



EAAP 2016 - 67th Annual Meeting of the
European Federation of Animal Science
Belfast UK, 29 Aug – 2 Sept 2016

A family-based **imputation**
algorithm
for high- and low-
coverage **sequence data**

Mara Battagin*
Gregor Gorjanc
Serap Gonen
Roberto Antolin
John Hickey

* mara.battagin@roslin.ed.ac.uk



THE UNIVERSITY *of* EDINBURGH



Introduction

Genomic Selection 2.0



Sequence pigs and chicken animals at high- and low-coverage

Phasing and imputation for all individuals

> 300,000 pigs and 250,000 chicken with sequence information

Who?

AlphaFamSeq

Genomic Selection

Distribution of sequencing resources

Unequal

Equal

i.e. All population at low coverage

Key sires

Focal individuals

Individuals whose **haplotypes** are shared with many other individuals in a population.

AlphaSeqOpt (by Gonen et al.)



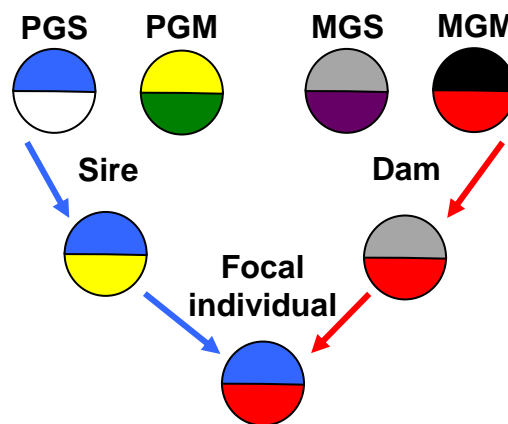
**“heuristic phasing algorithm for
high- and low-coverage sequence data”**

