



**HAL**  
open science

## **Genotyping-by-sequencing of pooled drone DNA for the management of living honeybee (*Apis mellifera*) queens in commercial beekeeping operations in New Zealand**

Gertje E. L. Petersen, Peter F. Fennessy, Tracey C. van Stijn, Shannon M. Clarke, Ken G. Dodds, Peter K. Dearden

### ► **To cite this version:**

Gertje E. L. Petersen, Peter F. Fennessy, Tracey C. van Stijn, Shannon M. Clarke, Ken G. Dodds, et al.. Genotyping-by-sequencing of pooled drone DNA for the management of living honeybee (*Apis mellifera*) queens in commercial beekeeping operations in New Zealand. *Apidologie*, 2020, 51 (4), pp.545-556. <10.1007/s13592-020-00741-w>. <hal-03175962>

**HAL Id: hal-03175962**

**<https://hal.science/hal-03175962v1>**

Submitted on 22 Mar 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization



# Genotyping-by-sequencing of pooled drone DNA for the management of living honeybee (*Apis mellifera*) queens in commercial beekeeping operations in New Zealand

Gertje E. L. PETERSEN<sup>1,2</sup>, Peter F. FENNESSY<sup>2</sup>, Tracey C. VAN STIJN<sup>3</sup>,  
Shannon M. CLARKE<sup>3</sup>, Ken G. DODDS<sup>3</sup>, Peter K. DEARDEN<sup>1,4</sup>

<sup>1</sup>Department of Biochemistry, University of Otago, Dunedin, New Zealand

<sup>2</sup>AbacusBio Limited, Box 5585, Dunedin, PO, New Zealand

<sup>3</sup>AgResearch, Invermay Agricultural Centre, Mosgiel, New Zealand

<sup>4</sup>Genomics Aotearoa and Biochemistry Department, University of Otago, Dunedin, New Zealand

Received 10 October 2018 – Revised 29 October 2019 – Accepted 22 January 2020

**Abstract** – The absence of a full pedigree can hinder selective breeding efforts. In honeybees, definitive maternity and especially paternity of queens is difficult to determine, even under managed mating schemes (e.g. using artificial insemination) due to the negative effects of single-drone mating on colony fitness. Here we genotyped 388 living queens from two beekeeping operations using Genotyping-by-Sequencing (GBS). We evaluate two methods to call single-nucleotide polymorphism (SNPs), Tassel 5 and Stacks, for their ability to supply SNPs that can recover known relationships. While Stacks discovered more SNPs (29,433), SNPs called with Tassel 5 (16,757) were found to be more accurate for the derivation of relationships. This methodology presents a low-cost genotyping approach and can be used to support commercial honeybee breeding schemes.

## Breeding / Genotyping-by-sequencing / Genotyping / SNPs / Queen

### 1. INTRODUCTION

Beekeeping is an agricultural sector of growing importance, but genetic improvement in honeybees (*Apis mellifera*) currently lags behind the achievements of other livestock industries. This is due to both the biological characteristics of honeybees, such as their mating habits, and a lack of tools for generating reliable phenotype and genotype data. An additional factor in this is the predominance of beekeeping as a hobby or part-time occupation in many countries, which further complicates the design and implementation of

said tools, and is in contrast with the professional beekeeping industries in countries such as Canada, Australia, and New Zealand.

A successful marketing campaign for native Mānuka honey based on its perceived health benefits has resulted in high prices for honey and related products in New Zealand, leading to a national beekeeping industry that has gone through significant growth and maturation. As a result, it now shows signs of organization that can be found in other modern livestock industries: industry centralization, with the founding of Apiculture NZ in 2016; and the establishment of a defined stud breeder level, with breeders specializing in producing improved “elite” queens. These breeders of elite queens can afford to invest more into both phenotyping and genotyping and have been able to transfer some of the technologies

Corresponding author: G. Petersen,

[GPetersen@abacusbio.co.nz](mailto:GPetersen@abacusbio.co.nz)

Manuscript editor: David Tarpy

used with great success in other species, such as artificial insemination (AI).

