

Conclusions: We develop the first variation-aware reference graph for an agricultural animal (<https://doi.org/10.5281/zenodo.3759712>). Our novel reference structure improves sequence read mapping and variant genotyping over the linear reference. Our work is a first step towards the transition from linear to variation-aware reference structures in species with high genetic diversity and many sub-populations.

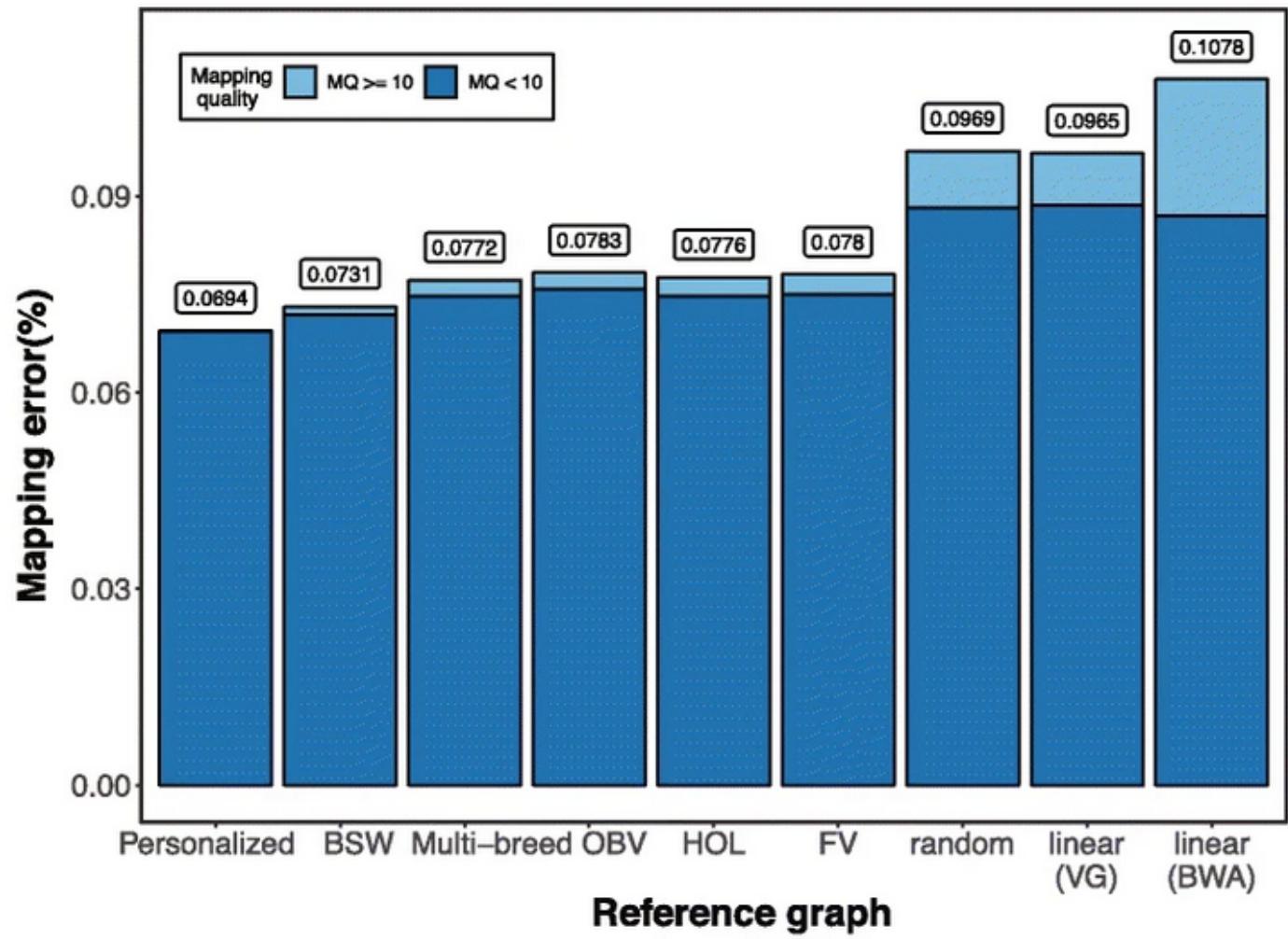
Keywords: Variation-aware genome graph, Sequence variant genotyping, Reference allele bias



