# Utilising the genome analysis toolkit (GATK) to identify single nucleotide polymorphisms for use as genetic markers
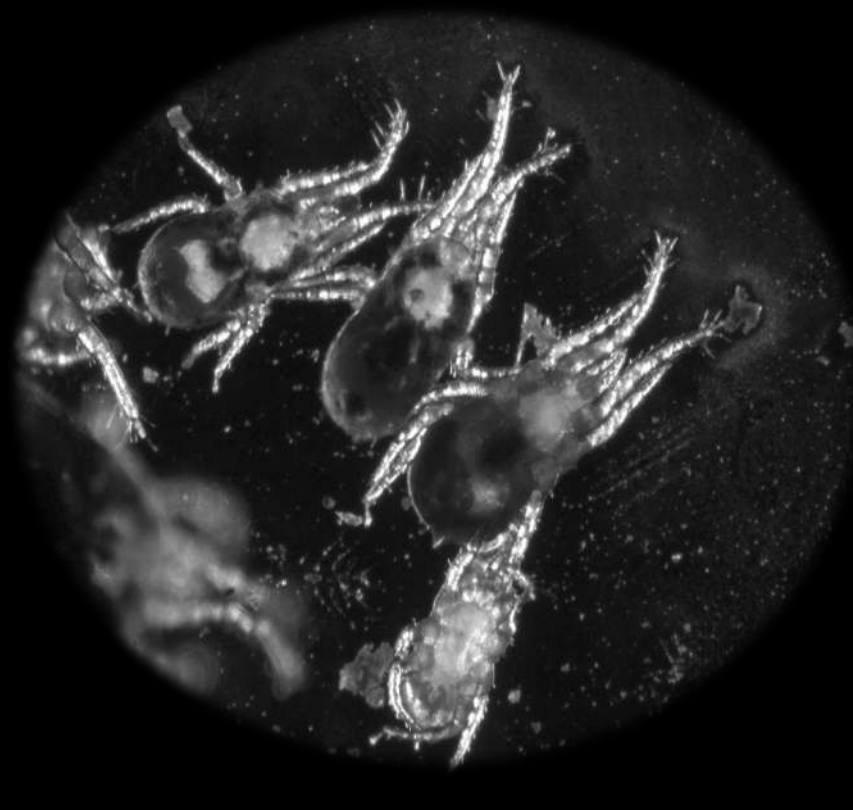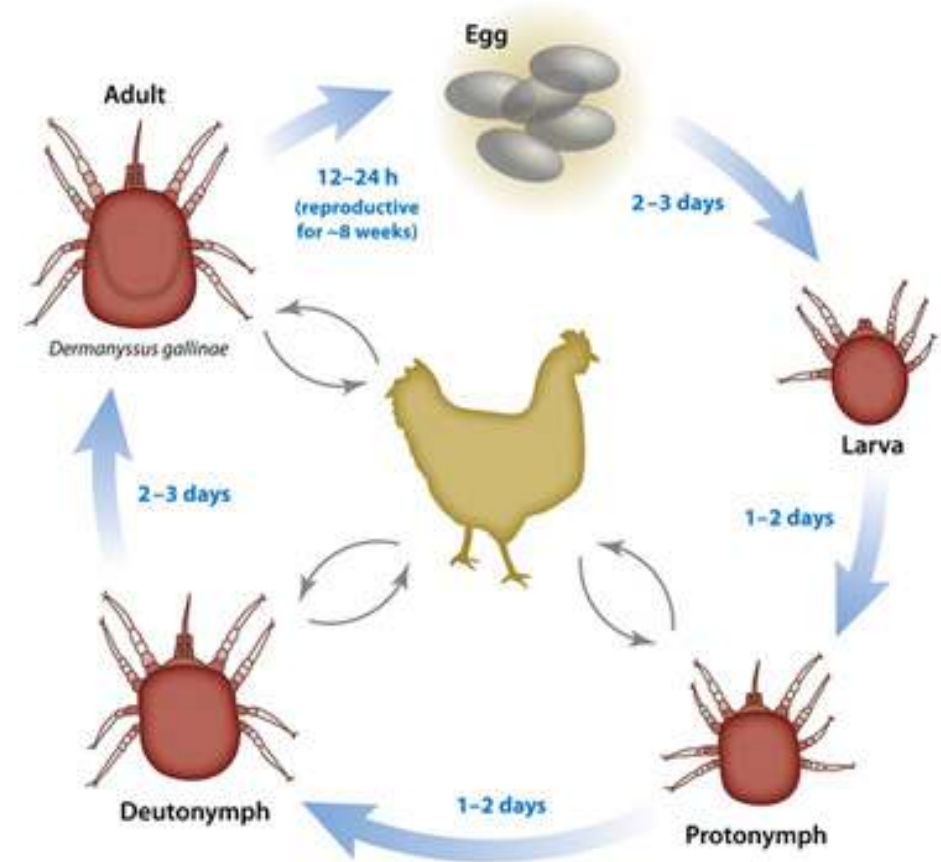
Eleanor Karp-Tatham