AlphaFamSeq

**AlphaSuite
program**

**Family-based
method**

**Developed
for sequence data**



Method AlphaFamSeq



Reads ● — Snp —>

nRef	0	20	0	0	2	0	14	0	0
nAlt	0	12	12	0	37	0	1	0	0

AlphaFamSeq

A family-based imputation algorithm using high- and low-coverage sequence data



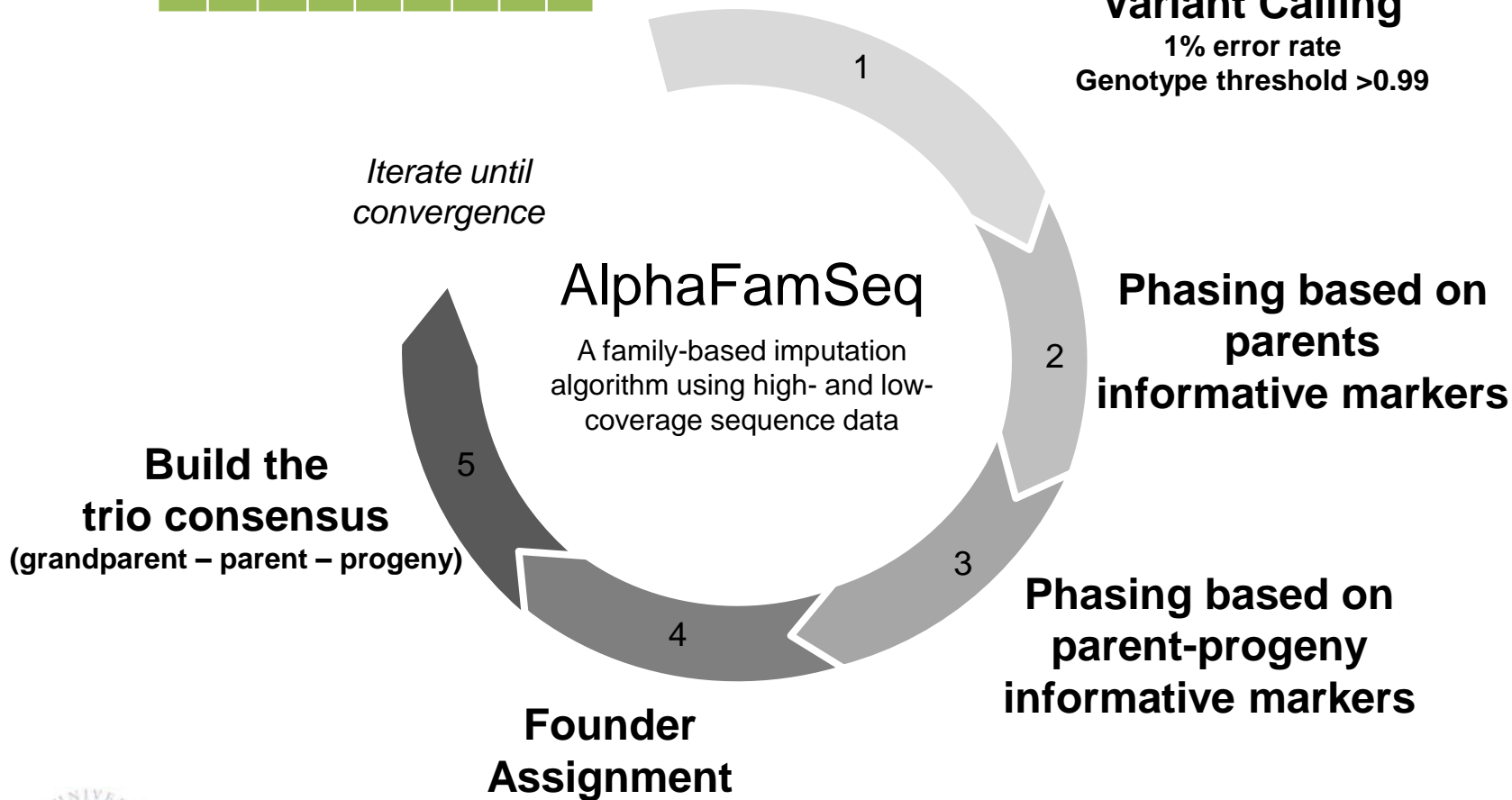
Method AlphaFamSeq



Reads ● — Snp —>

nRef	0	20	0	0	2	0	14	0	0
nAlt	0	12	12	0	37	0	1	0	0

Variant Calling
1% error rate
Genotype threshold >0.99



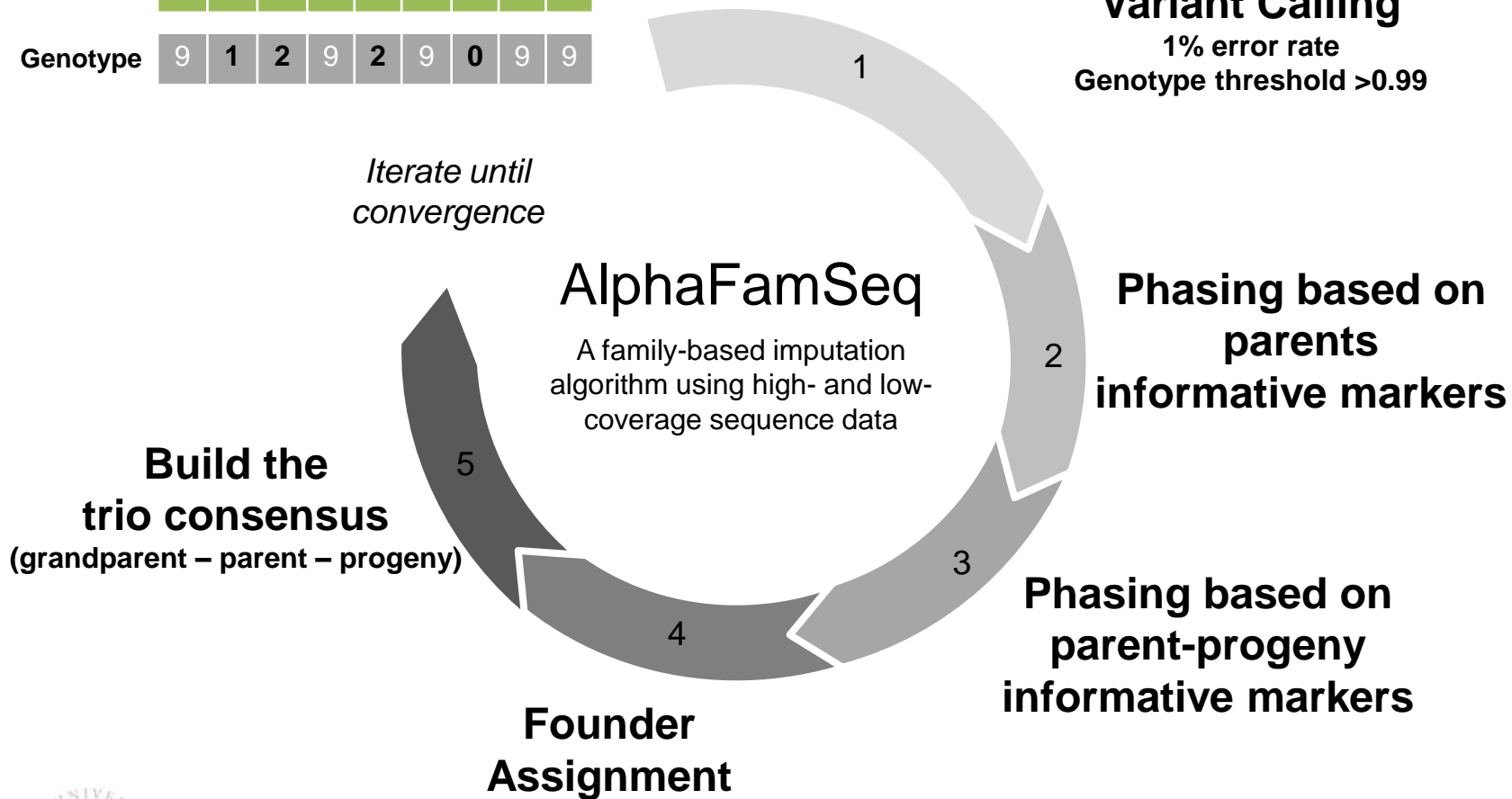
Method AlphaFamSeq



Reads ● — Snp —>

nRef	0	20	0	0	2	0	14	0	0
nAlt	0	12	12	0	37	0	1	0	0
Genotype	9	1	2	9	2	9	0	9	9

Variant Calling
1% error rate
Genotype threshold >0.99



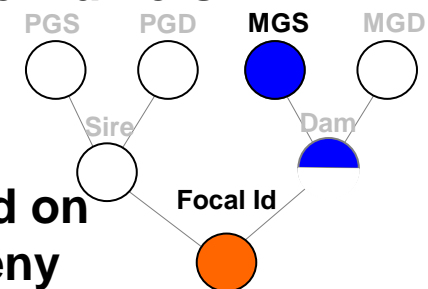
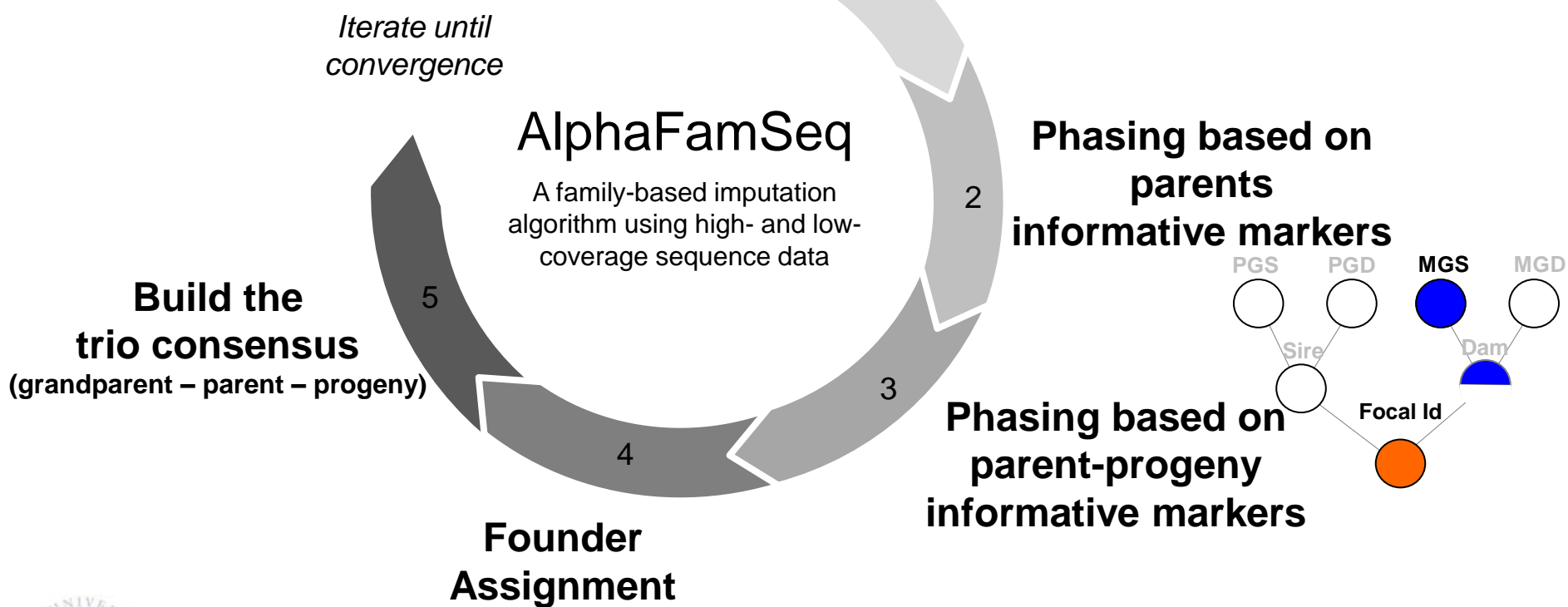
Method AlphaFamSeq



