

**Next Generation Sequencing –
The Role of New Sequence Technologies in Shaping the
Future of Veterinary Science**

Hosted by the RCVS Charitable Trust





Sequencing animal genomes

Alan Archibald

The Roslin Institute and R(D)SVS

University of Edinburgh

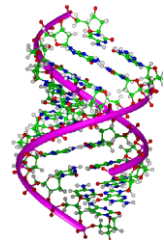


THE UNIVERSITY of EDINBURGH





***A sequenced genome is a requirement
for 21st Century biological research***



1953
Watson and
Crick



1977
DNA
sequenced
 Φ X174
5,386 nt

1990
Human
Genome
Project
launched



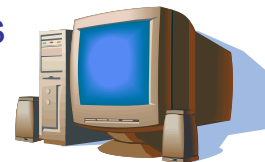
1991
PiGMAP
project
starts
'Halothane'
gene test

2001
Draft human
genome sequence



1920s and 30s
Fisher, Lush
and others
Population
Genetics

1970s +
Advances in
quantitative
analysis



1990s +
Quantitative
trait locus
(QTL)
mapping

2001
Genomic selection
proposed

2002
Mouse
draft
genome
sequence



2003
Human genome
sequence
"finished"
\$3 billion

2004
Chicken
genome
sequenced



2005
Dog
genome
sequenced

2007
Cat
genome
sequenced

2008
Human
1000
Genomes
Project
launched

2009
Cattle genome
sequenced

Horse genome
Sequenced

Mouse genome
"finished"



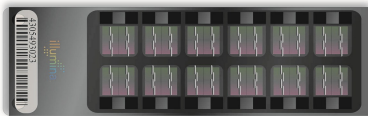
2012
Pig genome
sequenced
\$20 million



2008
Bovine 50K
SNP chip

2009
Pig 60K SNP
chip ~\$150

Sheep 60K
SNP chip



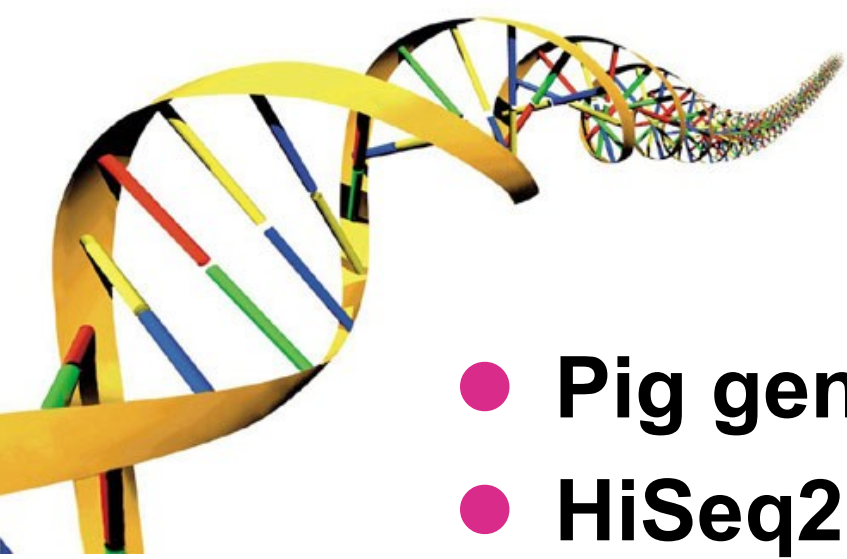
2010
750K bovine
SNP chips



2012
Chicken 600K
SNP chip

Goat 60K
SNP chip





DNA sequencing

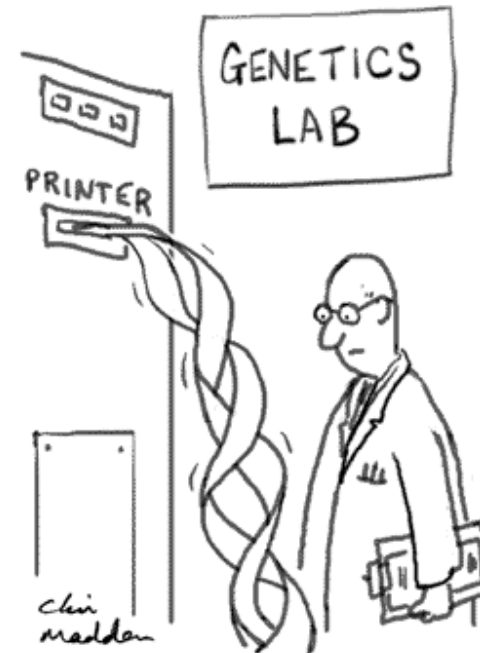
- **Pig genome = ~2.7 Gb**
- **HiSeq2000 generates per run:**
 - **600 Gb raw sequence (100 bp PE)**
 - **i.e. ~ 200x pig genome**
 - **in ~ 3×10^9 paired pieces**
 - **70% expected to pass chastity filter**
 - **Q score > 20 for full 100 bp**
- **But**

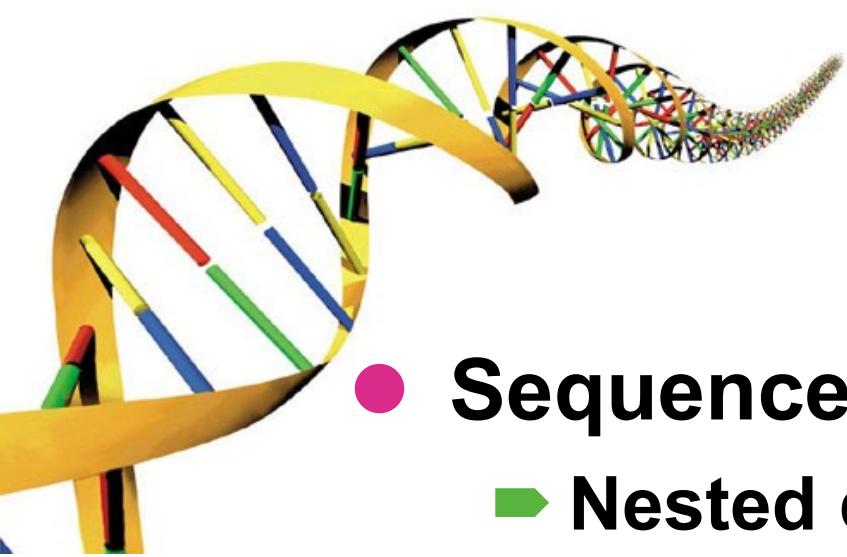




DNA sequencing

- **Generating DNA sequence data is:**
 - Easier, faster, cheaper
- **All DNA sequencing technologies:**
 - Generate errors
- **Pig nucleus**
 - Contains two haploid genomes
 - ~ 1 in 250 nt polymorphic
 - SNPs, indels, CNV





DNA sequencing

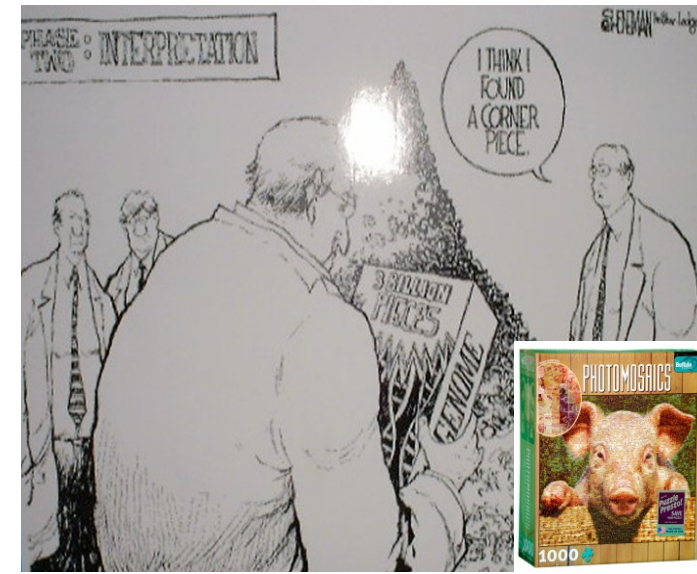
- **Sequence multiple cloned molecules**
 - **Nested deletions**
 - **2-300 bp reads**
 - **5 kbp gene = PhD**
- **Shotgun approaches**
 - **'random' fragmentation**
 - **high read depth**





DNA sequencing

- **Sequence assembly is not:**
 - Like solving a jigsaw puzzle
 - 40% repetitive DNA sequence
 - cf. cloudless blue sky
- **Building a high quality reference genome sequence takes:**
 - Time, money, effort



Pig Genome Sequencing Project

Swine Genome Sequencing Consortium

Comparative and Functional Genomics

Comp Funct Genom 2005; **6**: 251–255.