**Supervisors**: Prof Damer Blake, Prof Fiona Tomley, Dr Dong Xia & Dr Alasdair Nisbet (Moredun Institute)

# Dermanyssus gallinae

- Blood-feeding ectoparasite
- Significant welfare impact at moderate infestation levels
- €130 million loss to European poultry industry a year
- 28 avian hosts
- Vector for multiple pathogenic agents
- Five stage life-cycle
- Feed during hours of darkness (~30-90 minutes)
- Under optimal conditions (high relative humidity and 20-25°C) the complete cycle can occur in 7-10 days
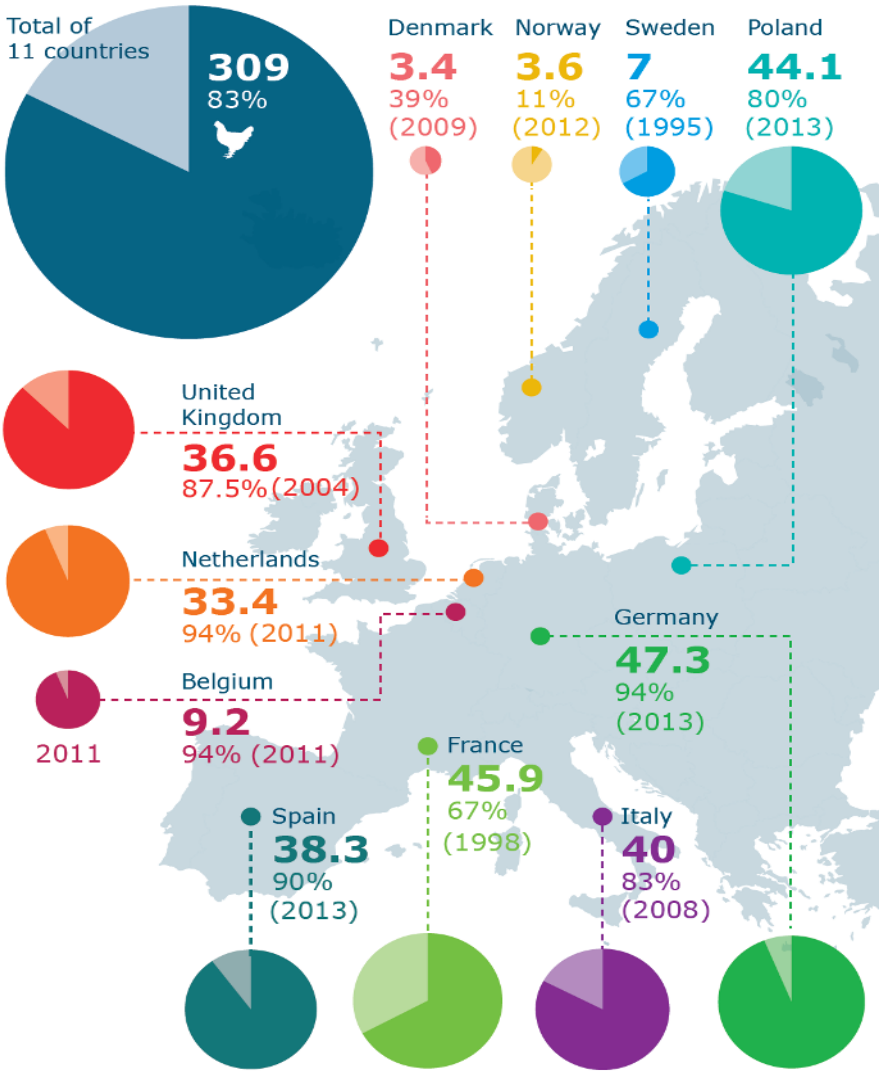


Sparagano OAE, et al. 2014.
Annu. Rev. Entomol. 59:447–66

# Current control of *Dermanyssus gallinae*



INFESTATION OF POULTRY RED MITE IN EUROPE

Number of laying hens per country in millions (2012) and poultry red mite prevalence in percentages.

Total of 11 countries: **309** 83%

Denmark: **3.4** 39% (2009)

Norway: **3.6** 11% (2012)

Sweden: **7** 67% (1995)

Poland: **44.1** 80% (2013)

United Kingdom: **36.6** 87.5% (2004)

Netherlands: **33.4** 94% (2011)

Belgium: **9.2** 94% (2011) — 2011

Germany: **47.3** 94% (2013)

France: **45.9** 67% (1998)

Italy: **40** 83% (2008)

Spain: **38.3** 90% (2013)













Exzolt®
BREAKTHROUGH MITE TREATMENT

# Genetics of *Dermanyssus gallinae*

Whole transcriptome analysis of the poultry red mite *Dermanyssus gallinae* (De Geer, 1778)

SABINE SCHICHT[1], WEIHONG QI[2], LUCY POVEDA[2] and CHRISTINA STRUBE[1]*

[1] Institute for Parasitology, University of Veterinary Medicine Hannover, Buenteweg 17, 30559 Hannover, Germany
[2] Functional Genomics Centre Zurich, Winterthurerstrasse 190, 8057 Zurich, Switzerland