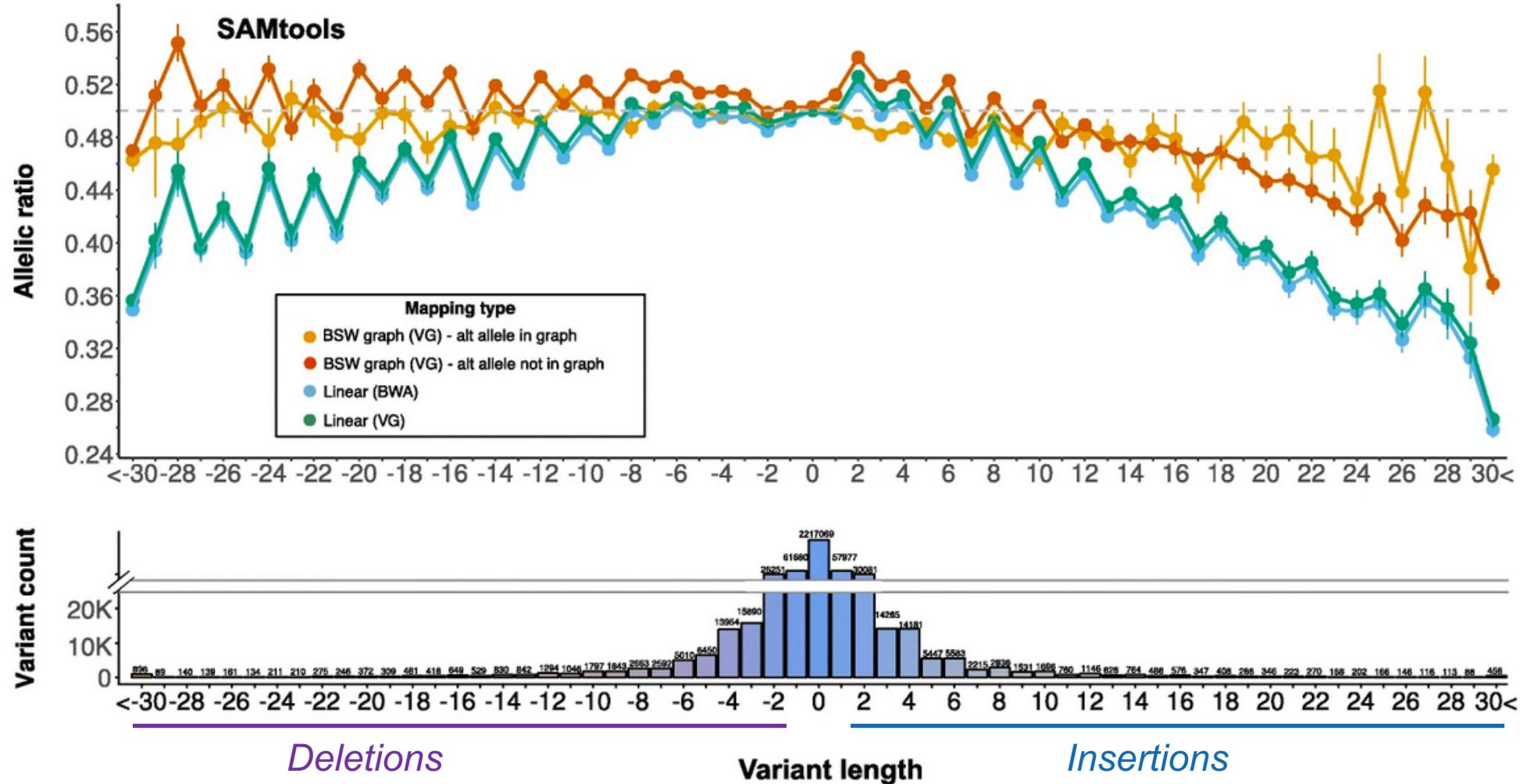
<https://github.com/danangcrysanto/bovine-graphs-mapping>

Linear mapping is outperformed by graph-based mapping

Breed-specific graph reduced mapping errors by 30% compared to current gold standard (*BWA mem*)