Published online in Wiley InterScience (www.interscience.wiley.com). DOI: 10.1002/cfg.479

Conference Review

Swine Genome Sequencing Consortium (SGSC): a strategic roadmap for sequencing the pig genome

Lawrence B. Schook^{1,2*}, Jonathan E. Beever^{1,2}, Jane Rogers³, Sean Humphray³, Alan Archibald⁴, Patrick Chardon⁵, Denis Milan⁶, Gary Rohrer⁷ and Kellye Eversole⁸

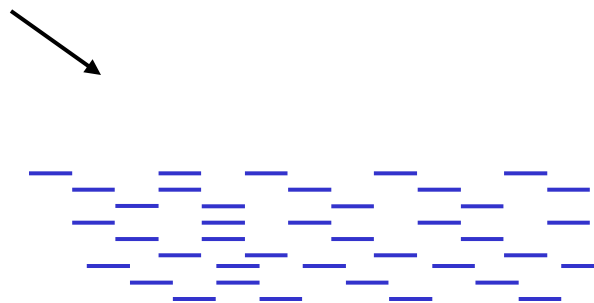
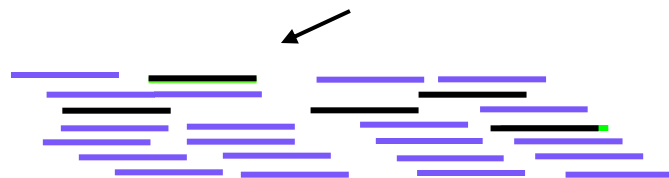




Hybrid Shotgun Sequencing Strategy

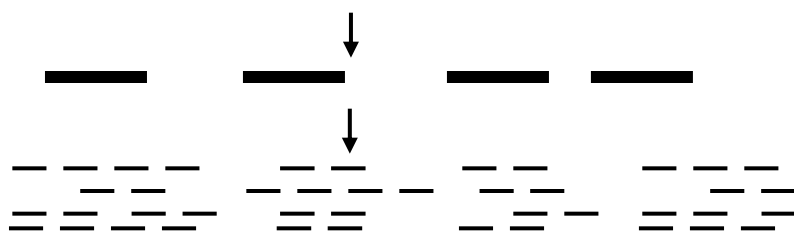


Minimal set of overlapping BACs selected from physical map

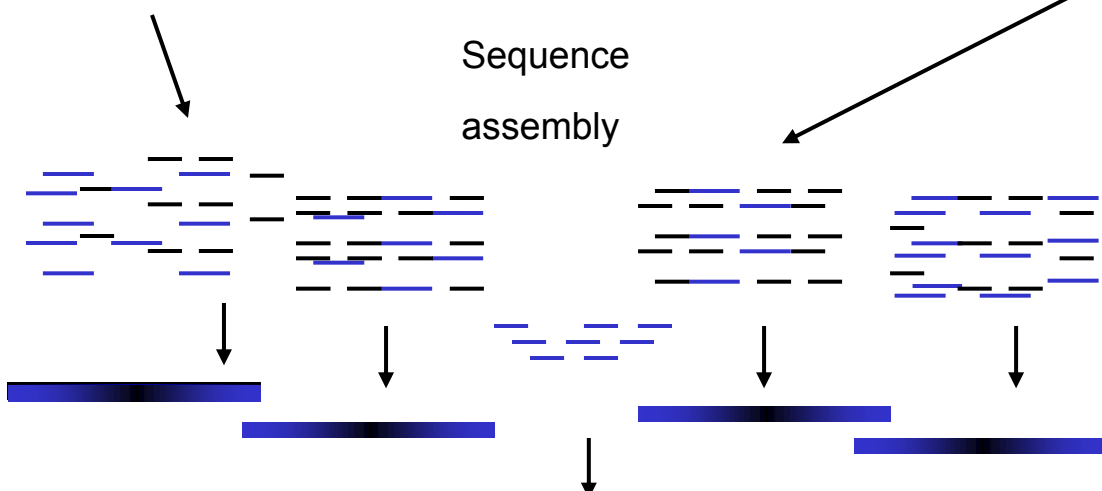


Whole-genome shotgun reads

BAC shotgun reads



Sequence assembly

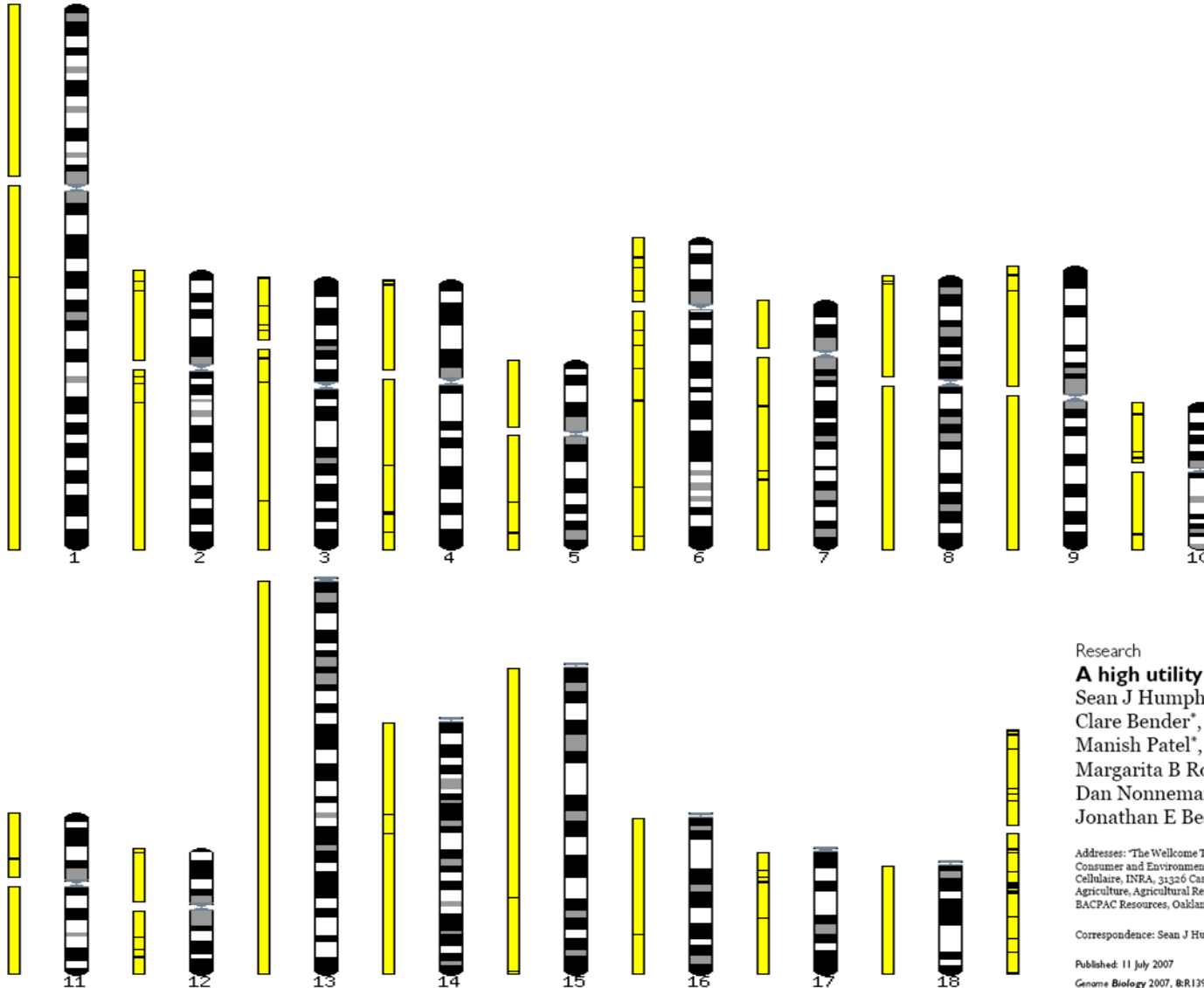


Combine overlapping whole-genome and BAC-derived reads

Assemble clone sequences to represent chromosomes and annotate using Ensembl automated pipeline



The best clone-based physical map of a mammal



- 172 placed contigs
- average length 15 Mb
- covering 2.58Gb
- >98% of euchromatin

Humphray *et al.*, 2007.
Genome Biology

Open Access

Research

A high utility integrated map of the pig genome

Sean J Humphray^{*}, Carol E Scott^{*}, Richard Clark^{*}, Brandy Marron[†], Clare Bender^{*}, Nick Camm^{*}, Jayne Davis^{*}, Andrew Jenks^{*}, Angela Noon^{*}, Manish Patel^{*}, Harminder Sehra^{*}, Fengtang Yang^{*}, Margarita B Rogatcheva[‡], Denis Milan[‡], Patrick Chardon[§], Gary Rohrer[¶], Dan Nonneman[¶], Pieter de Jong[¶], Stacey N Meyers[‡], Alan Archibald[#], Jonathan E Beever[†], Lawrence B Schook[†] and Jane Rogers^{*}

Addresses: ^{*}The Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA UK. [†]College of Agriculture, Consumer and Environmental Sciences, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801 USA. [‡]Laboratoire de Génétique Cellulaire, INRA, 31326 Castanet-Tolosan, France. [§]INRA-CEA, Domaine de Vilvert, 78352, Jouy en Josas cedex, France. [¶]US Department of Agriculture, Agricultural Research Service, US Meat Animal Research Center, Clay Center, NE 68933-0166, USA. [‡]Children's Hospital Oakland Research Institute, Oakland, California 94609, USA. [#]Roslin Institute, Roslin, Midlothian EH25 9PS, UK.