Reads ● — Snp —>

nRef	0	20	0	0	2	0	14	0	0
nAlt	0	12	12	0	37	0	1	0	0
Genotype	9	1	2	9	2	9	0	9	9

Variant Calling
1% error rate
Genotype threshold >0.99



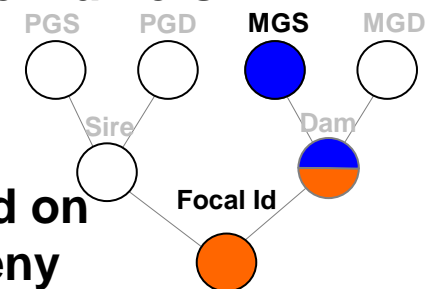
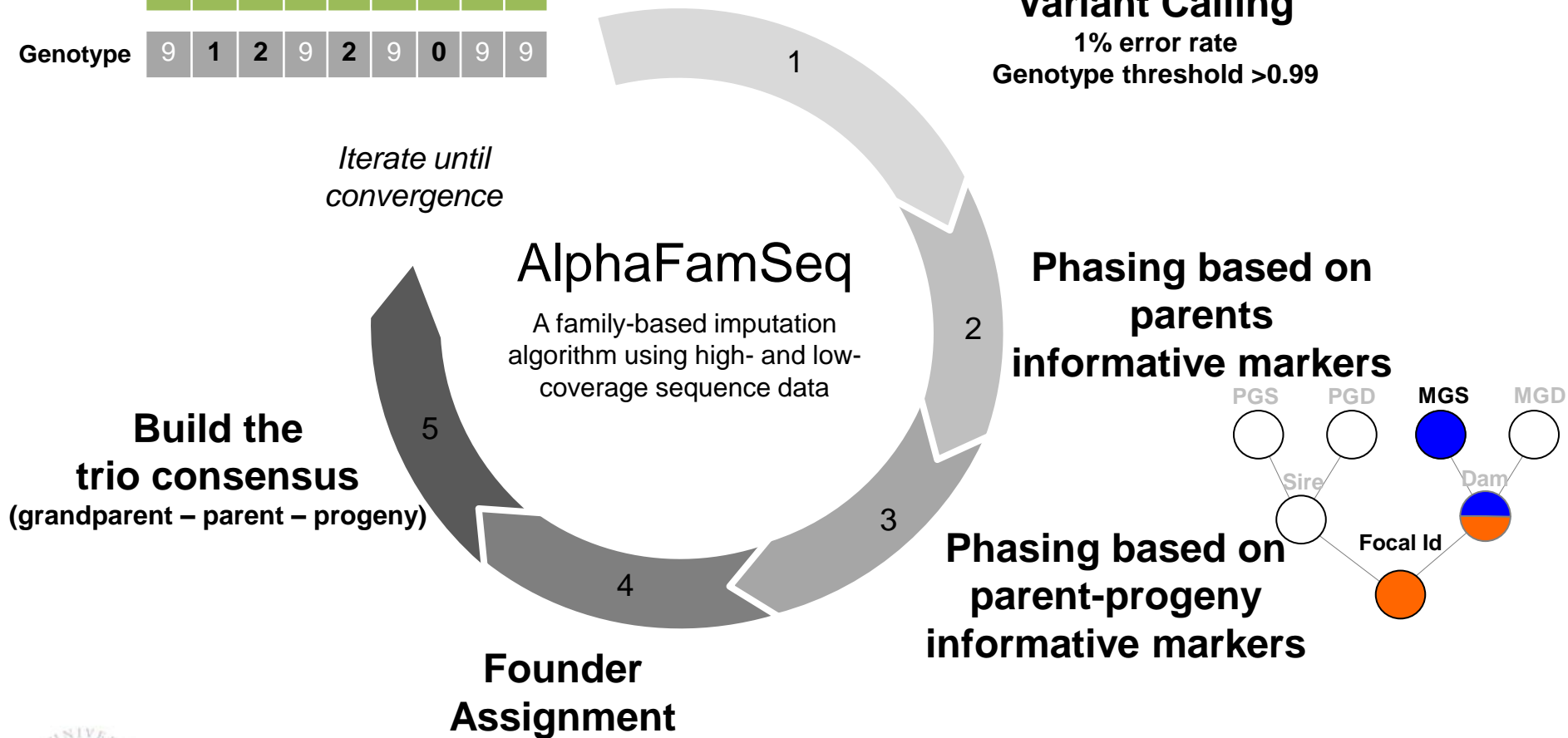
Method AlphaFamSeq



Reads ● — Snp —>

nRef	0	20	0	0	2	0	14	0	0
nAlt	0	12	12	0	37	0	1	0	0
Genotype	9	1	2	9	2	9	0	9	9

Variant Calling
1% error rate
Genotype threshold >0.99



Method AlphaFamSeq



Reads ● — Snp —>

nRef	0	20	0	0	2	0	14	0	0
nAlt	0	12	12	0	37	0	1	0	0
Genotype	9	1	2	9	2	9	0	9	9

Variant Calling
1% error rate
Genotype threshold >0.99

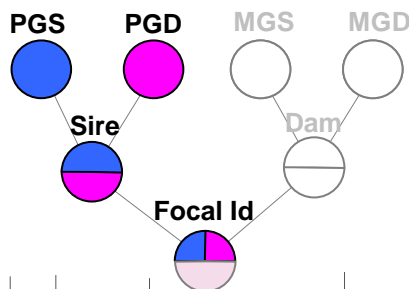
Iterate until convergence

AlphaFamSeq

A family-based imputation algorithm using high- and low-coverage sequence data

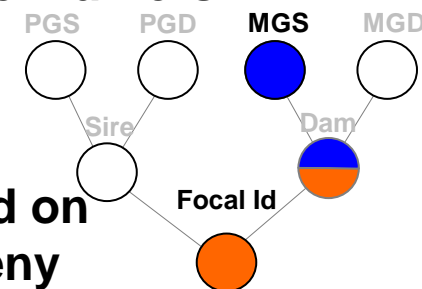
Phasing based on parents informative markers