Graph-based variant genotyping is unbiased



Options to mitigate reference allele bias

- *Personalized reference genomes*
 - Easy to implement
 - Standard genomic analysis tools work well
 - Coordinates remain unchanged
- *Graph-based references*
 - Relatively easy to implement
 - Standard genomic analysis tools work with workarounds
 - Variant prioritization is crucial
 - Coordinates remain unchanged
- *Multi-assembly graphs*
 - Full spectrum of sequence diversity accessible
 - New paradigm that requires novel tools (*minigraph*, *progressive cactus*, *pggb*)

Largely blind to structural variations



The assembly of reference-quality assemblies is possible *at scale*

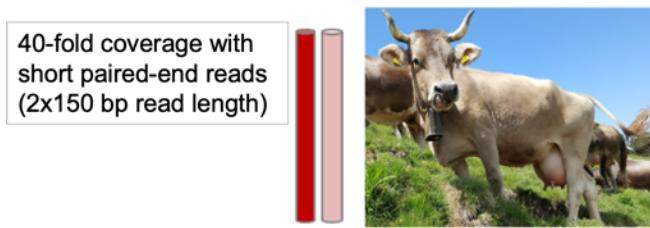
- Long-read sequencing is affordable
 - ONT, PacBio
- Sophisticated assembly algorithms have been developed
 - Trio-Binning (*Canu*, *hifiasm*, *GraphUnzip*)
 - Noisy reads (*Shasta*, *Raven*, *Ratatosk*)
- Long-read sequencing is highly accurate
 - PacBio Circular Consensus Sequencing (a.k.a. «HiFi-sequencing»)
 - PEPPER + Margin + DeepVariant for variant calling from ONT reads

Bovine HiFi-assemblies



contiguity *correctness* *completeness*

Animal	Size (gb)	Contigs	NG50 (mb)	k-QV	BUSCO (% complete)	Assembler
Original Braunvieh (primary)	3.17	765	91	50.9	96.1	Hifiasm
Brown Swiss (F1 haplotype)	3.07	1036	86.6	46.2	95.9	Hifiasm
Nellore (F1 haplotype)	2.95	1177	94.7	45.6	93.3	Hifiasm



Brown Swiss cow #238
Bos taurus taurus

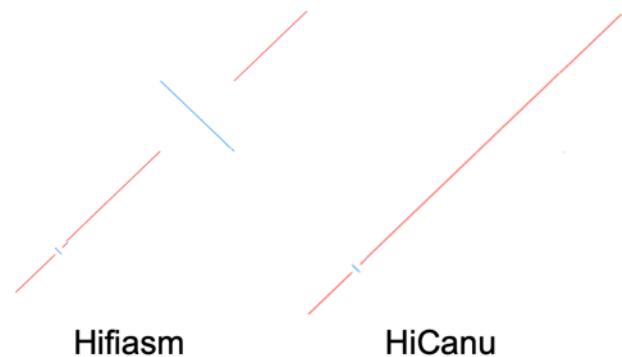


Nellore bull «Nuevo»
Bos taurus indicus



46-fold genome coverage with HiFi reads
(21 kb read length)
90-fold with ONT reads (50 kb read length)