Correspondence: Sean J Humphray. Email: sjh@sanger.ac.uk

Published: 11 July 2007

Genome Biology 2007, 8:R139 (doi:10.1186/gb-2007-8-7-r139)

Received: 12 March 2007

Revised: 21 June 2007

Accepted: 11 July 2007

BAC Contigs / Fragments

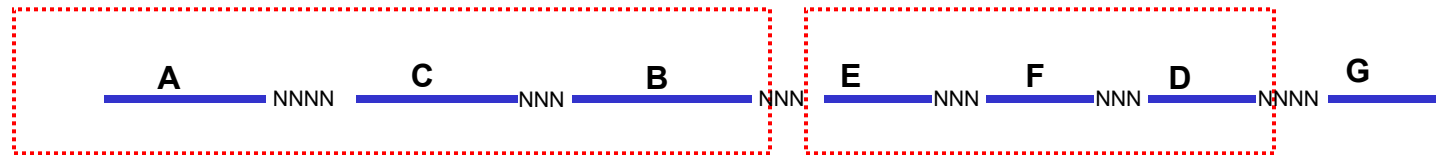
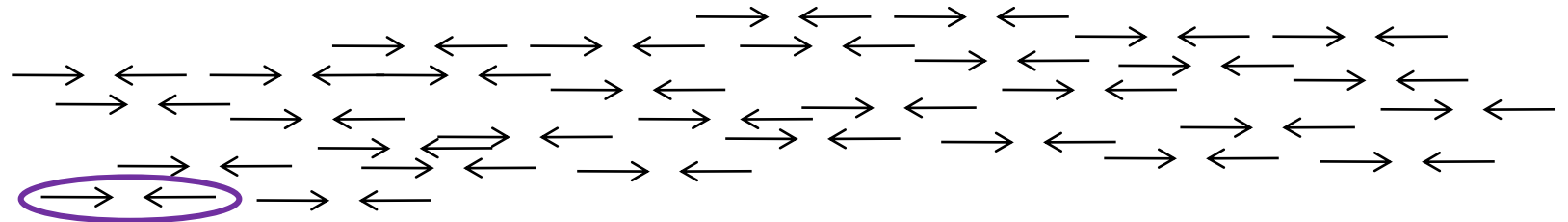
end sequences
of subclone libraries

768 subclones / BAC
~3-4x coverage

phrap

create fragment chains

Submission to EMBL/Genbank



fragment chain 1

fragment chain 2

BAC Contigs / Fragments

end sequences
of subclone libraries

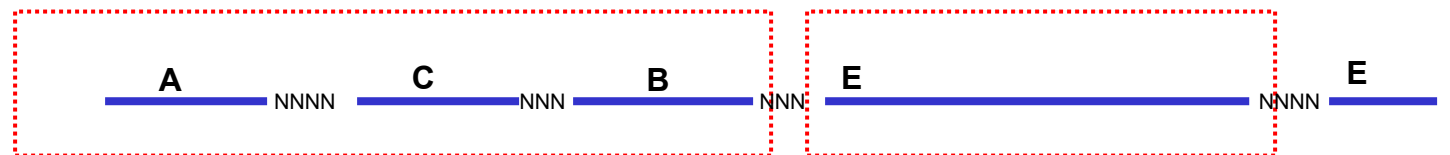
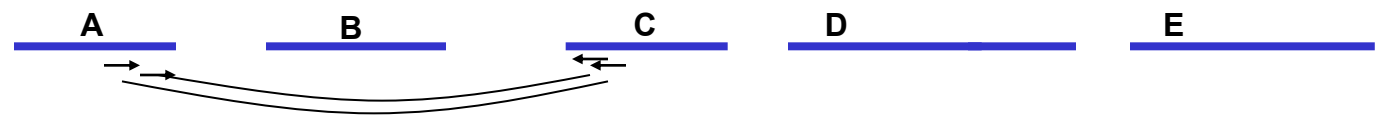
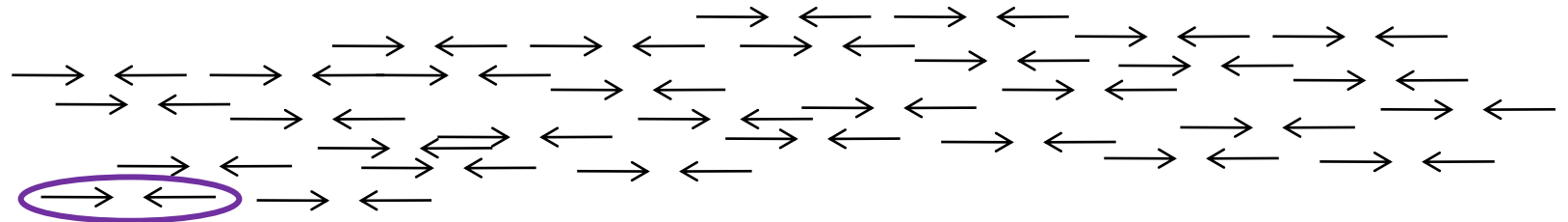
768 subclones / BAC
~3-4x coverage