Build the trio consensus
(grandparent – parent – progeny)



Founder Assignment

Phasing based on parent-progeny informative markers



Method

Algorithm tested through Simulation



Simulate true genotypes (**AlphaSim**)

10 chromosomes
1 Morgan length
100,000 segregating sites

2 scenarios

Within Family imputation

Selected 5 families as replicates

Multiple Families imputation

External **cattle** pedigree (14,814 ids)

Simulate sequence data different scenarios
Number of family members sequenced (from 1 to 7)
Sequencing depth (from 0 to 30x)

Extracted 3 generations for 500 sires

Total 1,614 animals

Key sires

Equal distribution

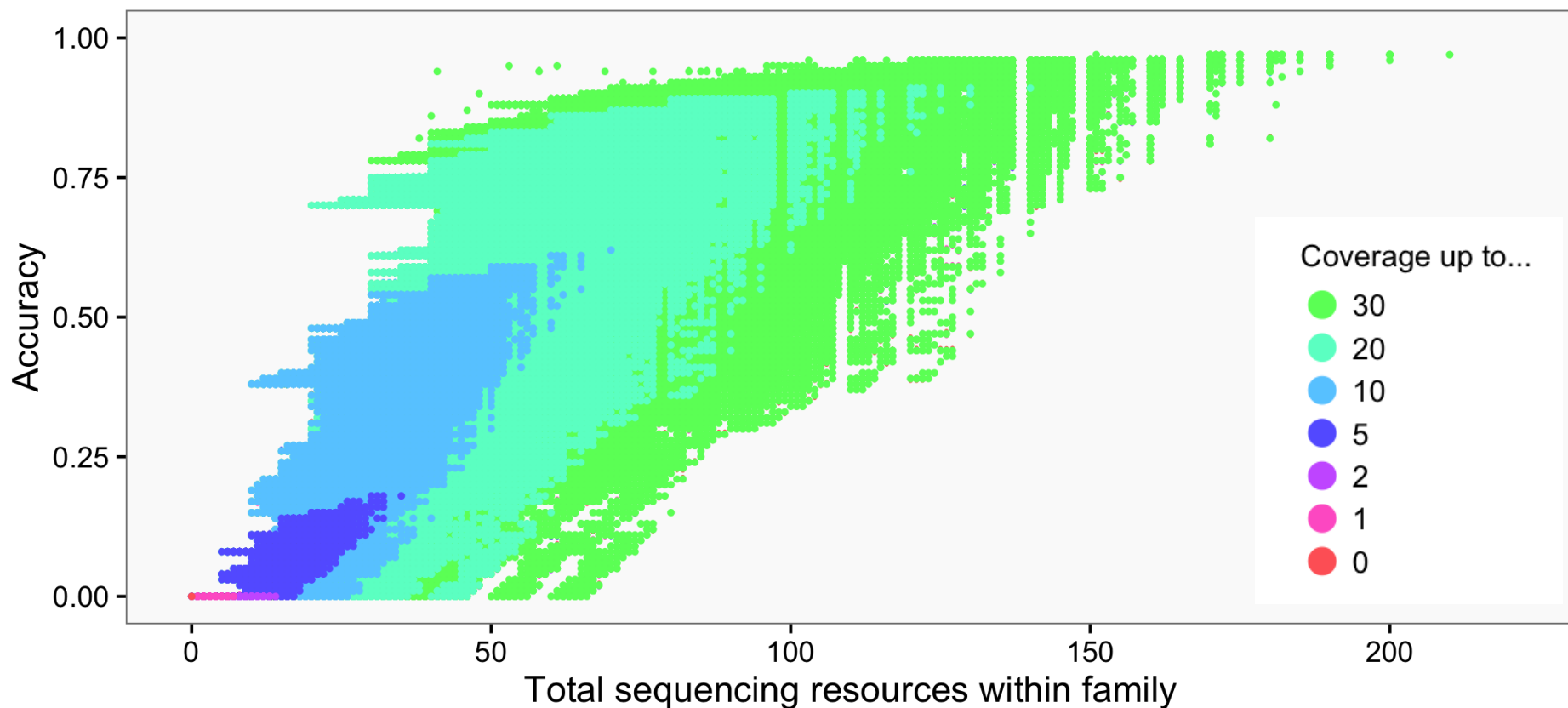
**Accuracy of phasing the focal individual
for 823,543 scenarios**

Budget
£ 420,000

Focal individuals

AlphaSeqOpt
(by Gonen et al.)





Cost of the library/individual = £40

Cost of 1x whole genome = £80

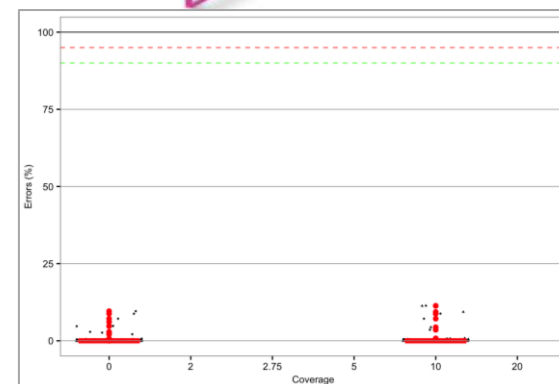
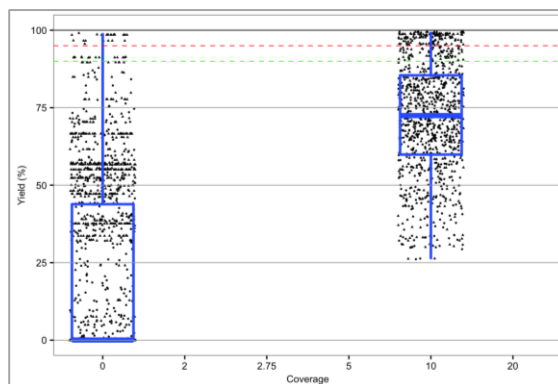
Results