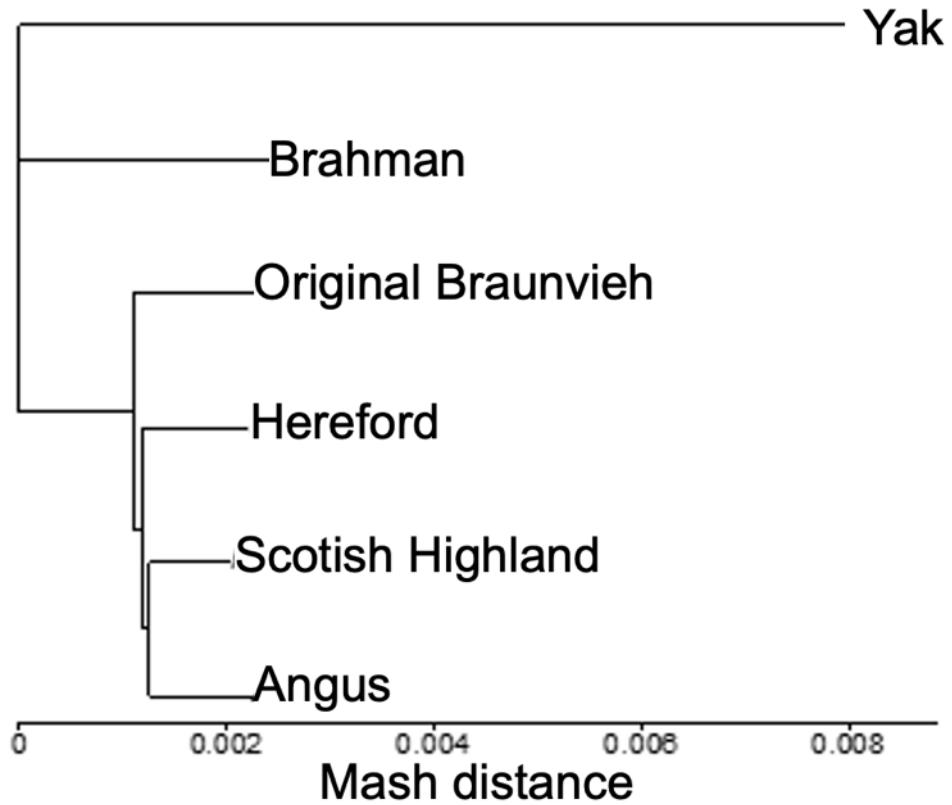
F1
*25.02.2020



<https://github.com/AnimalGenomicsETH/bovine-assembly>

Establishing a bovine pangenome

(*Bos taurus taurus*, *Bos taurus indicus*, *Bos grunniens*)



The pangenome uncovers novel functional sequences



bioRxiv

THE PREPRINT SERVER FOR BIOLOGY

[HOME](#) | [ABOUT](#) | [SUBMIT](#) | [NEWS & NOTES](#) | [ALERTS / RSS](#) | [CHANNELS](#)

Search



[Advanced Search](#)

bioRxiv is receiving many new papers on coronavirus SARS-CoV-2. A reminder: these are preliminary reports that have not been peer-reviewed. They should not be regarded as conclusive, guide clinical practice/health-related behavior, or be reported in news media as established information.

New Results

[Comment on this paper](#)

[Previous](#)

[Next](#)

Novel functional sequences uncovered through a bovine multi-assembly graph

Danang Crysanto, Alexander S. Leonard, Zih-Hua Fang, Hubert Pausch

doi: <https://doi.org/10.1101/2021.01.08.425845>

This article is a preprint and has not been certified by peer review [what does this mean?].

Posted January 08, 2021.

[Download PDF](#)

[Supplementary Material](#)

[XML](#)

[Email](#)

[Share](#)

[Citation Tools](#)

[Abstract](#)

[Full Text](#)

[Info/History](#)

[Metrics](#)

[Preview PDF](#)

[Tweet](#)