phrap

1 round automated primer
design and walk

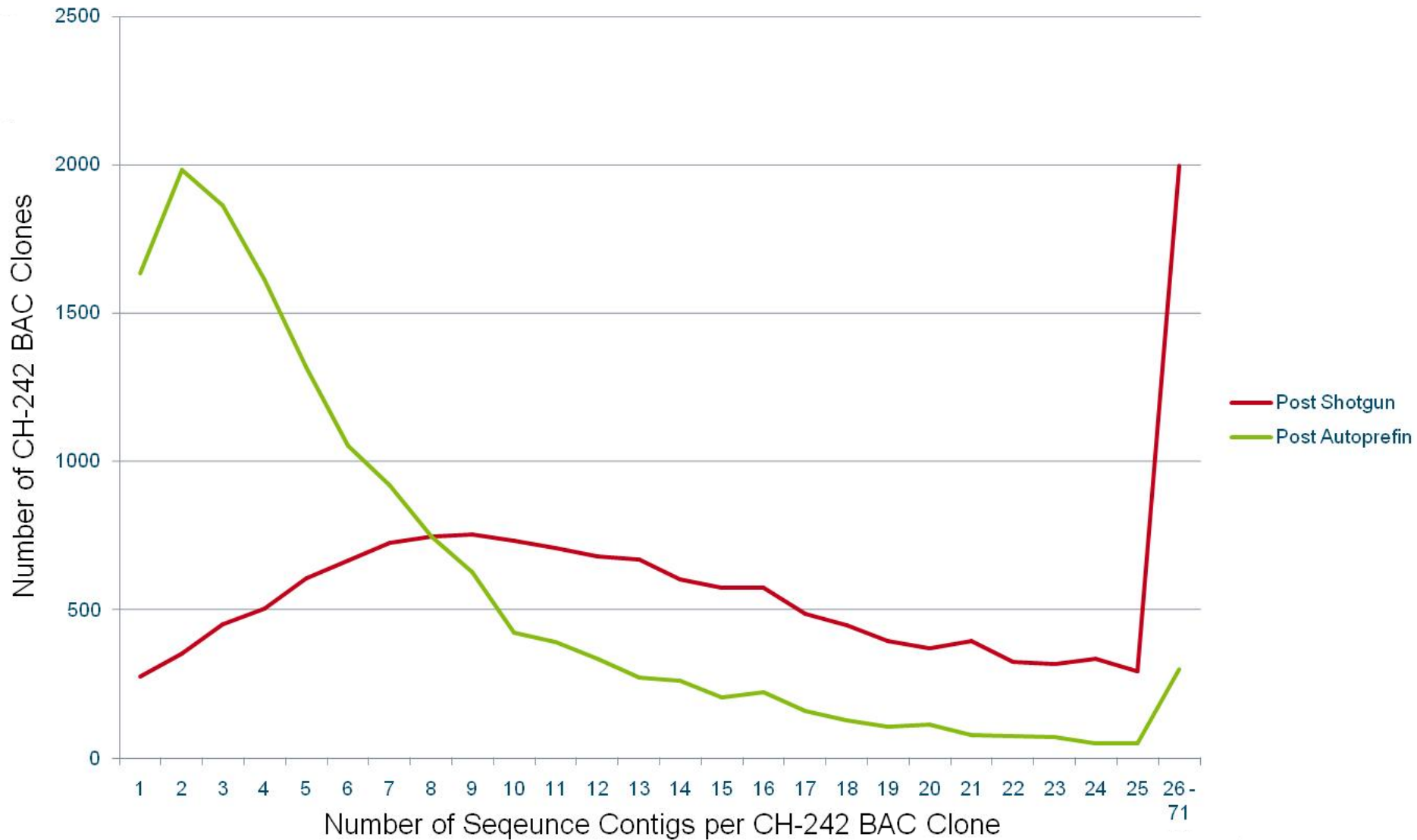
create fragment chains

Submission to EMBL/Genbank



fragment chain 1

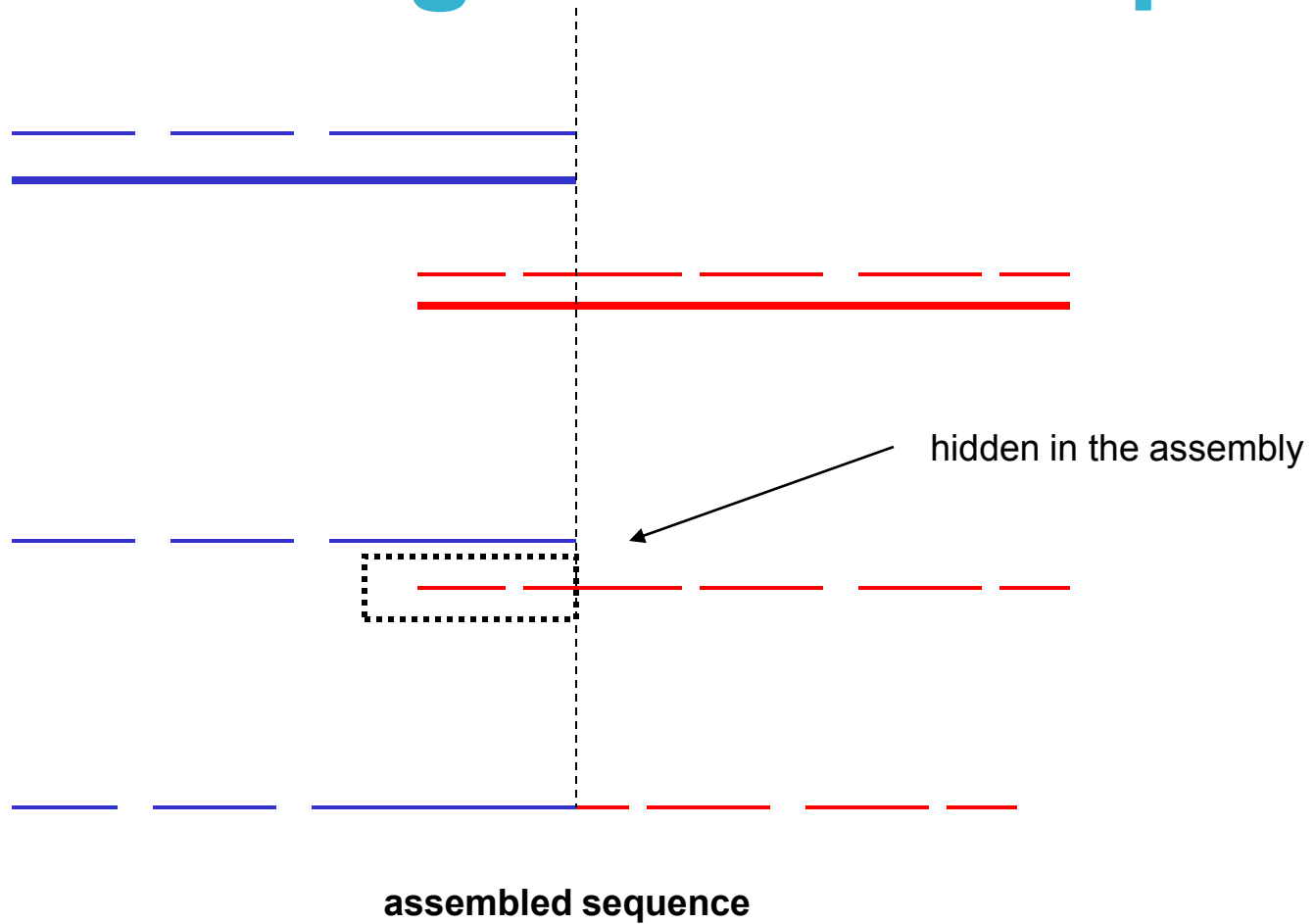
fragment chain 2



Sequence improvement

	Sequence coverage	Sequence accuracy	Number of sequence contigs	Average contig length bp	Genes with incorrect structures
<i>4x draft sequence</i>	<i>97% (~80Mb missing)</i>	<i>99.99% accurate (1 error in 10kbp)</i>	<i>160,000 (7k / chr)</i>	<i>17,000</i>	<i>30%</i>
<i>Improved draft sequence</i>	<i>99% (~32Mb missing)</i>	<i>99.99% accurate (1 error in 10kbp)</i>	<i>65,000 (3k / chr)</i>	<i>42,000</i>	<i>5%</i>
<i>Gold standard finished sequence</i>	<i>99.9% (~5Mb missing)</i>	<i>99.999% (1 error in 100kbp)</i>	<i><200</i>	<i>14,000,000</i>	<i>0%</i>

Resolving BAC overlaps



Add WGS data

Same Duroc individual as CHORI-242

BGI

- **66.5 Gb of sequence (24-fold)**
- **Read length: 44**

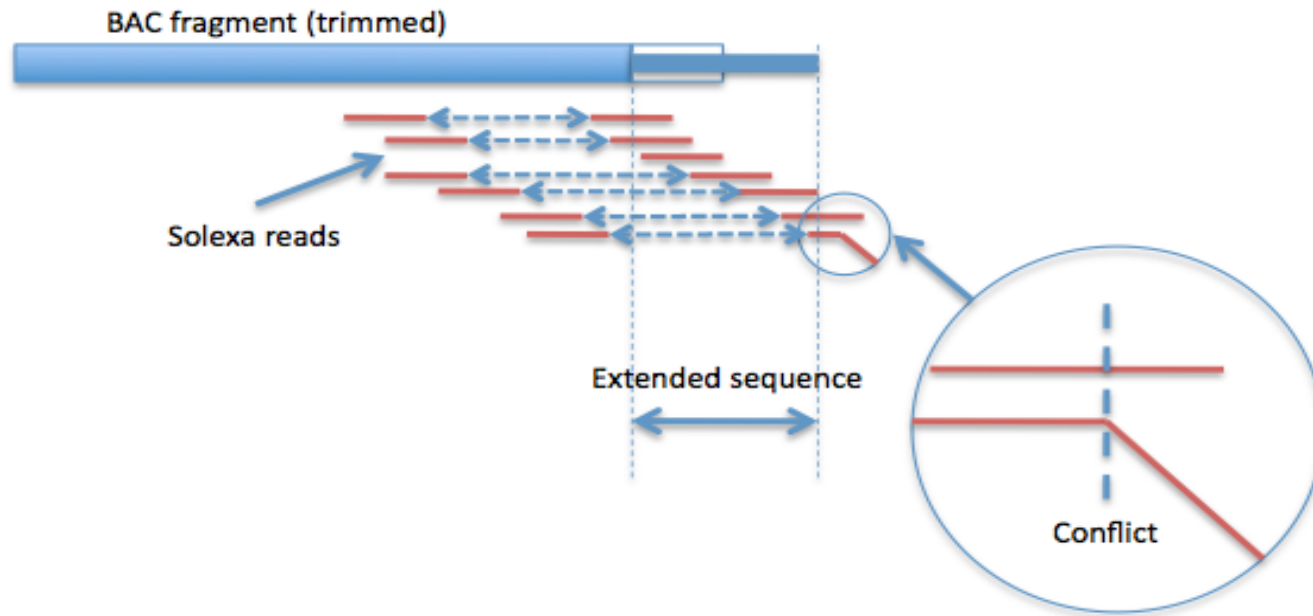
WTSI

- **~40 Gb of sequence (14-fold)**
- **Read length: 108**

Assemblies

- **SOAPdeNovo (Shenting Li – Pig Genome III - Hinxton)**
 - **18,409 gaps closed (using reads)**
 - **31,359 gaps closed (using scaffolds)**
 - **252Mb in new scaffolds (BRCA1 gene)**
- **Cortex (Caccamo & Iqbal)**
 - **~3000 gaps closed – use to correct issues within BAC assemblies**