Accuracy of phasing multiple families



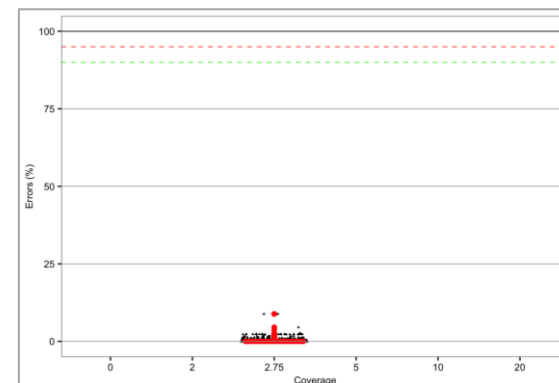
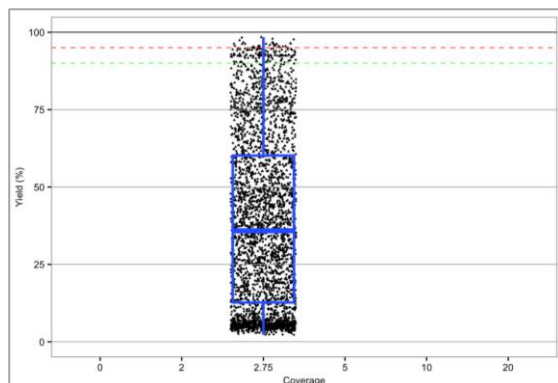
Percentage of **Correct** and **Wrong** phased gamete by read depth

Key Sires
Min 0.00
Mean 35.41
Max 99.50



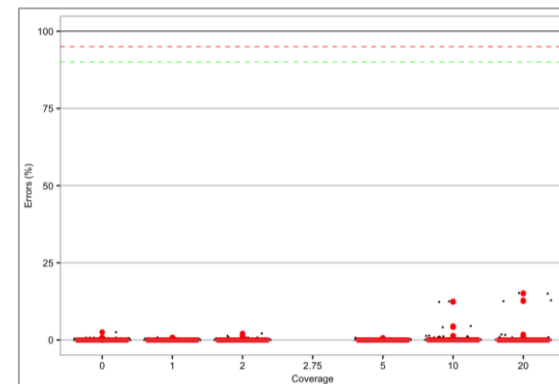
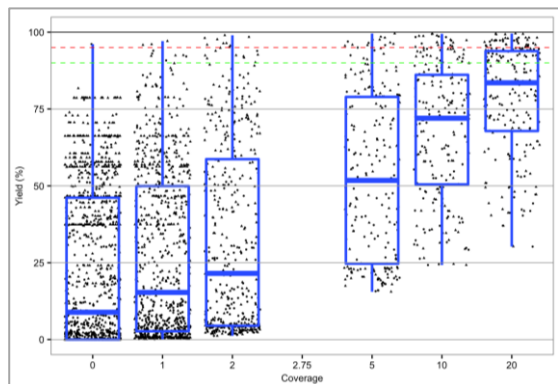
All population

Min 2.25
Mean 39.05
Max 98.60



Focal individuals

Min 0.00
Mean 34.28
Max 99.55



Heuristic phasing algorithm for
high- and low-coverage sequence data

Able to increase the **sequence information**
in livestock population

Good accuracy of can be achieved

ToDoList:
Test it with the **real** data!

Acknowledgements



Funding Bodies

**Thank
you!**



Industrial Partners

mara.battagin@roslin.ed.ac.uk
John Hickey
<http://www.alphagenes.roslin.ed.ac.uk/>
Gregor Gorjanc

Roberto Antolin

Serap Gonen

Diarmaid De Burca

Andrew Derrington

Chris Gaynor

Paolo Gottardo

Janez Jenko

Daniel Money

Roger Ros Freixedes

Maria Sanchez Perez

Masoud Shirali

Adriana Somavilla



KWS



Method multiple-families sequence data



Coverage	Unequal distribution of resources		Equal distribution of resources
	Key Sires	Focal Individuals	
0x	916	635	
1x		420	
2x		224	
2.75x			1,614
5x		110	
10x	500	116	
20x		109	
30x			
Cost (£)	420,000	419,640	419,640

Cost of the library/individual = £40

Cost of 1x whole genome = £80



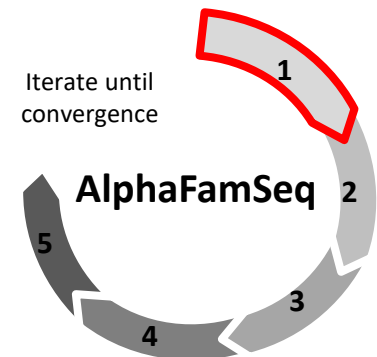
1

Variant calling

• ————— Snp —————>

Reads	Snp								
nRef	12	20	0	26	2	1	14	0	11
nAlt	0	2	12	2	37	0	0	2	0
Genotype	1	0	1	0	2	1	1	2	1
PatHap	9	0	0	0	1	9	1	1	0
MatHap	9	0	1	0	1	9	0	1	1

Parameters:
Error Rate = 0.01
Genotype threshold = 0.999



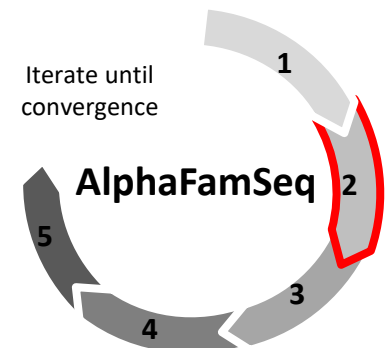
2

Phasing based on parents informative markers

Phase the markers inherited from homozygous parents

Sire	
PatHap	0 1 1 1 0 0 1
MatHap	0 1 1 1 0 0 1

Progeny	↓	↓	↓					
	PatHap	0	1	1	1	0	0	1
MatHap	1	1	1	1	0	1	1	
Genotype	2	1	2	2	0	1	1	2



3

Phasing based on parent-progeny informative markers

Phase the heterozygous markers of individuals that have a parent and progeny with opposite alleles