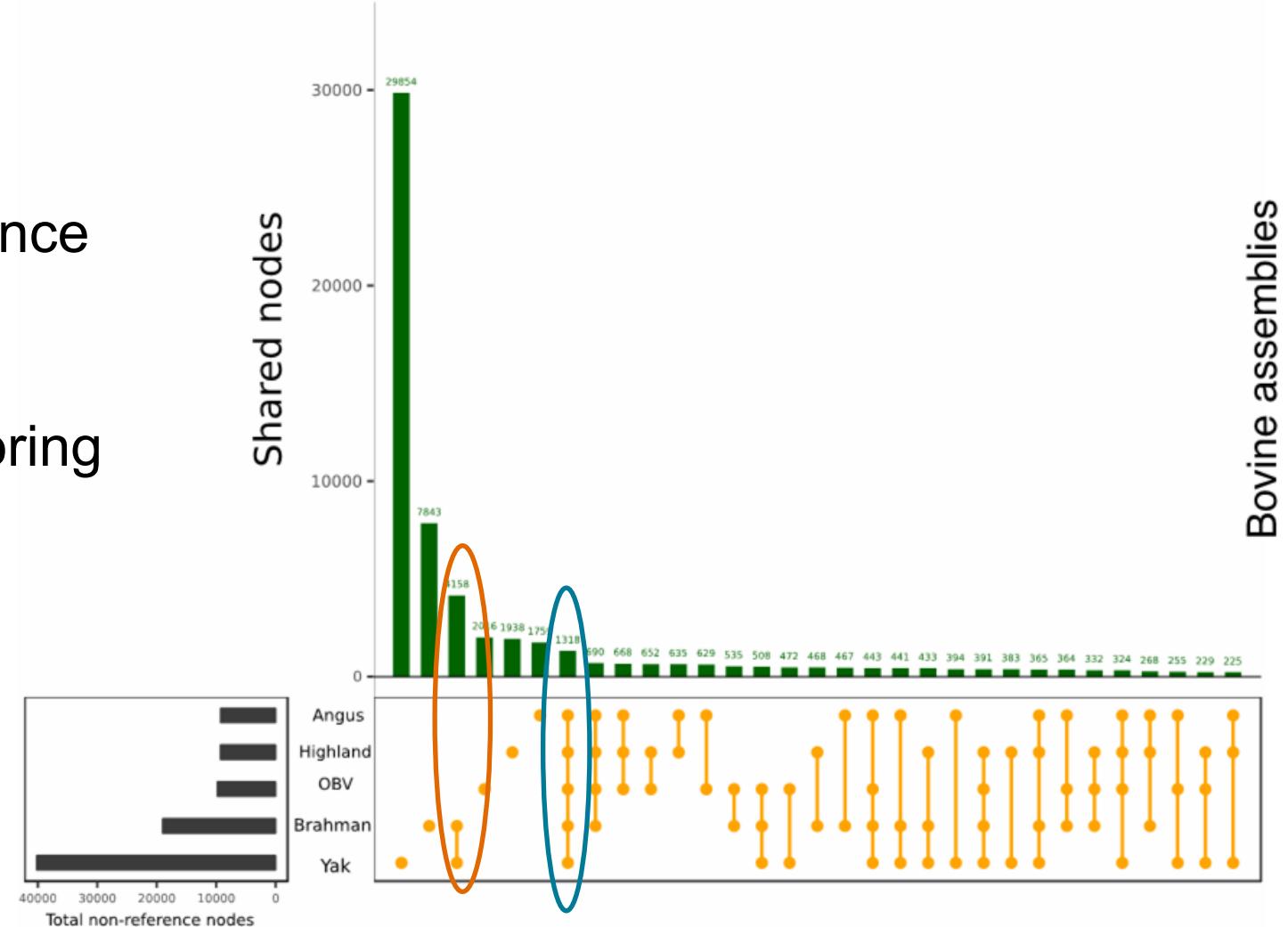
[Like 0](#)