Illumina Assembly – using reads



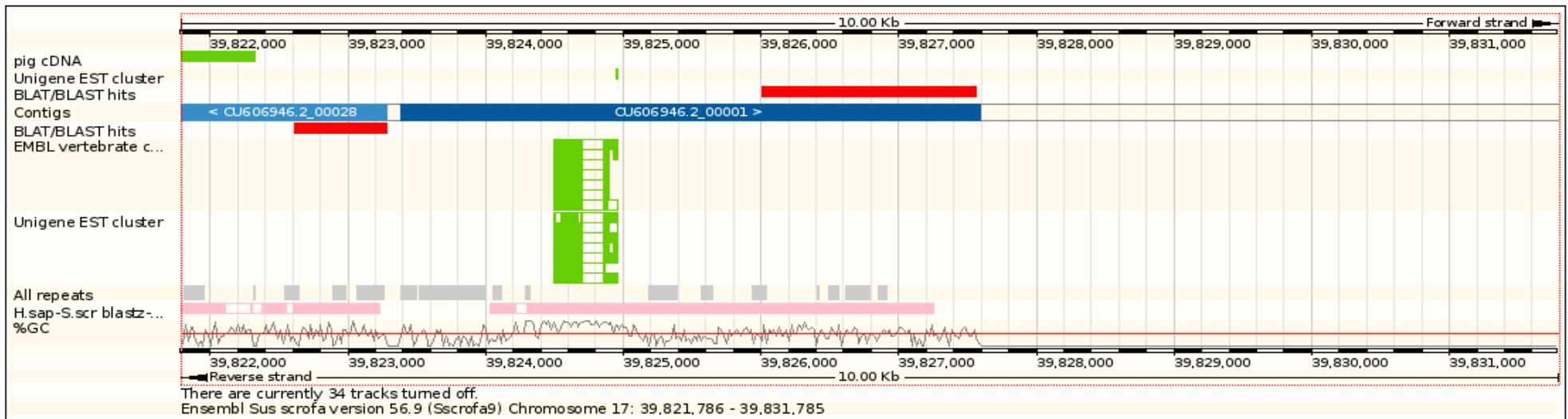
Illumina Assembly – using contigs

Spanners



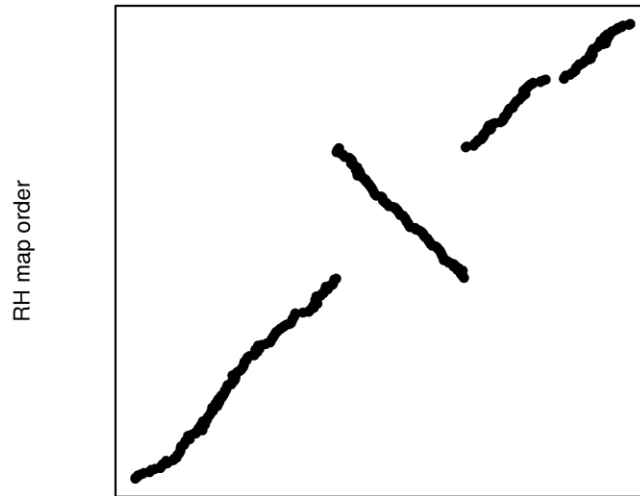
Hangers



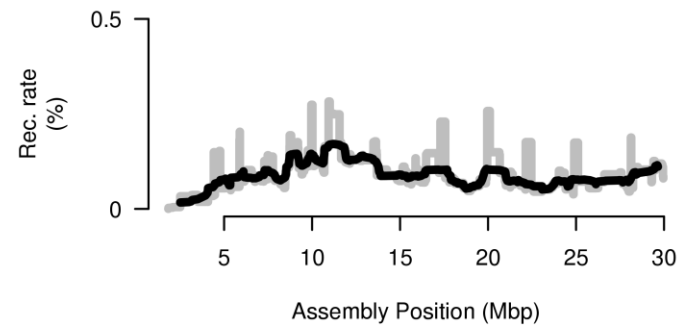
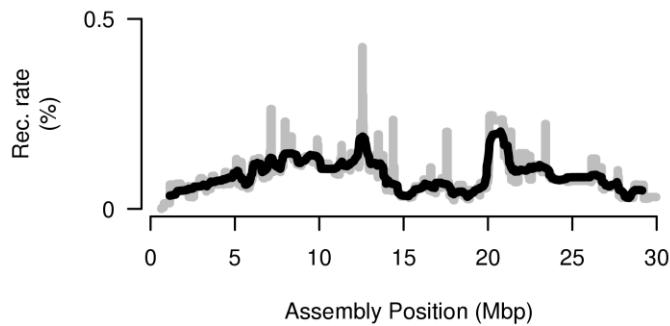
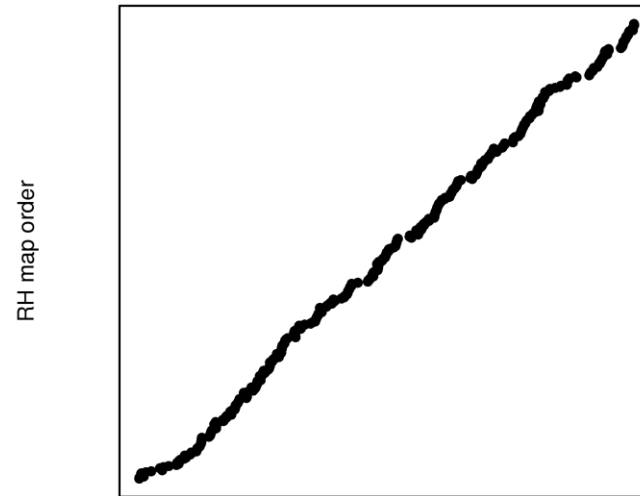


Order-orientation Sequence vs RH maps

a - Build 9 assembly



b - Current assembly



Annotated assembly in Ensembl

The screenshot shows the Ensembl website interface for the Pig (Sus scrofa) genome assembly. The top navigation bar includes links for BLAST/BLAT, BioMart, Tools, Downloads, Help & Documentation, Blog, and Mirrors. The main content area is titled "Pig (Sus scrofa)" and features a search bar with a "Go" button. Below the search bar, there are sections for "Description", "Assembly", "Previous assemblies", and "Annotation". The "Assembly" section includes a photograph of a pig and a link to "Download Pig genome sequence (FASTA)". The "Previous assemblies" section shows a dropdown menu for "Sscrofa9 (Release 66, Feb2012)" and a "Go to archive" button. The "Annotation" section provides information about the annotation process and a link to "Detailed information on genebuild (PDF)". The footer of the page includes the Ensembl release date (67 - May 2012) and copyright information (WTSI / EBI), along with links for "About Ensembl", "Privacy Policy", and "Contact Us".

e!Ensembl BLAST/BLAT | BioMart | Tools | Downloads | Help & Documentation | Blog | Mirrors Login · Register



Pig (Sscrofa10.2)

e.g. ENSSSCG00000004244 or 7:60107914-60305245 or apoptosis

Description

Assembly

The Sscrofa10.2 assembly of the pig genome was produced in August 2011 by the Swine Genome Sequencing Consortium (SGSC). It consists of 20 chromosomes (1-18, X and Y) and 4562 unplaced scaffolds. This genome assembly has GCA_000003025.4 as its GenBank assembly accession. The genome assembly represented here corresponds to GenBank Assembly ID [GCA_000003025.4](#)



[Download Pig genome sequence \(FASTA\)](#)

Previous assemblies

Sscrofa9 (Release 66, Feb2012)

Annotation

Sscrofa10.2 was annotated using a standard Ensembl mammalian genebuild pipeline, incorporating RNA-Seq data provided by the (SGSC). The annotation process is described in this document. When the annotation was completed, the gene set contained the following: 21,640 protein coding genes, 380 pseudogenes and 2,965 ncRNAs.

- [Detailed information on genebuild \(PDF\)](#)

Ensembl release 67 - May 2012 © WTSI / EBI [About Ensembl](#) | [Privacy Policy](#) | [Contact Us](#)

[Permanent link - View in archive site](#)

Done Internet 100%

http://www.ensembl.org/Sus_scrofa/Info/Index

Sscrofa10.2 – assembly, genes

Assembly	Annotation*			
	Placed	Unplaced		
Total length	2,596,639,456	211,869,922	Protein-coding genes:	21,627
Ungapped length	2,323,671,356	195,490,322	Pseudogenes:	380
Scaffolds	5,343	4,562	ncRNA genes**:	2,965
Contigs	73,524	168,358	Gene exons:	197,675
Scaffold N50	637,332	98,022	Gene transcripts:	26,487
Contig N50	80,720	2,423		

Data sharing and publication

- **Bermuda agreement**
- **Fort Lauderdale agreement**

nature

Vol 461|10 September 2009

OPINION

Prepublication data sharing

Rapid release of prepublication data has served the field of genomics well. Attendees at a workshop in Toronto recommend extending the practice to other biological data sets.

Pig genomes provide insight into porcine demography, domestication and evolution

Martien A.M. Groenen^{1*}, Alan L. Archibald^{2*}, Hirohide Uenishi³, Christopher K. Tuggle⁴, Yasu Takeuchi⁵, Max F. Rothschild⁴, Claire Rogel-Gaillard⁶, Chankyu Park⁷, Denis Milan⁸, Hendrik-Jan Megens¹, Shengting Li⁹, Denis Larkin¹⁰, Heebal Kim¹¹, Laurent A. F. Frantz¹, Mario Caccamo¹², Hyeonju Ahn¹¹, Bronwen L. Aken¹³, Anna Anselmo¹⁴, Christian Anthon¹⁵, Loretta Auvil¹⁶, Bouabid Badaoui¹⁴, Craig W. Beattie¹⁷, Christian Bendixen¹⁸, Daniel Berman¹⁹, Frank Blecha²⁰, Jonas Blomberg²¹, Lars Bolund⁹, Mirte Bosse¹, Sara Botti¹⁴, Zhan Bujie¹⁸, Megan Bystrom⁴, Boris Capitanu¹⁶, Denise Carvalho-Silva²², Patrick Chardon⁶, Celine Chen²⁴, Ryan Cheng⁴, Sang-Haeng Choi²⁵, William Chow¹³, Richard C. Clark¹³, Christopher Clee¹³, Richard P.M.A. Crooijmans¹, Harry D. Dawson²⁴, Patrice Dehais⁸, Floravante De Sapio², Bert Dibbits¹, Nizar Drou¹², Zhi-Qiang Du⁴, Kellye Eversole²⁶, João Fadista¹⁸, Susan Fairley¹³, Thomas Faraut⁸, Geoffrey J. Faulkner², Katie E. Fowler²⁷, Merete Fredholm¹⁵, Eric Fritz⁴, James G.R. Gilbert¹³, Elisabetta Giuffra¹⁴, Jan Gorodkin¹⁵, Darren K. Griffin²⁷, Jennifer L. Harrow¹³, Alexander Hayward²⁸, Kerstin Howe¹³, Zhi-Liang Hu⁴, Sean J. Humphray¹³, Toby Hunt¹³, Henrik H. Jensen¹⁸, Patric Jern²⁸, Matthew Jones¹³, Jerzy Jurka²⁹, Hiroyuki Kanamori³⁰, Ronan Kapetanovic², Jaebum Kim^{31,23}, Jae-Hwan Kim³², Kyu-Won Kim³³, Tae-Hun Kim³⁴, Greger Larson³⁵, Kyooyeol Lee⁷, Kyung-Tai Lee³⁴, Richard Leggett¹², Harris A. Lewin³⁶, Yingrui Li⁹, Wansheng Liu³⁷, Jane E. Loveland¹³, Yao Lu⁹, Joan K. Lunney¹⁹, Jian Ma³⁸, Ole Madsen¹, Katherine Mann¹⁹, Lucy Matthews¹³, Stuart McLaren¹³, Takeya Morozumi³⁰, Michael Murtaugh³⁹, Jitendra Narayan¹⁰, Dinh Truong Nguyen⁷, Peixiang Ni⁹, Song-Jung Oh⁴⁰, Suneel Onteru⁴, Frank Panitz¹⁸, Eung-Woo Park³⁴, Hong-Seog Park²⁵, Geraldine Pascal⁴¹, Yogesh Paudel¹, Miguel Perez-Enciso⁴², Ricardo Ramirez-Gonzalez¹², James M. Reecy⁴, Sandra Rodriguez-Zas⁴³, Gary A. Rohrer⁴⁴, Lauletta Rund⁴³, Yongming Sang²⁰, Kyle Schachtschneider⁴³, Joshua Schraiber⁴⁵, John Schwartz³⁹, Linda Scobie⁴⁶, Carol Scott¹³, Stephen Searle¹³, Bertrand Servin⁸, Bruce R. Southey⁴³, Goran Sperber⁴⁷, Peter Stadler⁴⁸, Jonathan Sweedler⁴⁹, Hakim Tafer⁴⁸, Bo Thomsen¹⁸, Rashmi Wali⁴⁶, Jian Wang⁹, **Jun Wang⁹, Simon White¹³, Xun Xu⁹, Martine Yerle⁸, Jianguo Zhang⁹, Guojie Zhang⁹, Jie Zhang⁵⁰, Shuhong Zhao⁵⁰, **Jane Rogers¹²**, **Carol Churcher¹³** and **Lawrence B. Schook⁵¹**.**