- ➤ Transcriptomic data
- ➤ Illumina
- ➤ Not currently published

Draft Genome Assembly of the Poultry Red Mite, *Dermanyssus gallinae*

Stewart T. G. Burgess,[a] Kathryn Bartley,[a] Francesca Nunn,[a] Harry W. Wright,[a] Margaret Hughes,[b] Matthew Gemmell,[b] Sam Haldenby,[b] Steve Paterson,[b] Stephane Rombauts,[c,d,e] Fiona M. Tomley,[f] Damer P. Blake,[f] James Pritchard,[f] Sabine Schicht,[g] Christina Strube,[g] Øivind Øines,[h] Thomas Van Leeuwen,[i] Yves Van de Peer,[c,d,e,j] Alasdair J. Nisbet[a]

[a] Moredun Research Institute (MRI), Edinburgh, United Kingdom
[b] Centre for Genomic Research, Institute of Integrative Biology, University of Liverpool, Liverpool, United Kingdom
[c] Department of Plant Biotechnology and Bioinformatics, Ghent University, Ghent, Belgium
[d] VIB Center for Plant Systems Biology, Ghent, Belgium
[e] Bioinformatics Institute Ghent, Ghent University, Ghent, Belgium
[f] Department of Pathology and Population Sciences, Royal Veterinary College, Hatfield, Hertfordshire, United Kingdom
[g] Institute for Parasitology, Centre for Infection Medicine, University of Veterinary Medicine Hannover, Hannover, Germany
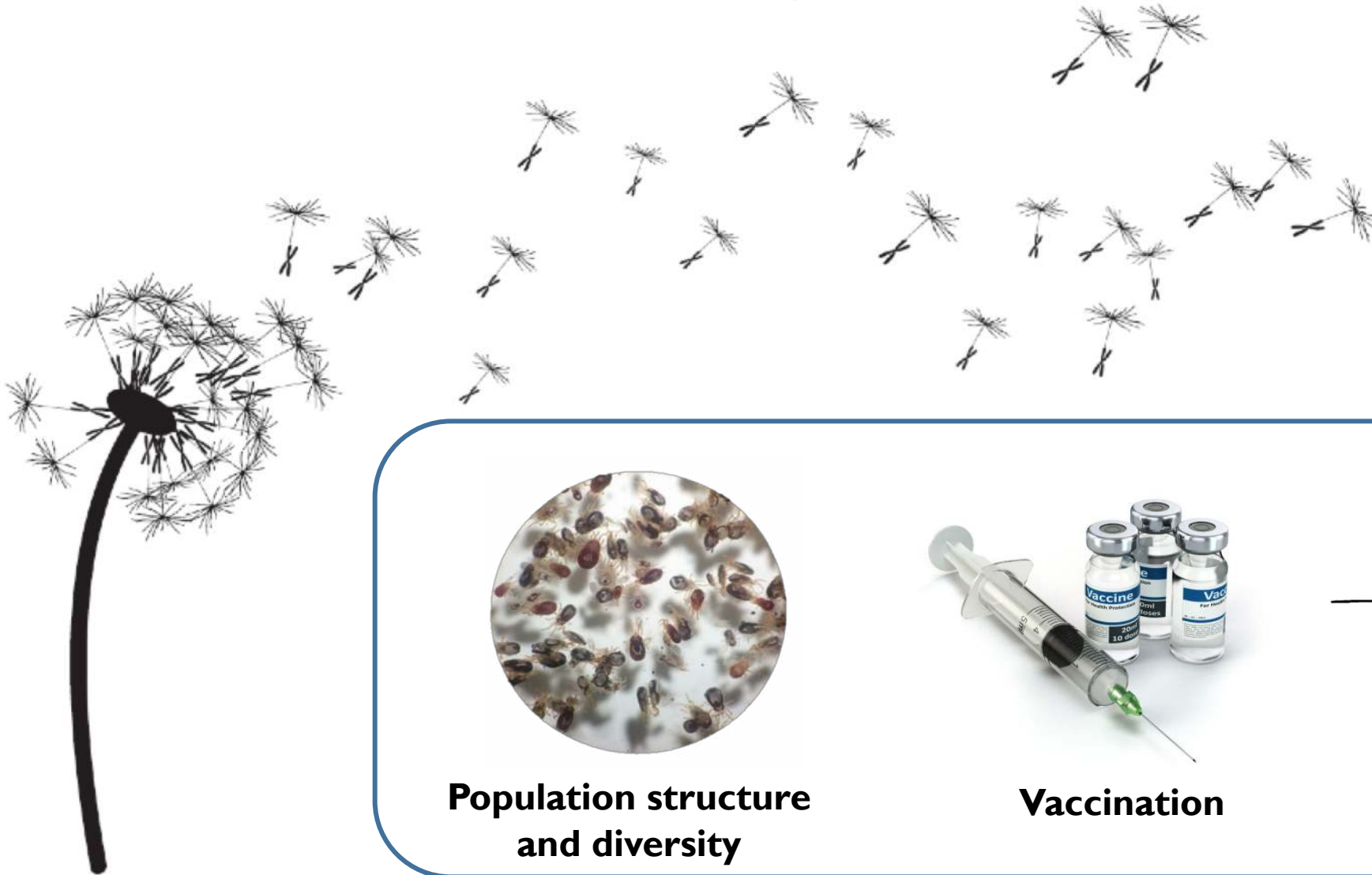[h] Norwegian Veterinary Institute, Oslo, Norway
[i] Department of Plants and Crops, Ghent University, Ghent, Belgium
[j] Department of Biochemistry, Genetics, and Microbiology, University of Pretoria, Pretoria, South Africa

# Population genetics

**Population genetics** is the study of genetic variation within and among populations and the evolutionary factors that explain this variation.