<https://github.com/AnimalGenomicsETH/bovine-graphs>

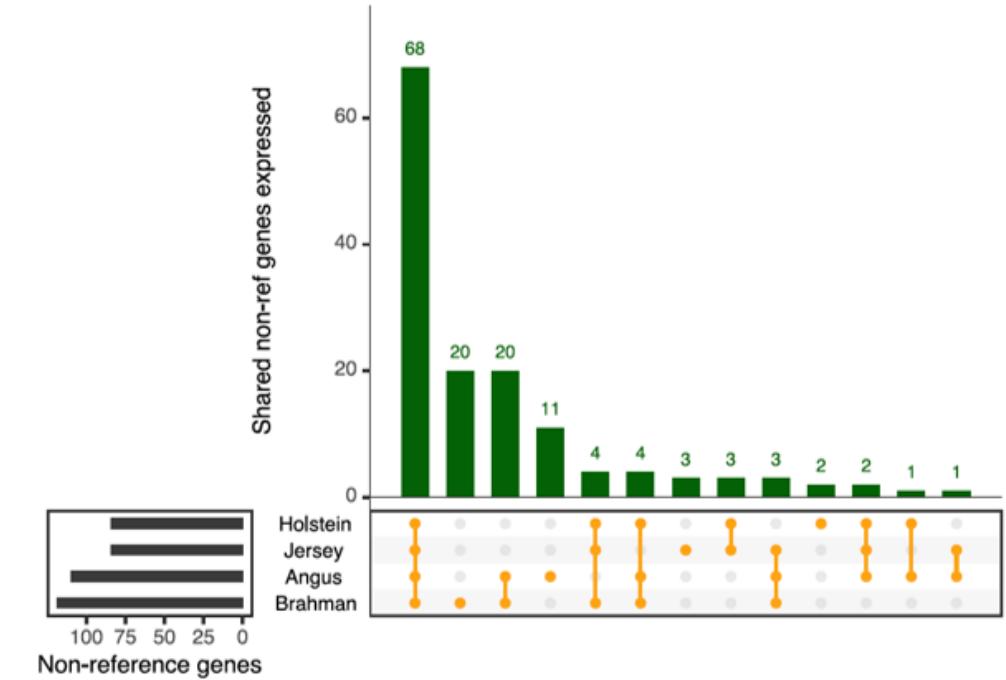
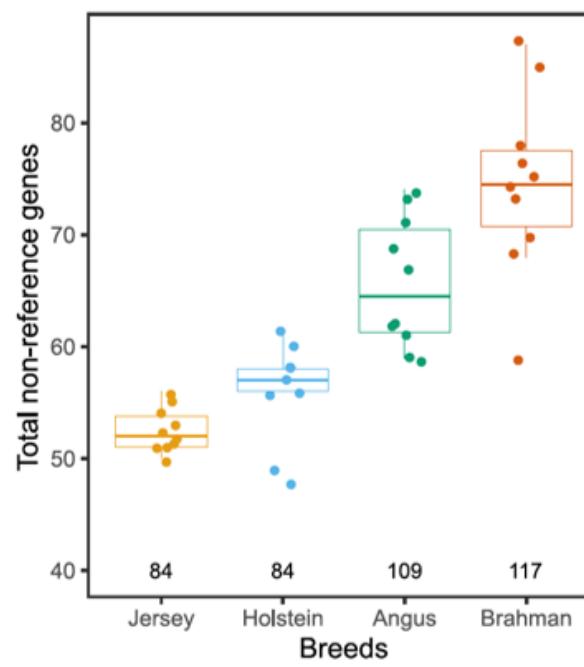
Structural variants (*non-reference nodes*) of the multi-assembly graph contain 70 million bases

- Identify nonreference sequence
 - Evolutionary relationships?
 - Reference bias?
- Highlights areas worth exploring



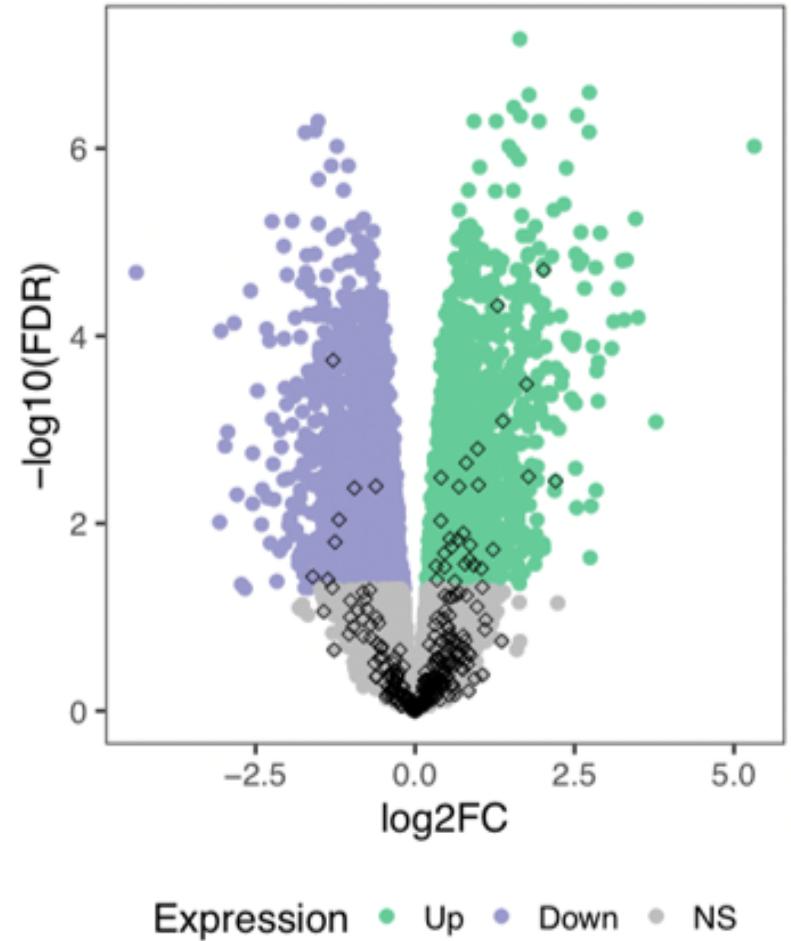
Novel sequences contain novel genes

- Gene models were predicted using Augustus & BLASTP
- Liver transcriptomes from 45 cattle provide additional experimental evidence