Pig Genome Sequence Project: A Blueprint for Agriculture, Life and Biomedical Sciences

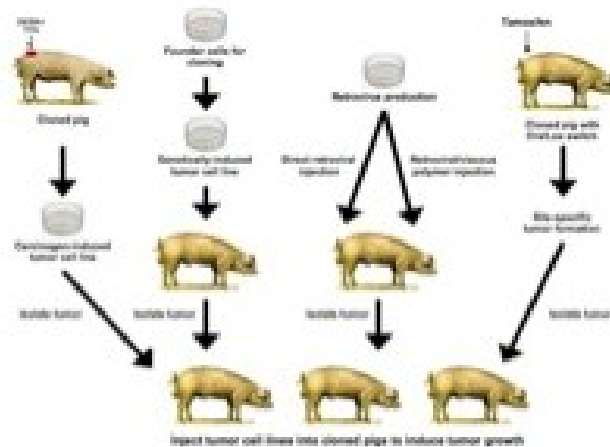
CGCTCATATCGATCGATTGCATCGATCCGATCGATCGATGCTAGCTAGCGAGTCGAA
 GCCTCTATCTCATCCCTTAACCTAGCCATATACCCCTAGCCCTAGCTCCATCTCAGATCGA
 TC
 TT
 CC

DNA-Based Models

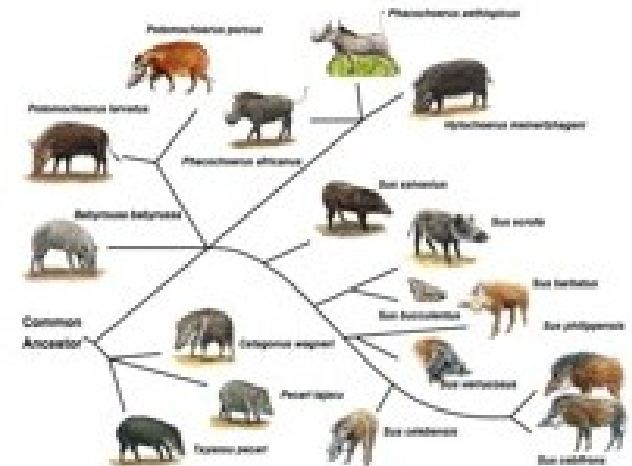
GGTCTATCTATCTCTTAACGTTAGCTATATCGCTAGGCTAGTCTTGATCTGAGATCGA
 TCAGGGCCATTATCGGCATCGATCGATTGATCAACCTGTAGAGCCATTGCTATTCGA



AGRICULTURAL SCIENCES
Genotypes and Phenotypes



BIOMEDICAL SCIENCES
Designing Cancer Models



LIFE SCIENCES
Mechanisms of Mammalian
Evolution and Diversity

Acknowledgements - funding

Funding Source	Funding Level	Activity Supported
CSREES USDA	\$10,000,000	Clone shotgun sequencing
USDA – ARS	\$1,000,000	Clone shotgun sequencing
N Carolina Agric Res Service	\$50,000	Sequencing
Iowa State University	\$200,000	Targeted sequencing
USA National Pork Board	\$750,000	WGS sequencing
Iowa Pork Producers Assn.	\$100,000	Targeted Sequencing
N. Carolina Pork Council	\$100,000	Sequencing
EU SABRE Project	€1,600,000	6X of SSC7 & SSC14
Danish government/BGI	€600,000	Next-gen seq WGS
INRA Genoscope, France	€1,000,000	SNP discovery (1 million reads)
BBSRC, UK	£2,000,000	Annotation, analysis; SSCX/Y
Wellcome Trust Sanger Institute	£1,600,000	Contribution to read costs
Wellcome Trust Sanger Institute	£300,000	Finish ENCODE + MHC
Dutch IPG	€600,000	6X of SSC4
Japan	£10,000	35 SSC7 clones
Korea	0.5 million reads	WGS sequencing



Genomics-enabled tools

Reference genome sequence as a key resource and framework for biological research

■ Genetics

- Variation (SNPs, indels, CNVs)
 - SNP chips, Genotype-by-Sequence
- Genome-Wide Association Studies (GWAS)
- Genetic improvement

■ Functional genomics

- incl. physiology, immunology,
- Genome-wide analysis of responses to perturbation
 - Gene expression, methylation,
 - Microarrays, Assay-by-sequence

HIGH DENSITY SNP GENOTYPING CHIP

OPEN ACCESS Freely available online

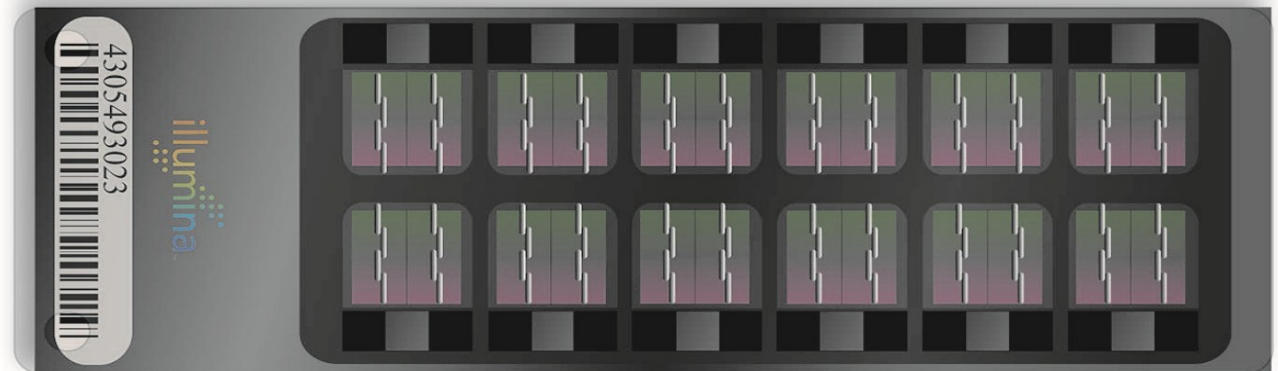
August 2009

 PLoS one

Design of a High Density SNP Genotyping Assay in the Pig Using SNPs Identified and Characterized by Next Generation Sequencing Technology

Antonio M. Ramos¹, Richard P. M. A. Crooijmans¹, Nabeel A. Affara², Andreia J. Amaral¹, Alan L. Archibald³, Jonathan E. Beever⁴, Christian Bendixen⁵, Carol Churcher⁶, Richard Clark⁶, Patrick Dehais⁷, Mark S. Hansen⁸, Jakob Hedegaard⁵, Zhi-Liang Hu⁹, Hindrik H. Kerstens¹, Andy S. Law³, Hendrik-Jan Megens¹, Denis Milan⁷, Danny J. Nonneman¹⁰, Gary A. Rohrer¹⁰, Max F. Rothschild⁹, Tim P. L. Smith¹⁰, Robert D. Schnabel¹¹, Curt P. Van Tassell¹², Jeremy F. Taylor¹¹, Ralph T. Wiedmann¹⁰, Lawrence B. Schook⁴, Martien A. M. Groenen^{1*}

¹ Wageningen University, Animal Breeding and Genor Kingdom, ³ Division of Genetics and Genomics, The R Kingdom, ⁴ Institute for Genomic Biology, University Denmark, ⁶ The Wellcome Trust Sanger Institute, Tl Cellulaire, Castanet Tolosan, France, ⁸ Illumina, Inc., S



- Human 1000 Genomes Project
 - ~4-6x coverage / individual
 - revealing genetic burden
 - ~1-200 potential Loss of Function mutations per person
- Human genetics studies
 - 10's of thousands per study
 - ICQG 2012
 - 30K sequenced genomes in a study

- Pooled samples
 - 10-15x coverage
 - Chickens, cattle, pigs,
 - SNP discovery
 - Signatures of selection
 - Signatures of domestication
- Individual genomes
 - 4-10x coverage
 - £1,000 per genome

- 1000 Bull Genomes Project
 - Collaborative, Cloud data repository
 - Nnn bulls, average coverage ~11x
 - Data analysis cycles for genomic prediction
- Pigs
 - Groenen (Wageningen) ~300 individual pigs
 - Korean ~60 individual pigs
 - China ?? pigs



Gene Expression Microarray

Gene Expression Atlas of the Pig

- Tool for monitoring gene expression
- Inferring function of unknowns
 - Inform genome annotation
- Comparative functional genomics
 - Is pig kidney more/less like human kidney than mouse kidney?
- Current arrays
 - Poorly annotated
 - Content elderly (e.g. Affymetrix 2004 design)