While AI allows control over the drones contributing to a mating, it must be used with prudence regarding both the number of drones and the number of drone-producing queens (drone sources) involved. Diversity within the workers of a colony has been shown to directly impact colony fitness and productivity (Mattila and Seeley 2007), nutritional status (Eckholm et al. 2014) as well as pathogen load (Desai and Currie 2015). As a consequence, restrictive mating (using single drones or brothers) should be avoided even in conditions where AI is performed, as colonies arising from restrictive mating are unlikely to perform at the desired level and survival could be compromised. Honeybee breeding is, therefore, not an exception and management of genetic diversity at a population level should be a key element of any honeybee breeding program to maintain a viable population (Bienefeld 2016; Meuwissen 1997). Traditional management strategies for inbreeding avoidance are based on pedigree information. Optimum contribution selection strategies seeking to increase genetic gain while minimizing inbreeding have in fact shown to realize more genetic gain than comparable schemes based purely on genomic information (Henryon et al. 2019). In addition, complete and correct pedigrees are already difficult to obtain in many livestock species, for honeybees this is even more problematic, i.e. when a queen is mated to multiple diverse drones. Genomic information has been recognized as second source of information to manage genetic diversity. Therefore, these problems can be addressed via genotyping of living queens, e.g. using wing clippings or larval exuviae (Chaline et al. 2004; Gregory and Rinderer 2004). However, these non-destructive sources of DNA require implementation of management techniques (timely collection of hatched queen cells, wing clipping) that are not necessarily in line with standard management practices in commercial beekeeping. An alternative affordable standard protocol for the genotyping of queens based on pooled drone samples (Petersen et al. 2017) does not require any manipulation of the queen herself, and could potentially be used across an entire industry.

The focus of this work was to find solutions to three practical issues which are key to the success of a commercial genotyping pipeline, namely (1) that genotyping may be performed by different providers or at different times, requiring a robust SNP calling pipeline that does not show bias between libraries, (2) sampling in the field, and (3) the high costs of DNA extraction and analysis including the significant computational power required for SNP calling. Therefore, we investigated the feasibility of Genotyping-by-Sequencing (GBS) (Elshire et al. 2011) on a set of samples from two beekeeping operations to support a genetic improvement program.

## 2. MATERIALS AND METHODS

### 2.1. Sampling and DNA extraction

Honeybee samples were collected during two New Zealand summers, 2015/16 and 2016/17, in January 2016 and November 2016, respectively, from Taylor Pass Honey Co Ltd. (TPH), a commercial beekeeping operation in Marlborough, New Zealand, and Bettabees Research Ltd. (BBRL), a queen breeder in Mosgiel, New Zealand. BBRL routinely supplies TPH with mated breeder queens (Italian type) that have been artificially inseminated. Both beekeeping operations are part of an ongoing research program.

As a previous trial on pools of adult drones had identified issues with mis-sampling of adult drones across colonies (Petersen et al. 2017), the sampling protocol was modified to drone larvae for the second round of sampling. A total of 196 sets of adult drones were collected from TPH, 19 sets of adult drones from BBR, and 173 sets of drone larvae from TPH, for a total of 388 sample sets (see Table 1).

Adult drones were collected and put on ice until transfer to a  $-20^{\circ}\text{C}$  freezer; drone larvae were placed in a portable  $-20^{\circ}\text{C}$  freezer immediately after removal in the field. At arrival in the laboratory, drones were transferred into a  $-80^{\circ}\text{C}$  freezer until required for DNA extraction.

Material from 10 drones/drone larvae from each hive was pooled prior to DNA extraction. Genomic DNA (gDNA) was extracted from pooled heads of adult drones and pooled larval

**Table 1.** Overview of sample sets of adult drones or drone larvae. Sample sets of haploid drones were pooled to create a representation of their diploid mother queen.