**Population structure and diversity**

**Vaccination**

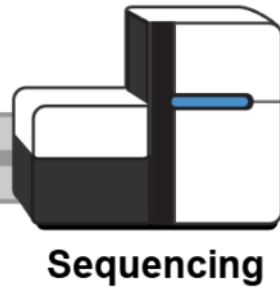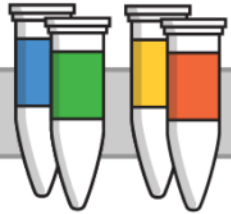**Acaricide/drug resistance**

# Overall plan

**1** **Hannover transcriptome –** 454 sequencing[1]
**Moredun genome assembly** used as reference genome

**2** **Genome Analysis Toolkit**
- GATK best practices followed for Germline SNP & Indel Discovery in Whole Genome and Exome



Sequencing    READS    gatk best practices™    VARIANTS

**4** **MassARRAY** Panel
- MassARRAY panel designed to capacitate 384 samples per plate with up to 40 SNPs per sample
- Option to run two plates - 768 samples with 40 SNPs
- Mites from the UK and Europe
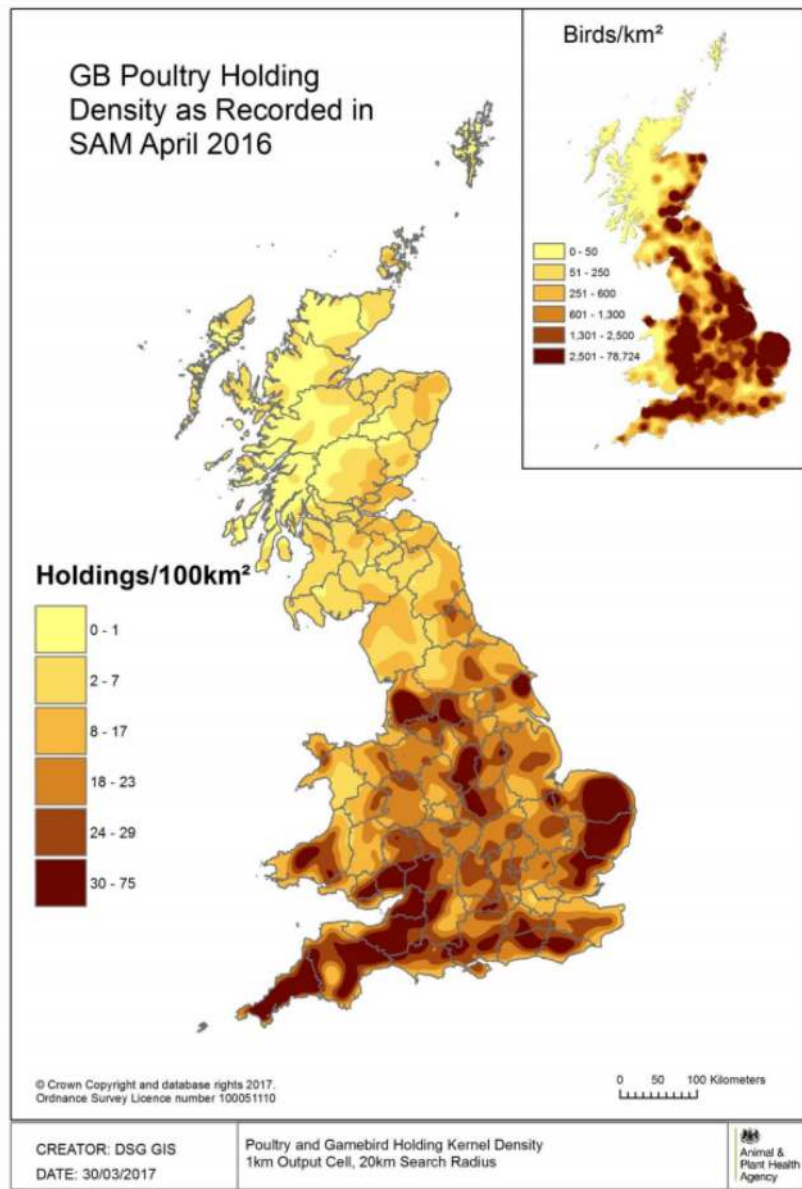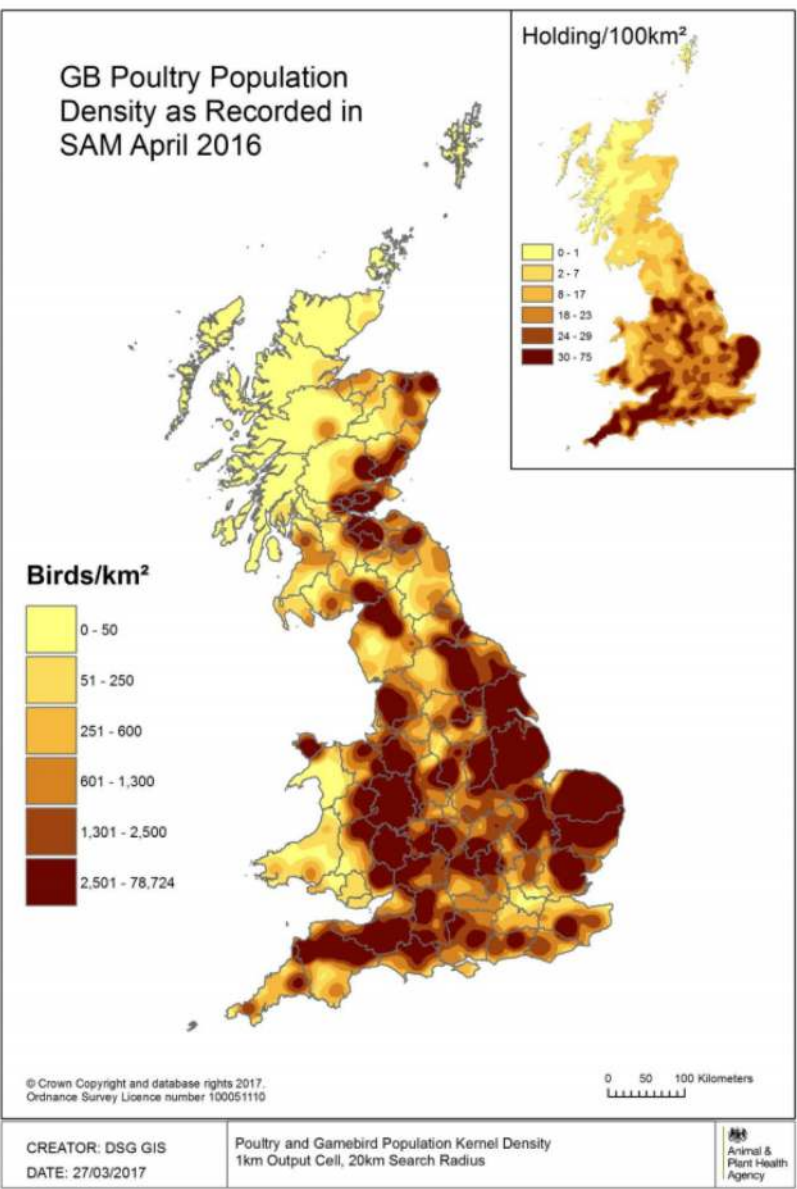- Single mite samples