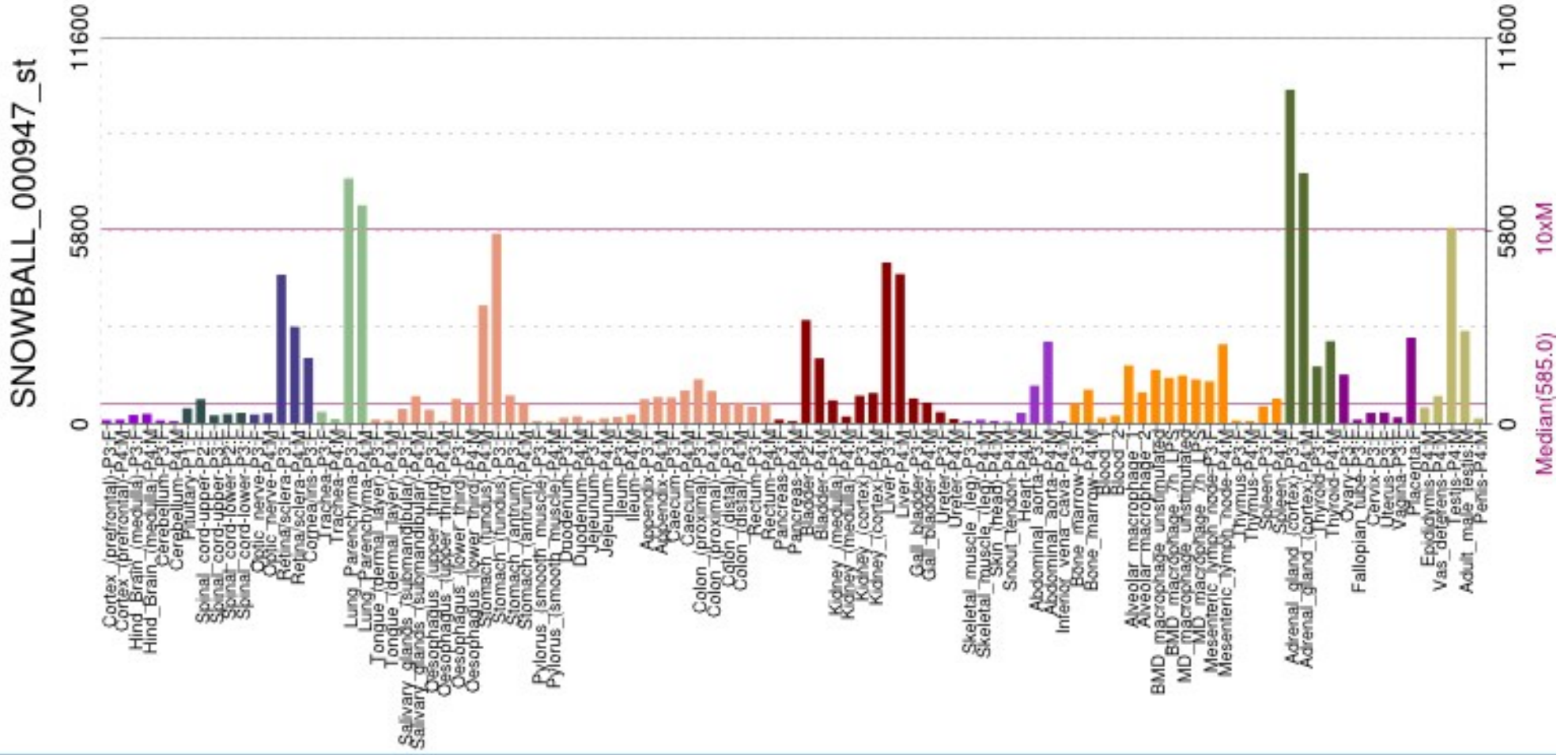
Affymetrix Porcine Snowball Array content



- 123 Affy controls
- 35 virus genomes (tiled 17 bp spacing)
- 1,857 miRNA probes
- 37 MT-mRNA
- 45,927 mRNA
 - 37,842 with annotation
 - 6,767 LOC annotations
 - 16,626 unique genes with official symbol/description

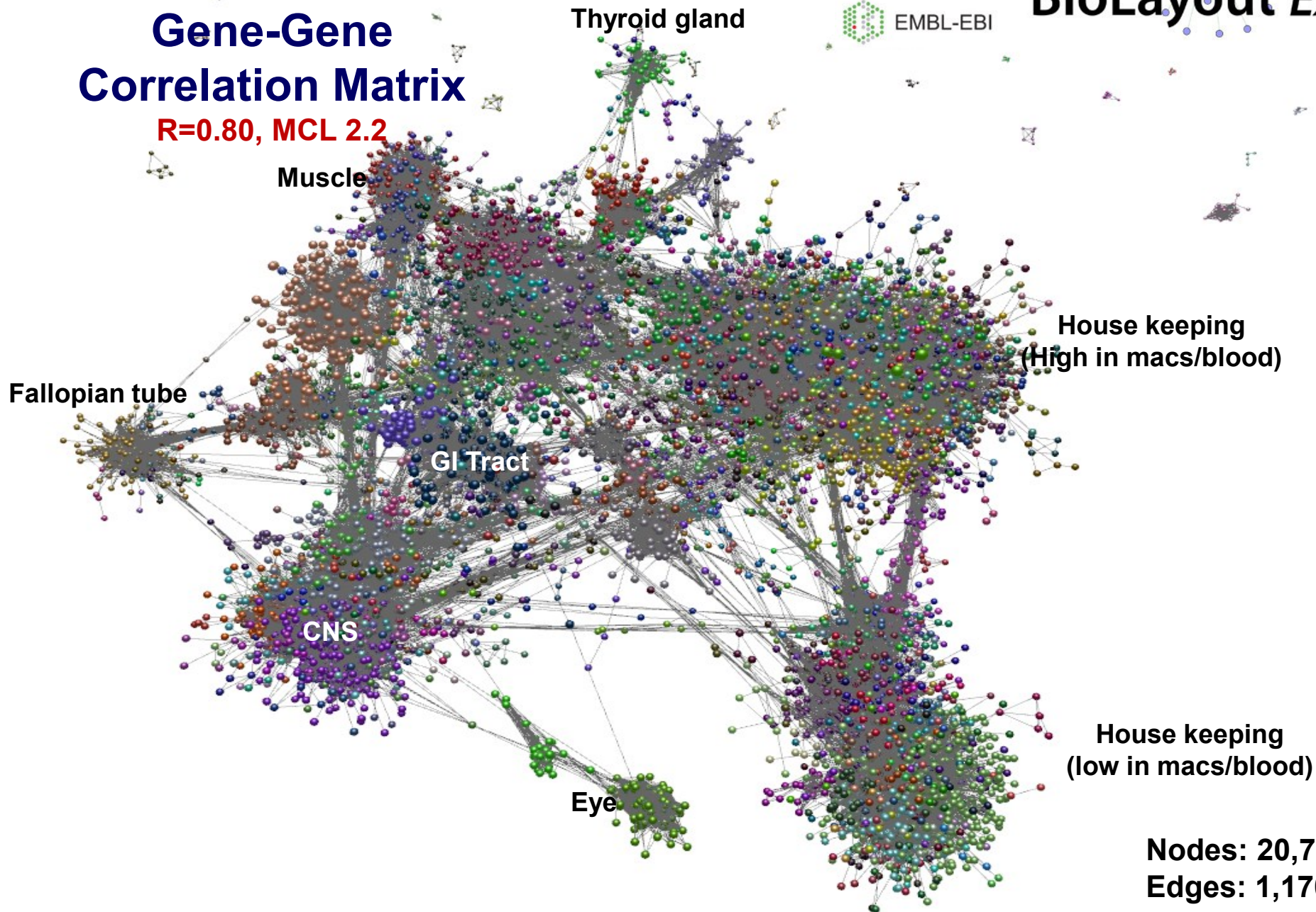


Expression profiles



Gene-Gene Correlation Matrix

R=0.80, MCL 2.2



Nodes: 20,703
Edges: 1,170,296

Those involved:

Tom Freeman

Array design/annotation

Fios Genomics

Chris Tuggle

Sanger Pig genome project

Affymetrix Inc.