| Beekeeping operation | Sampling time point | Sets of adult drones | Sets of drone larvae |
|----------------------|---------------------|----------------------|----------------------|
| Taylor Pass Honey Co | Early 2016          | 196                  | –                    |
| BettaBees Research   | Early 2016          | 19                   | –                    |
| Taylor Pass Honey Co | Late 2016           | –                    | 173                  |
|                      |                     | 205                  | 173                  |

material (avoiding the gut). The “Quick-DNA Tissue/Insect 96 Kit” (Zymo Research, Irvine, CA, USA) was used for all extractions, but the standard protocol was modified to accommodate the extraction from pooled larva samples. For the pools of adult drone heads, a 2010 GenoGrinder® (SPEX® SamplePrep, Metuchen, NJ, USA) was used at 1200 rpm for 8 min (note that full sample lysis was not achieved). The resulting DNA yields were measured by Qubit® dsDNA HS Assay Kit (ThermoFisher Scientific, Waltham, MA, USA). For the larval samples, the BashingBeads tissue lysis step was replaced with the following protocol: around 2 mg material each from 6 drone larvae was placed in a 1.5 ml centrifugation tube with a 3.2 mm stainless steel ball bearing and 1 ml of lysis buffer and homogenized in two 30 s steps in a QIAGEN® TissueLyser II (QIAGEN, Venlo, NL). 100 µl of the lysate were then transferred into the second step of the Quick-DNATM 96 kit.

## 2.2. Genotyping-by-sequencing

GBS (Elshire et al. 2011) is one of several genotyping methods using restriction enzymes to produce a reduced representation of the genome. For this study, a double digest with restriction enzymes ApeK1 and Msp1 was carried out on 388 samples. An optimized version of the original GBS protocol was followed that is routinely used for industry-wide GBS studies (Dodds et al. 2015). Using 100 ng of DNA per sample, 3 indexed GBS libraries were prepared and purified using the Pippin Prep (SAGE Science, Beverly, MA, USA) for size selection to generate DNA sequencing libraries with 150–500 bp fragment length. On an Illumina HiSeq2500, single-end

sequencing ( $1 \times 100$ ) was performed, resulting in around 27 Gb of raw sequence data per lane. Raw fastq files were quality checked using FastQC v0.10.1 (Andrews 2010) before further analysis.

## 2.3. Data analysis, SNP calling, and alignment

The sequencing data were analyzed in separate batches as sequencing was performed at different times (see Table 2). Two different GBS data processing pipelines, TASSEL 5 GBSv2 (Glaubitz et al. 2014) and Stacks (Catchen et al. 2013) were used to investigate the resulting SNP numbers as well as computational demands of both pipelines. TASSEL was run with default settings except for the minimum kmer count to be included as tag, which was set to 4. Stacks was run with default settings, which include the minimum number of reads to be included as a stack (the equivalent to a ‘tag’ in TASSEL) at 3, and the number of mismatches allowed between stacks to be merged at 2. In accordance with the method described by Paris et al. (Paris et al. 2017), stacks were built de novo and aligned afterwards (instead of aligning reads and then running the *refmap.pl* wrapper).

In both cases, short reads were aligned to the *A. mellifera* reference genome (Weinstock et al. 2006) using the Burrows-Wheeler Aligner (Li and Durbin 2010) and further processed by constructing a relationship matrix using the KGD R software (Dodds et al. 2015) after merging the VCF files containing SNP calls for each batch into one file. This relationship matrix is based on KGD unbiased estimates of relatedness, calculated via method 1 of VanRaden (VanRaden 2008) adjusted to

**Table II.** Batches of GBS data analyzed using Stacks and TASSEL 5 correspond to sequencing libraries. Since the mean read depth in 1 had been found to be high (50.9), subsequent sequencing runs were extended to include 2 96-well libraries each with 1 negative control per plate, resulting in 190 samples per sequencing run.

| Sequencing run ID | Beekeeping season | No of samples (no included in the analysis)        |
|-------------------|-------------------|--|
| 1*                | 2015/16           | 78 (8)   |
| 2                 | 2015/16 & 2016/17 | 190 (190)  |
| 3                 | 2016/17           | 190 (190)  |
|                   |                   | <i>Total samples included in the analysis: 388</i> |

\*Sequencing run 1 was the pilot project described in Petersen et al. 2017, which lead to a much higher mean read depth for samples originating from this run as it was limited to 78 samples in one lane of sequencing (compared with 190 samples per lane for sequencing runs 2 and 3). Only one replicate of the 8 queens that contributed drones and their own DNA were included in this work

account for read depth at each locus including SNPs with missing reads.