**3** **Variant calling file**
- VCF stands for Variant Call Format.
- It is a standardized text file format for representing SNP, indel, and structural variation calls
- Used to identify a subset of SNPs with the highest confidence

[1] Schict, S., Qi, W., Poveda, L. and Strube, C. (2014) 'Whole transcriptome analysis of the poultry red mite *Dermanyssus gallinae* (De Geer, 1778)', *Parasitology,* 141(3), pp. 336-46.

# Sample collection: UK

# Sample collection: Europe

# Genome Analysis Toolkit (GATK):

➢ Developed in the Data Sciences Platform at the Broad Institute, the toolkit offers a wide variety of tools with a primary focus on variant discovery and genotyping. Its powerful processing engine and high-performance computing features make it capable of taking on projects of any size.

# GATK on Galaxy

➢ Galaxy servers implement a common core set of tools and reference genomes, and are open to anyone to use. They also contain tools and genomes that are local to each server.

# GATK's Best Practices

➢ Germline short variant discovery (SNPS + Indels)

# GATK workflow



**1**

Initial rounds of SNP calling

**2**

Consolidate into one VCF file with highest quality SNPs

**3**

Run final round of SNP calling using GATK best practices pipeline

**4**

Final table of SNPs produced

**5**

Ready for primer design, panel preparation and everything else ☺

# Pre-processing

**Hannover 1**

Mapped Reads ⓘ

96.4%
1032293

**Hannover 2**

Mapped Reads ⓘ

96.4%
879935

**Raw reads**

Map to reference genome (Moredun)
BWA-MEM
Read groups information auto-assigned to
SAM/BAM specification
Analysis mode: Simply illumina
Default job source parameters

**Mark Duplicates
(Picard)**

**Count co-variates**
Standard set of co-variates used in addition to
those selected
Covariates used: ReadGroupCovariate,
QualityScoreCovariate, CycleCovariate,
DinucCovariate
Basic GATK options
Basic Advanced analysis options

**Analyse covariates**

**Table recalibration**
Basic GATK options selected
Basic analysis options selected

**Count co-variates**
Standard set of co-variates used in addition to
those selected
Covariates used: ReadGroupCovariate,
QualityScoreCovariate, CycleCovariate,
DinucCovariate
Basic GATK options
Basic Advanced analysis options