GrandSire

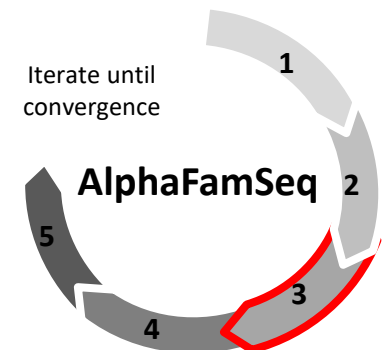
hPat	1	0		1		0	0		
hMat	1	0		1		0	0		

Dam

hPat				1		0			
hMat				0		1			

Progeny

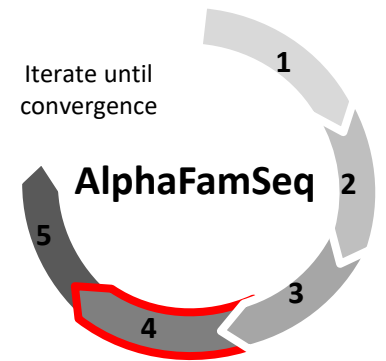
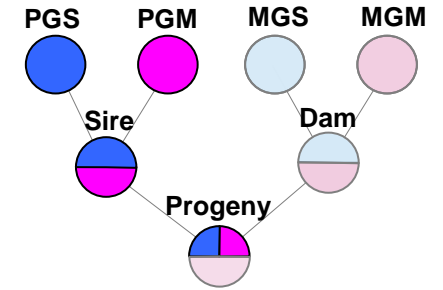
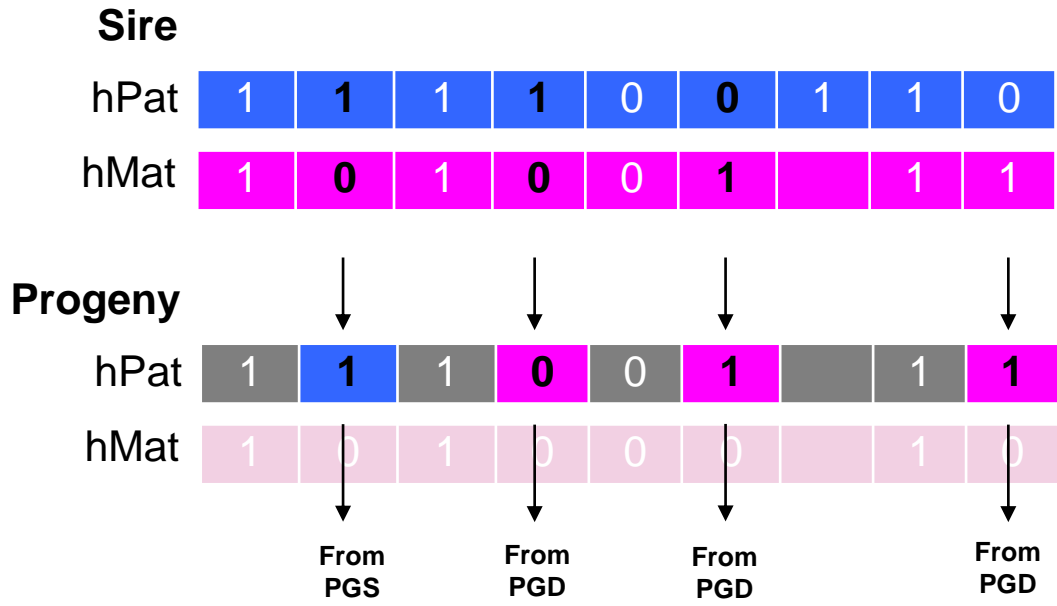
hPat	0	1	1	0		1	0		1
hMat		1	0	0		1	1		1



4

Founder Assignment

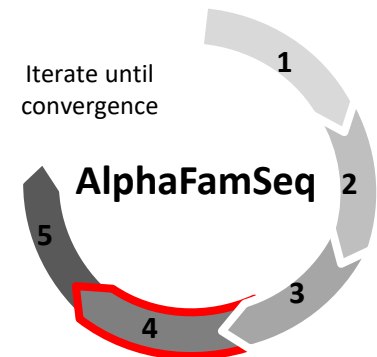
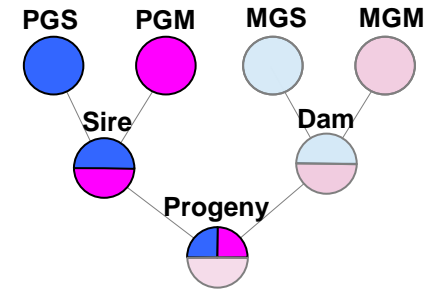
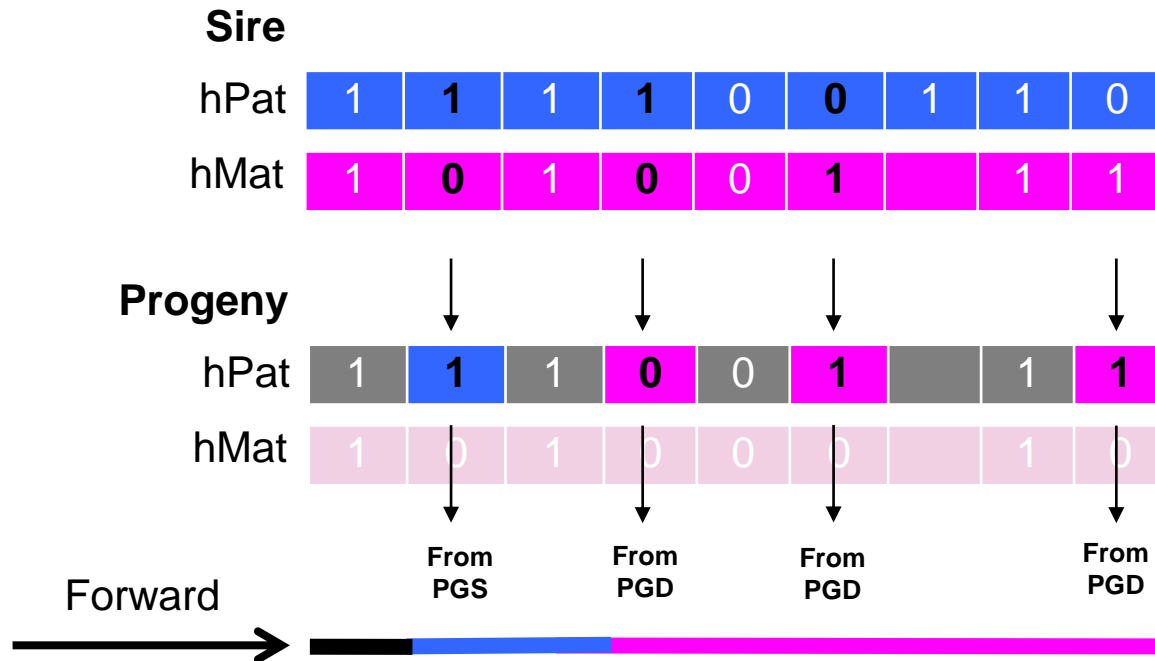
Identify trios (grandparent, parent, individual) by assigning grandparents to individuals using the heterozygous markers of their parents