Novel genes are differentially expressed

- We revisit an analysis from McLoughlin et al (*Frontiers in immunology* 5, 396. (2014))
 - Transcriptome data of peripheral blood leukocytes of 16 cattle
 - Eight *Mycobacterium bovis*-infected
 - Eight unaffected controls
- The novel sequences contain 36 putatively novel genes that are differentially expressed
 - The top down-regulated novel gene encodes leukocyte immunoglobulin-like receptor A5 (LILRA5) which is missing in ARS-UCD1.2



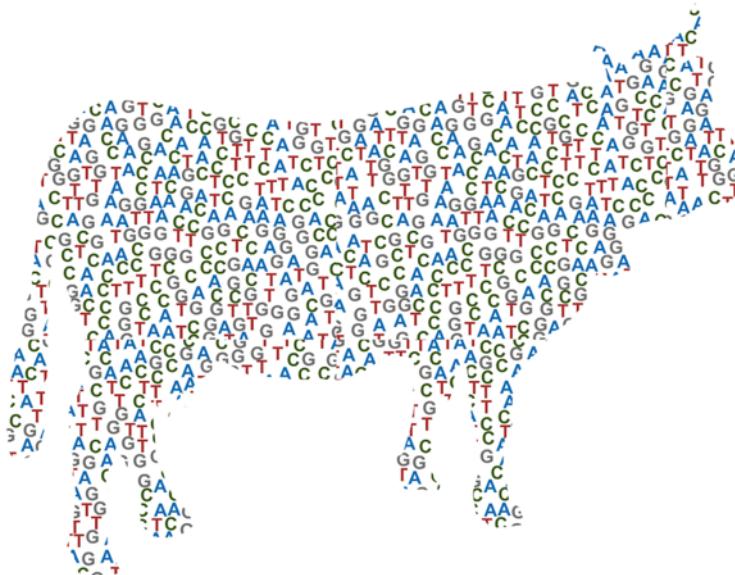
Conclusions

- Linear reference genomes are incomplete and susceptible to reference allele bias
- Haplotype-resolved assemblies are required to accurately detect alleles at structural variations
- Novel reference structures such as multi-assembly graphs make «non-reference sequences» amenable to genetic investigations

Contributors, collaborators & funding

- Danang Crysanto, Alexander Leonard, Zih-Hua Fang, [ETH Animal Genomics](#)
- Ben Hayes, Hans Daetwyler, Michael Goddard, Brisbane / Melbourne, Australia, for the [*1000 Bull Genomes Consortium*](#)
- Ben Rosen, Derek Bickhart, Tim Smith, USDA, USA, for the [*Bovine Pangenome Consortium*](#)

Thank you for your attention



information variance average genome sperm
findings mutations genes chromosome validation
sequencing highly identified region genotyping
calves traits homozygous sequenced results
accuracy traits association male bulls performed approach
variation effect dna mutation milk low frequency
genomic study trait revealed analysis
significantly number mutation male bulls results
control spps variant reference sequence imputed affected
homozygosity significant number variants population studies
breeding identification growth artificial trait male effects
genomic located spps variants cattle SNP protein phenotypes
however bovine haplotype QTL bta breeds also obtained illumina resulting
candidate causal allele imputation alleles genetic figure high segment
regions minimac fleckvieh data animal
chromosomes might populations bull gene genotyped
samples respectively genomewide recessive breed insemination
higher different phenotypic hypoplasia

ETH Zürich
Hubert Pausch
Animal Genomics
Eschikon 27 | EHB E 21
CH-8315 Lindau

hubert.pausch@usys.ethz.ch
www.ag.ethz.ch