**Analyse covariates**

**Analysis-
Ready reads**

# Variant Discovery



**Analysis Ready Reads**

**Unified Genotyper**
Reference: Moredun
Table calibrated BAM file
BOTH option selected for genotype likelihood calculation model
Minimum phred-scaled confidence threshold at which variants not at 'trigger' track sites should be called set to 20.0
Minimum phred-scaled confidence threshold at which variants not at 'trigger' track sites should be emitted set to 20.0
Advanced analysis option selected and the following annotations selected; FisherStrand, HaplotypeScore, HomopolymerRun, MappingQualityRankSumTest, QualbyDept and ReadPosRankSumTest
Standard selected for annotation interfaces/groups

~~Variant filtration~~

**Combine variants**
Input file - SNP
Input file - Indel
Reference genome: Moredun
Basic GATK selected
Basic Analysis options selected

**Variant Annotation**
Reference: Moredun
Variant file: Combined variants
Bam file: Table recalibrated
All possible annotations selected except MVLikelihoodRatio
No additional annotations
No Binding for reference-ordered comparison data
Do not set snpEff
No Expression
Basic GATK options used
All annotation interfaces and groups applied
Genotype quality set to 0
Annotations to exclude set to MVLikelihoodRatios

**Select variants**
Unified genotype file
Reference: Moredun
Advanced analysis option - under the select only a certain type of variant from input file' - **SNP** chosen
No other settings changed

**Select variants**
Unified genotype file
Reference: Moredun
Advanced analysis option - under the select only a certain type of variant from input file' - **INDEL** chosen
No other settings changed

**Initial round of SNP calling complete**
**Output: VCF file**

# Initial SNP Tables

| Dataset | Read sets | File size | Sequencing platform | Read type | Quality scores | GATK run1 complete | Mapping results | No. of SNPs | VCF intersect |
|---------|-----------|-----------|---------------------|-----------|----------------|--------------------|-----------------|-------------|---------------|
| Hannover | 1 | 12.9MB | 454 | Single | Y | Y | 96.4% | 63,592 | **39,396** |
| | 2 | 13.5MB | 454 | Single | Y | Y | 96.4% | 69,440 | |

# GATK workflow



**39,396**

**1** Initial rounds of SNP calling

**2** Consolidate into one VCF file with highest quality SNPs

**3** Run final round of SNP calling using GATK best practices pipeline

**4** Final table of SNPs produced

**5** Ready for primer design, panel preparation and everything else ☺

Hannover reads 1

Hannover reads 2

# Variant Discovery

# In reality…

# Variant Filtration

# Variant Filtration

| Dataset | Read sets | GATK run1 complete | Mapping results | No. of SNPs | VCF intersect | GATK run 2 complete | No. of SNPs total | No. of SNPS PASS | No. of excluded SNPS | VCF intersect | VCF intersect PASS |
|---------|-----------|--------------------|-----------------|-------------|---------------|---------------------|-------------------|------------------|----------------------|---------------|--------------------|
| Hannover | 1 | Y | 96.4% | 63,592 | 39,396 | Y | 66,296 | 65,248 | 1048 | 32,940 | **32,599** |
| | 2 | Y | 96.4% | 69,440 | | Y | 69,440 | 68294 | 1146 | | |