4

Founder Assignment

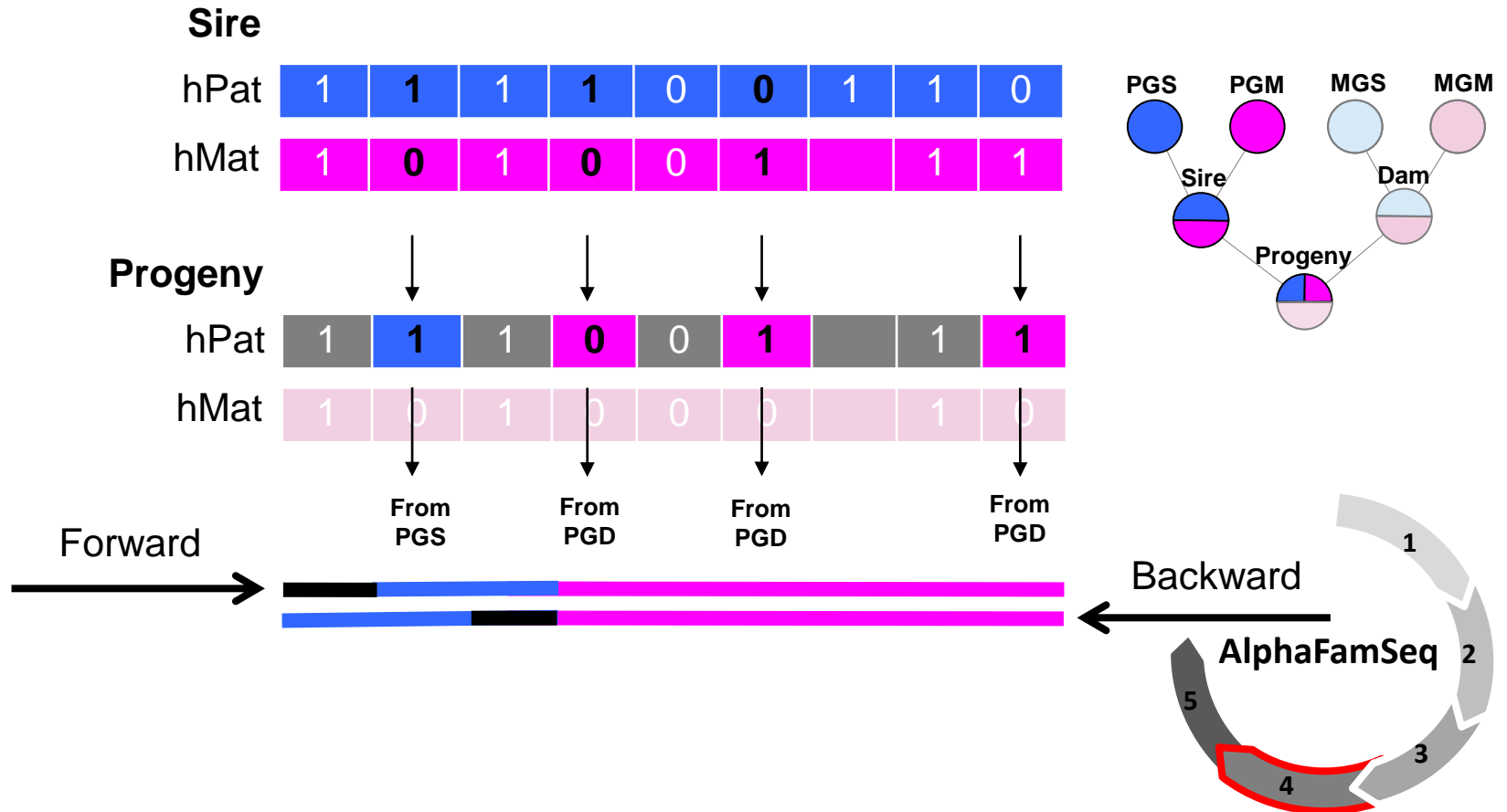
Identify trios (grandparent, parent, individual) by assigning grandparents to individuals using the heterozygous markers of their parents