Within KGD, both clustering to construct a heatmap and a Principal Component Analysis (PCA) are performed as part of the standard analysis, which vaguely suggested the presence of several clusters (see Fig. 1).

Further clustering was based on the relationship matrix called ‘G5’ as described by Dodds et al. (2015) and based on kinship using GBS with depth adjustment. Clustering was performed with the R package dendextend (Galili 2015) to investigate the potential structure of the population and to find out if any relationships between breeder queens and production queens (their putative daughters) could be recovered.

### 3. RESULTS AND DISCUSSION

#### 3.1. DNA extraction from adult drones and drone larvae

A trial investigating the use of GBS of drone pools as proxies for their queen mothers had previously shown promising results. In that work, most of the variation in how well a particular drone pool reflected their mother’s genome was traced back to problems surrounding drone sampling as some adult drones proved not to be sons of the resident queen (Petersen et al. 2017). After consultation with beekeepers, the sampling protocol was changed to sampling of drone larvae. To avoid accidental sampling of diploid drones, a

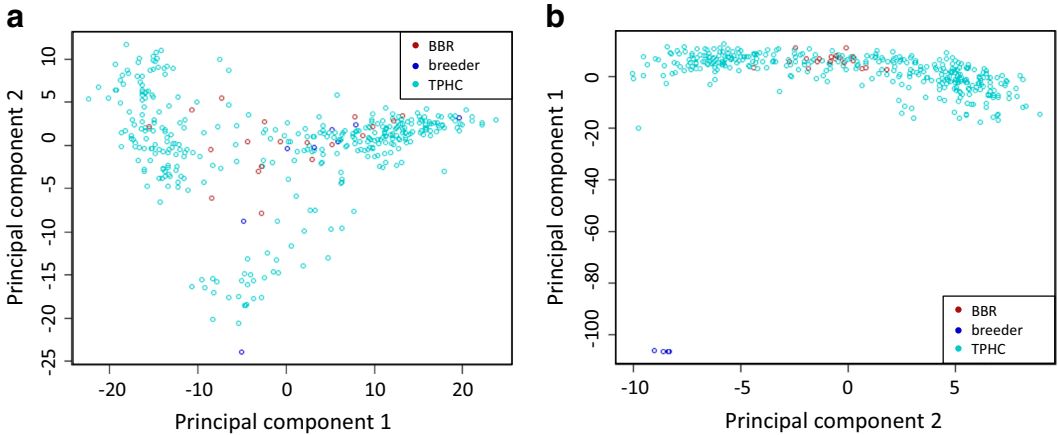
standard operating procedure was established, stating that drone larvae had to be sampled from drone cells around the edges of the brood nest.

gDNA yields for extraction from adult drones were generally lower than from pooled larvae and ranged from 10 µg/ml to 110 µg/ml. However, since adult drone samples were handled on an individual basis, this was to be expected, and genotype quality did not differ between adult drones and larvae. With gDNA yields from 87.6 µg/ml to >120 µg/ml (the upper limit of the Qubit® dsDNA HS Assay Kit), yields from pooled larvae were both higher and more consistent, in addition to getting around the mis-sampling problems.

#### 3.2. SNP calling using TASSEL and stacks

SNP calling with TASSEL 5 resulted in 21,951 SNPs that could be called across all three of the separate batches and processed further in R. The proportion of missing genotypes was 9.39%, while the mean sample depth was 14.84. Applying a – 0.05 Hardy-Weinberg disequilibrium cut-off resulted in 16,757 SNPs (proportion of missing genotypes 12.22%, mean sample depth 10.85) that were used to calculate the KGD relatedness measures and construct a relationship matrix.

SNP calling with Stacks resulted in a higher number of SNPs (29,433) and mean sample depth (33.04), but also a substantially higher rate of missing genotypes (43.87%). Applying a – 0.05 Hardy-Weinberg disequilibrium cut-off resulted



**Figure 1.** Principal component analysis of KGD relatedness measures for 388 honeybee queen from two beekeeping operations based on SNPs called with TASSEL 5 (a) and Stacks (b). While there is clustering present among the 388 queens in the sample, this cannot be explained by the source population of the samples (BBR = BettaBees Research, breeder = Taylor Pass Honey Co breeder queens, TPHC = Taylor Pass Honey Co production queens)

in 28,882 SNPs for the calculation of the KGD relationship matrix, with 43.32% missing genotypes and a mean sample depth of 29.07.