## -6,797 SNPs

# Selecting SNPs

## 32,599 SNPs

**1** High quality scoring **SNPs** in coding regions

**2** Bioinformatics tools to predict **SNP** effect

**3** SNPs related to vaccine targets

**4** SNPs relating to acaricide/drug resistance

# Selecting SNPs



**VCF Intersect File**

Hannover read sets 1 +2

32,599 SNPs

↓

Read depth set to a minimum of 50

↓

PHRED quality score set to a minimum of 500

↓

Local BLAST alignment against genome
*No. of threads: 1, Filter low complexity*
*Expect: 10, Word size: 11, Match 1, Mismatch -3, Existence 5, Extension 2, Max no. of hits: 1,500*
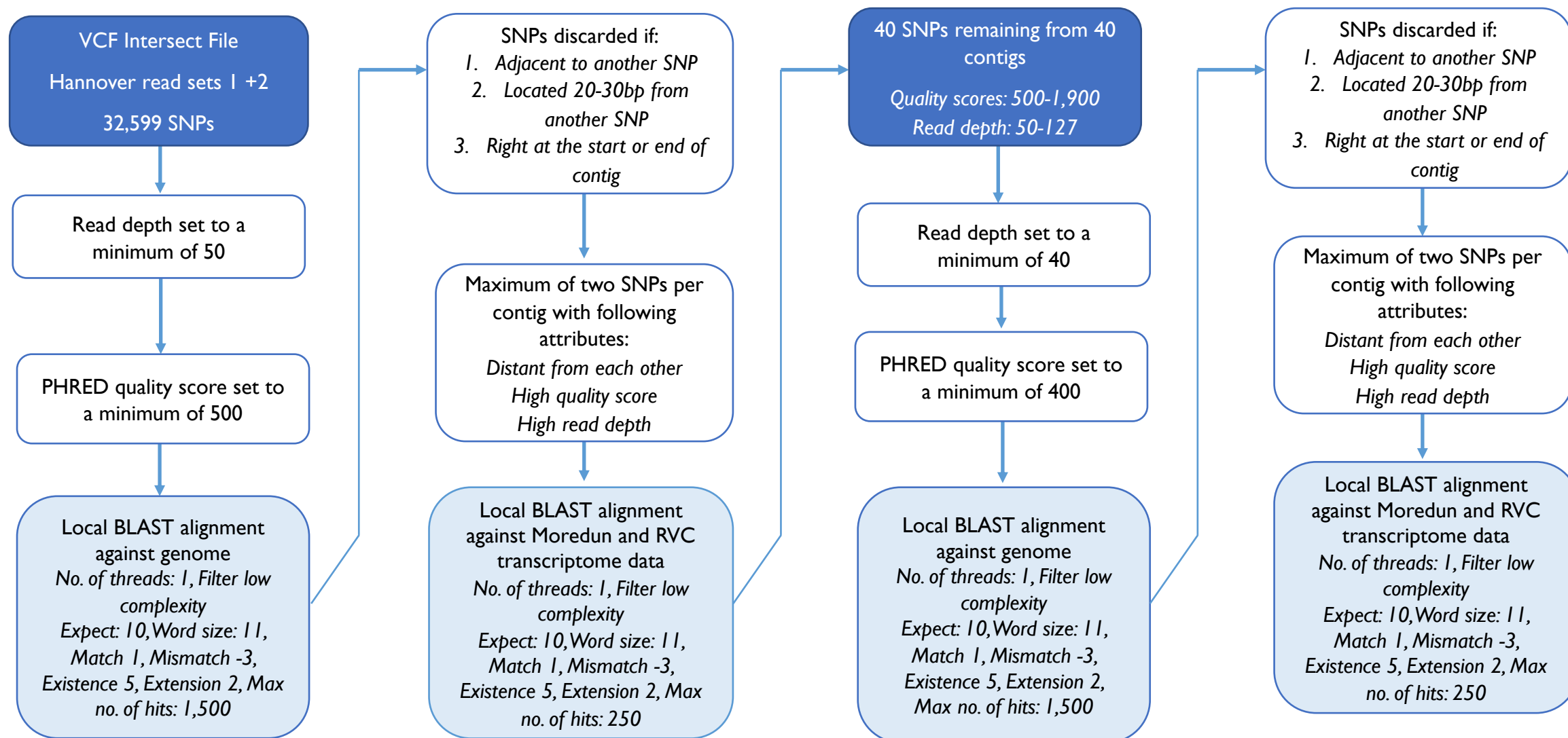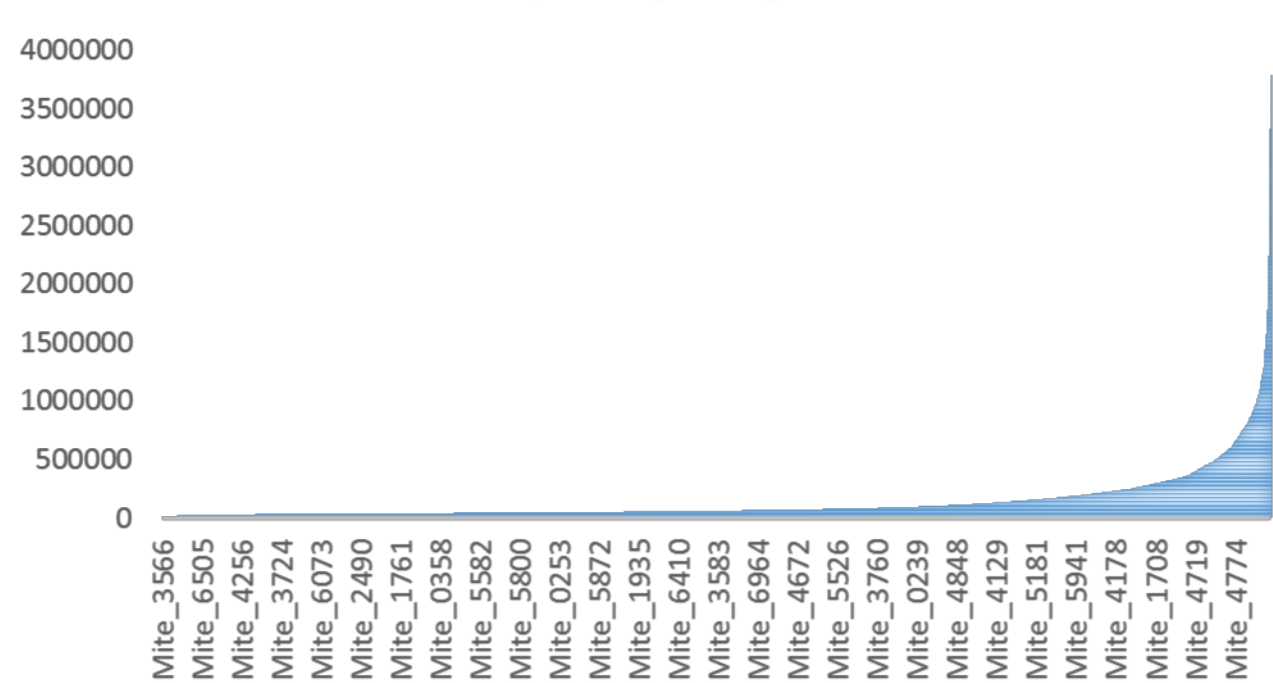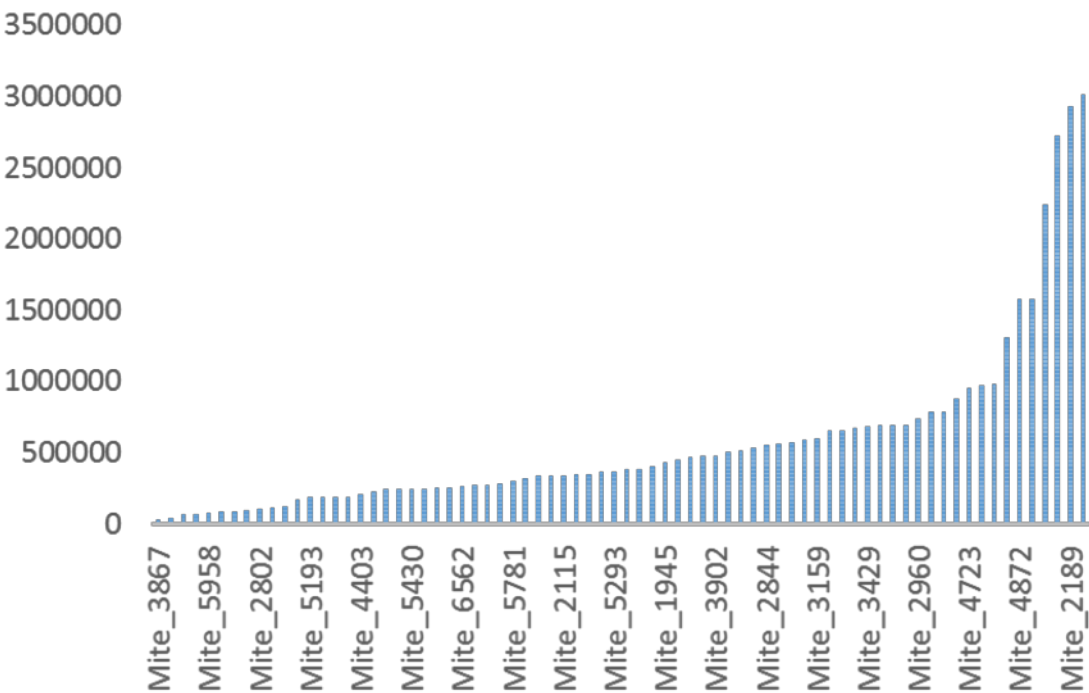
SNPs discarded if:
1. *Adjacent to another SNP*
2. *Located 20-30bp from another SNP*
3. *Right at the start or end of contig*