4

Founder Assignment

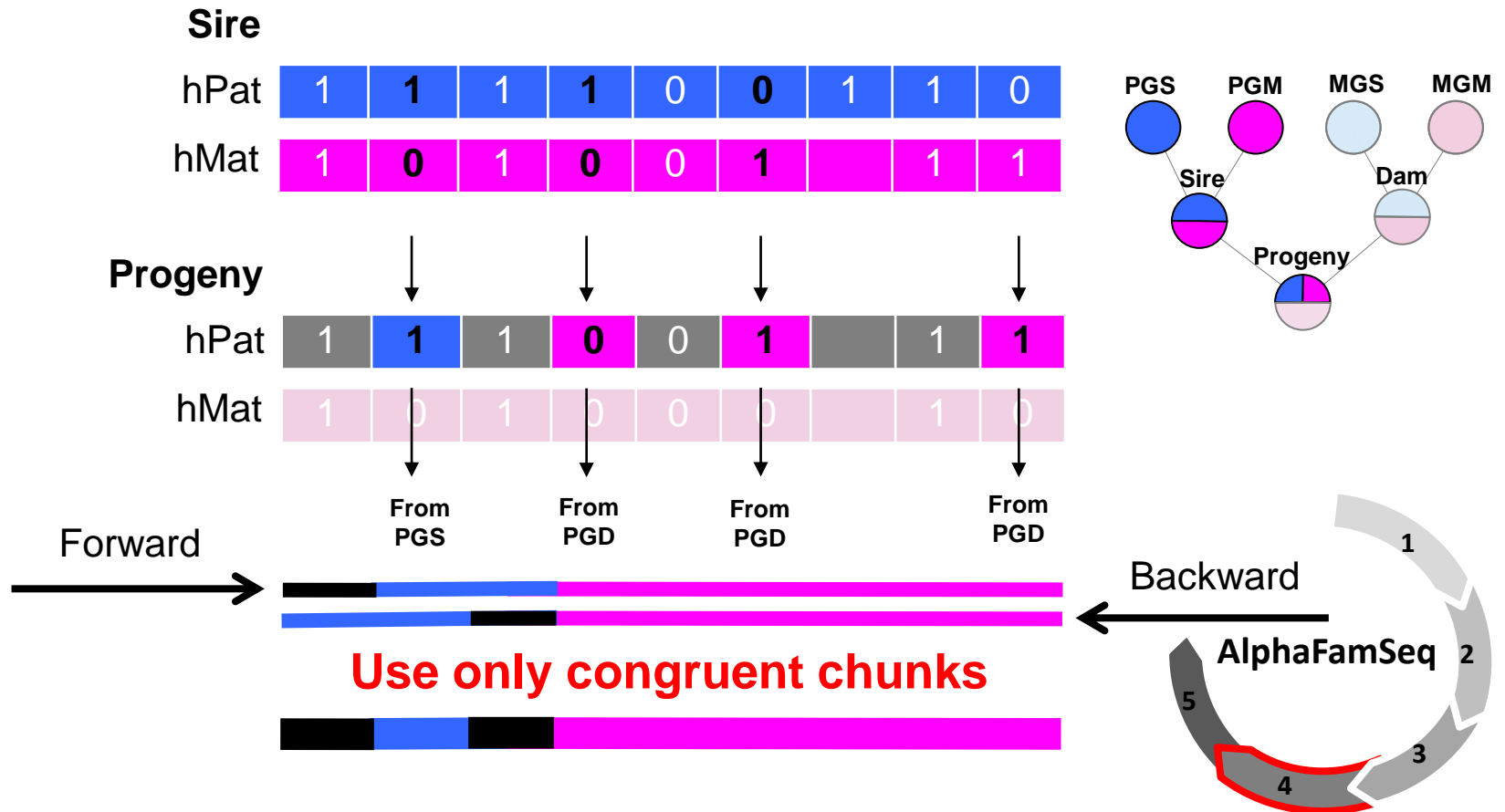
Identify trios (grandparent, parent, individual) by assigning grandparents to individuals using the heterozygous markers of their parents



4

Founder Assignment

Identify trios (grandparent, parent, individual) by assigning grandparents to individuals using the heterozygous markers of their parents



5

Build the trio consensus

Identify chunks of haplotypes within trios and use the chunks to impute missing markers between members of the trio.

GrandDam

hPat	1	0		0	0	0	0	0
------	---	---	--	---	---	---	---	---

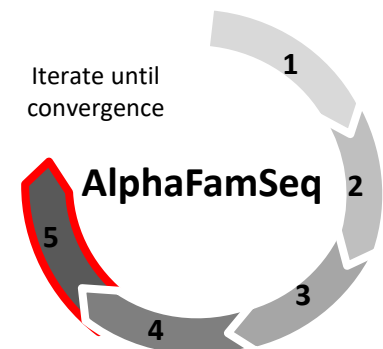
hMat	1	0		0	0	<u>1</u>	0	1	<u>1</u>
------	---	---	--	---	----------	----------	---	----------	----------

Sire

hMat	1	0	1	0	0	<u>1</u>	0	1	<u>1</u>
------	---	---	---	---	---	----------	----------	---	----------

Progeny

hPat	1	1	1	0	0	1	0	1	1
------	---	---	---	---	---	---	----------	---	---



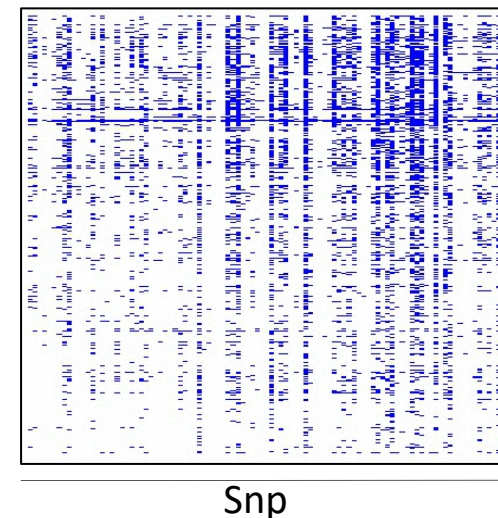
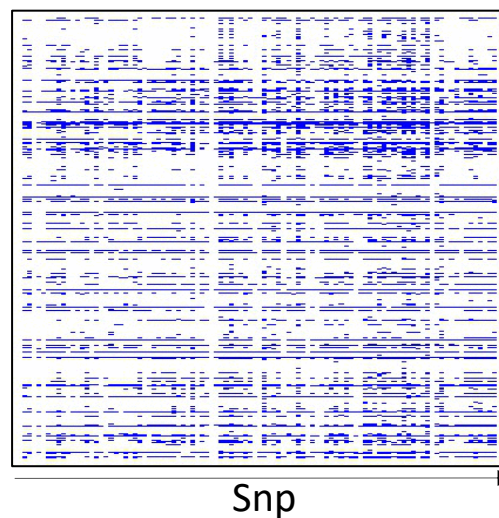
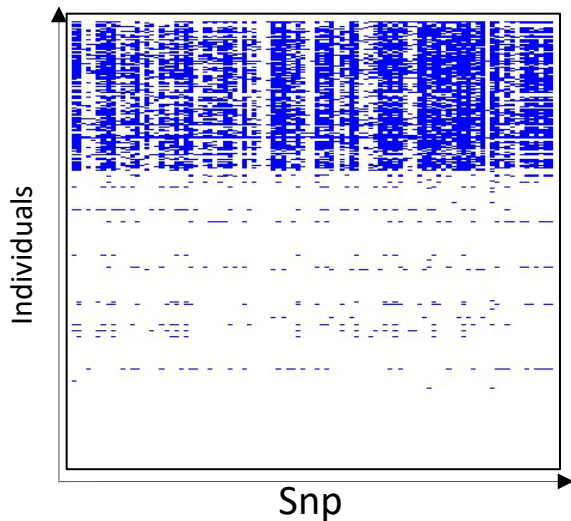
Chunk of a paternal gamete (100 SNP) for 1,614 individuals sequenced according to three different strategies

Unequal distribution of sequencing resources

Key Sires

Focal Individuals

Equal distribution of sequencing



Yield: ■ Correct ■ Wrong Missing