Dario Beraldi

Tissue Atlas

David Hume

Alan Archibald

Mark Barnett

Kenny Bailie

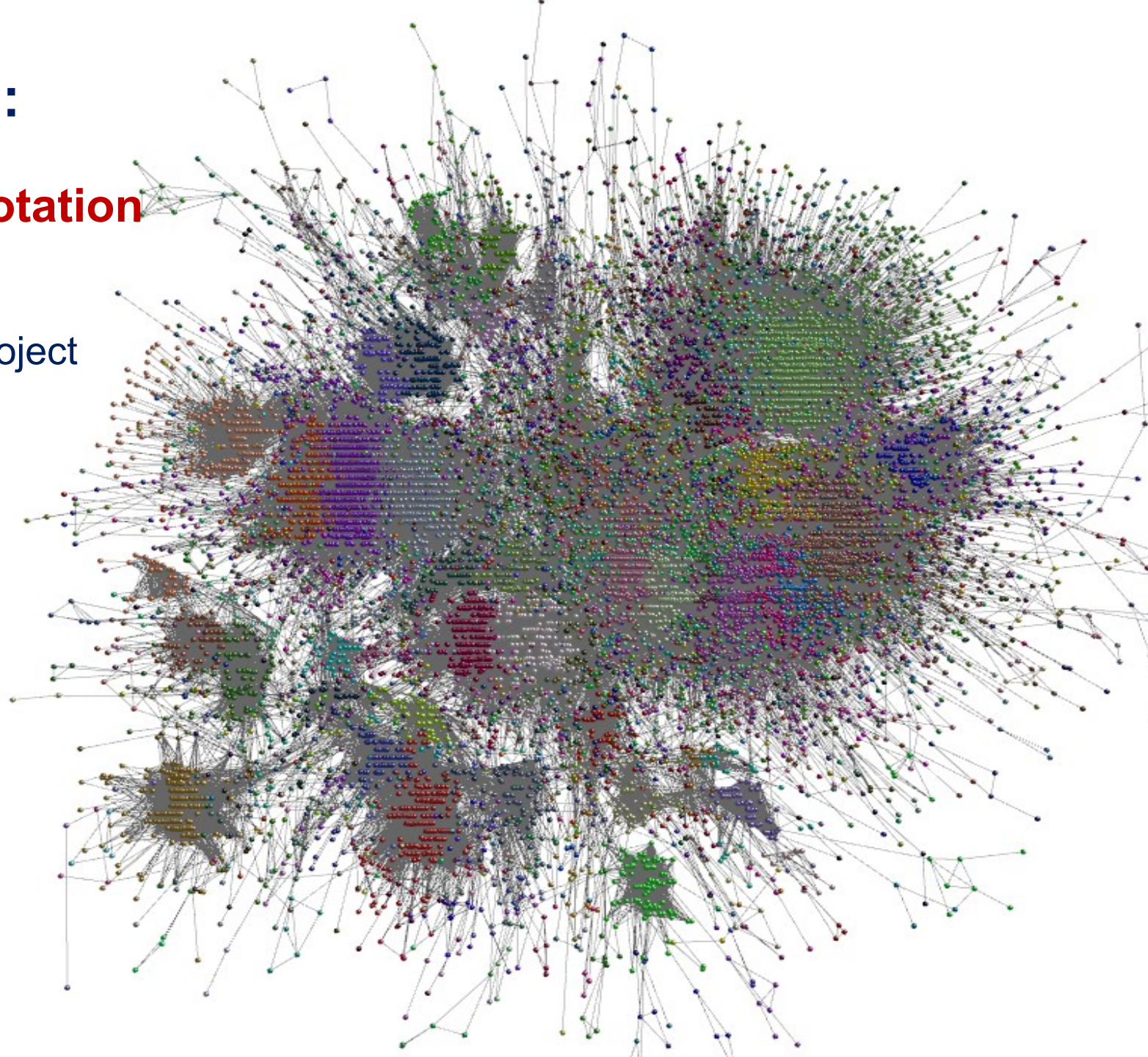
Ronan Kapetanovic

Kim Summers

Lynsey Fairbairn

Andru Tomoiu

ARK-Genomics



Roslin's contribution



- Leadership
 - PiGMaP, ChickMap,SABRE, Quantomics, 3SR,..
 - SGSC, ISGC,....
- Genome sequence data
 - De novo: sheep (70x)
 - Re-sequence: pig, cattle, chicken
- Analysis, annotation (Ensembl)
 - Chicken, cattle, turkey,..pig,....
- SNP chip development
 - Pig, cattle, chicken, salmon,
- Microarray (ARK-Genomics)
 - cDNA to Affymetrix Snowball



ARK-Genomics
Centre for Comparative & Functional Genomics



ARK-GENOMICS



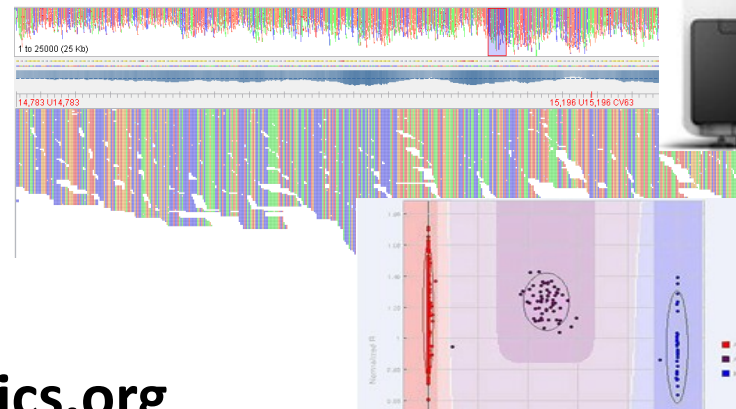


- **ARK-Genomics**

- High-throughput facility focusing on genetics and genomics
- BBSRC National Capability based at The Roslin Institute
- Collaborators on every continent

- Offering research, collaborations and service provision

- Investing in the latest genomics technologies
 - Sequencing
 - Genotyping
 - Transcriptomics
 - Comparative Genomics
 - Bioinformatics



<http://www.ark-genomics.org>



Technologies

DNA Sequencing

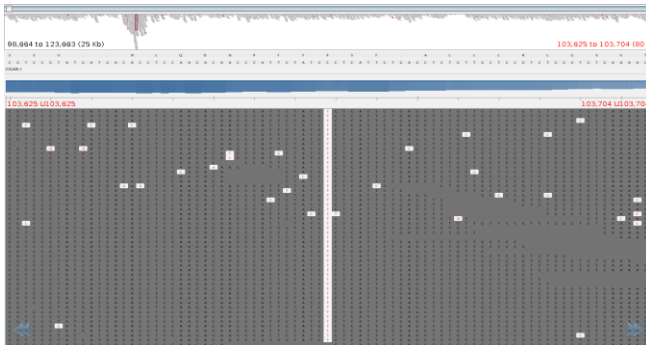
Illumina Sequencing

- Up to 150bp paired
- Novel genomes
- Resequencing
- RNA-Seq
- ChIP-Seq
- Epigenetics

Illumina HiSeq 2000

Illumina GAIIx

Sanger 3730



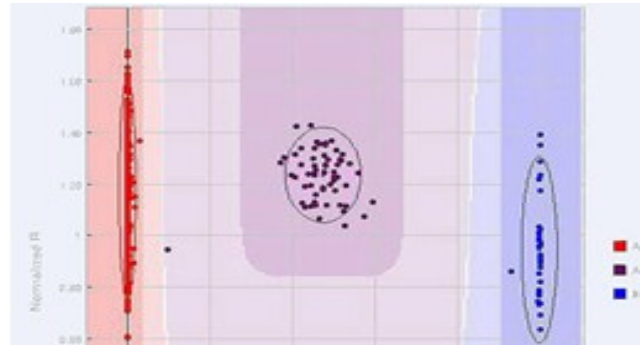
Genotyping

Illumina - from HD to custom chips

- iScan, Infinium
- BeadXpress, Goldengate
- BeadChip

Affymetrix

- GeneTitan, Axiom
- Process 96 arrays / run



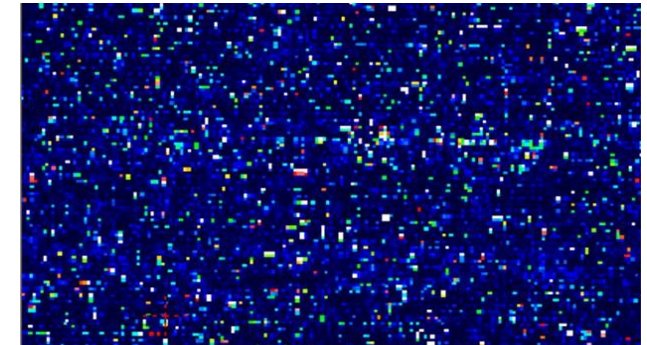
Microarrays

Gene Expression

- Affymetrix
- Agilent
- Illumina
- Whole genome
- Exon-level
- microRNA

CGH, ChIP-Chip, MeIP

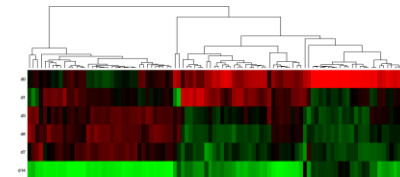
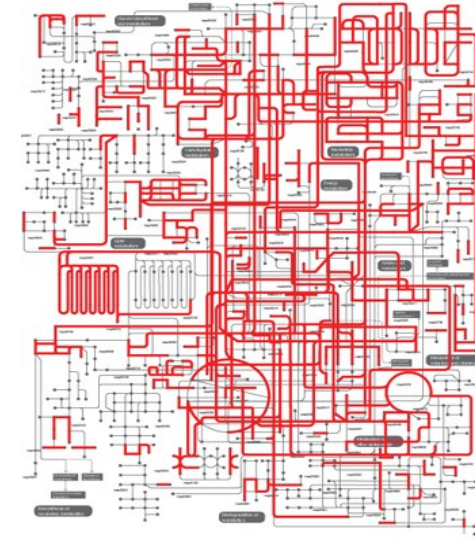
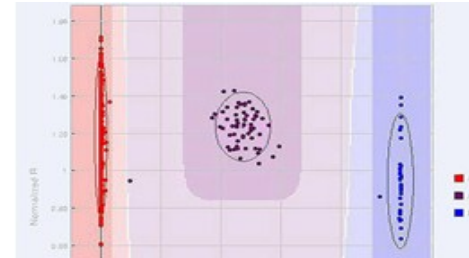
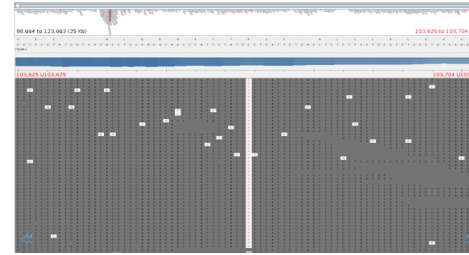
- Nimblegen





Examples of Current research

- Genome (re)sequencing
 - Sheep genome with BGI
 - 21 chicken genomes
 - With IAH: disease resistance
 - With breeders: design SNP chip
 - Bacterial genomics
- Functional genomics
 - Host-pathogen interactions, gene regulation, transcription factor binding, microRNAs, epigenetics
- Metagenomics
 - Gut microbiome: : the “forgotten organ”
 - Ruminants, chickens, others



**Next Generation Sequencing –
The Role of New Sequence Technologies in Shaping the
Future of Veterinary Science**

Hosted by the RCVS Charitable Trust