↓

Maximum of two SNPs per contig with following attributes:
*Distant from each other*
*High quality score*
*High read depth*

↓

Local BLAST alignment against Moredun and RVC transcriptome data
*No. of threads: 1, Filter low complexity*
*Expect: 10, Word size: 11, Match 1, Mismatch -3, Existence 5, Extension 2, Max no. of hits: 250*

**40 SNPs remaining from 40 contigs**

*Quality scores: 500-1,900*
*Read depth: 50-127*

↓

Read depth set to a minimum of 40

↓

PHRED quality score set to a minimum of 400

↓

Local BLAST alignment against genome
*No. of threads: 1, Filter low complexity*
*Expect: 10, Word size: 11, Match 1, Mismatch -3, Existence 5, Extension 2, Max no. of hits: 1,500*

SNPs discarded if:
1. *Adjacent to another SNP*
2. *Located 20-30bp from another SNP*
3. *Right at the start or end of contig*

↓

Maximum of two SNPs per contig with following attributes:
*Distant from each other*
*High quality score*
*High read depth*

↓

Local BLAST alignment against Moredun and RVC transcriptome data
*No. of threads: 1, Filter low complexity*
*Expect: 10, Word size: 11, Match 1, Mismatch -3, Existence 5, Extension 2, Max no. of hits: 250*

# Selecting SNPs

- 75 SNPS
- Contig size: 27,617-3,015,868
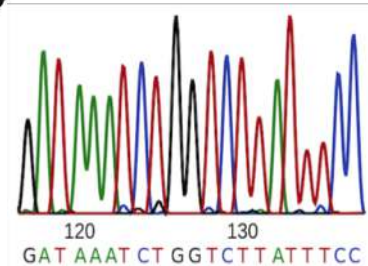- Quality score: 400-2552
- Read depth: 40 - 127

# What's next?

**1** **Validation of SNPs**

- Selected SNPs tested in house
- Primer design followed by PCR and sequencing to further validate GATK workflow

**2** **Send to MassARRAY**

- Final panel of SNPS sent to MassARRAY company
- Primers designed to encompass target SNPs
- Sample preparation undertaken
- Wait for processing time

**3** **Computational analysis**

- Diversity assessed across UK and European samples
- Inter-farm and intra-farm diversity investigated
- Phylogenetic analysis conducted

# Acknowledgements

COREMI