In addition to the differences in output, TASSEL 5 and Stacks differed markedly in terms of hardware requirements and processing time. While TASSEL 5 could be run on a desktop computer with 8 cores and 128GB RAM, Stacks required a much higher number of cores and had to be run on a dedicated bioinformatics server, accessing 20–30 cores throughout processing. Even though, the full TASSEL 5 pipeline took about 26 h (~11 h per full batch (as batch 1 was smaller)) and the Stacks pipeline took up to a week per full batch.

It should be noted that both of the SNP-calling pipelines offer room for optimisation of the parameter space as per Paris et al. (2017), which was not fully explored since the objective was to find a protocol that would need minimal set-up time (including deviation from default parameters) and optimisation of Stacks would have required additional weeks of computing time.

### 3.3. Relationships

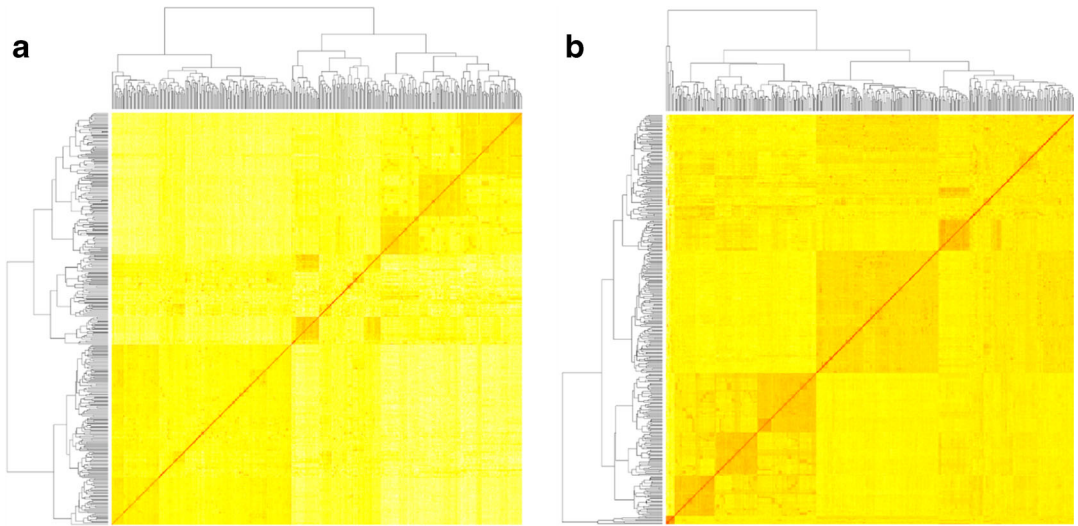
Some queens sampled within the TPHC population had known or suspected relationships with others in the sample, as noted below:

- The first samples taken in early 2016 came from breeder queens present in the TPHC queen unit at the time, and thus were expected to be (some of) the mothers of the early 2016 production queen sample.
- A group of sisters were sampled in the TPHC queen multiplying unit in early 2016, with their mother included in a).
- The late 2016 sample consisted of daughters of three breeder queens present in the queen unit by late 2016 (which were unfortunately not identical with the ones in a)); in addition, the production queens were suspected to be descendants of sisters of queens sampled in the BBR population.

The relationships discovered in the downstream analysis were different based on SNPs called using Tassel and Stacks, and not all known relationships could be recovered using both approaches.

The relationships discovered in the downstream analysis differed based on SNPs called using TASSEL 5 and Stacks, and not all known or suspected relationships could be recovered using both approaches.

Heatmaps and PCA generated as part of the routine KGD output suggested that there was underlying structure within the population. Both



**Figure 2.** Heatmaps based on KGD relatedness measures for 388 honeybee queens from 2 beekeeping operations. The heatmaps suggests the presence of structure within the sampled population

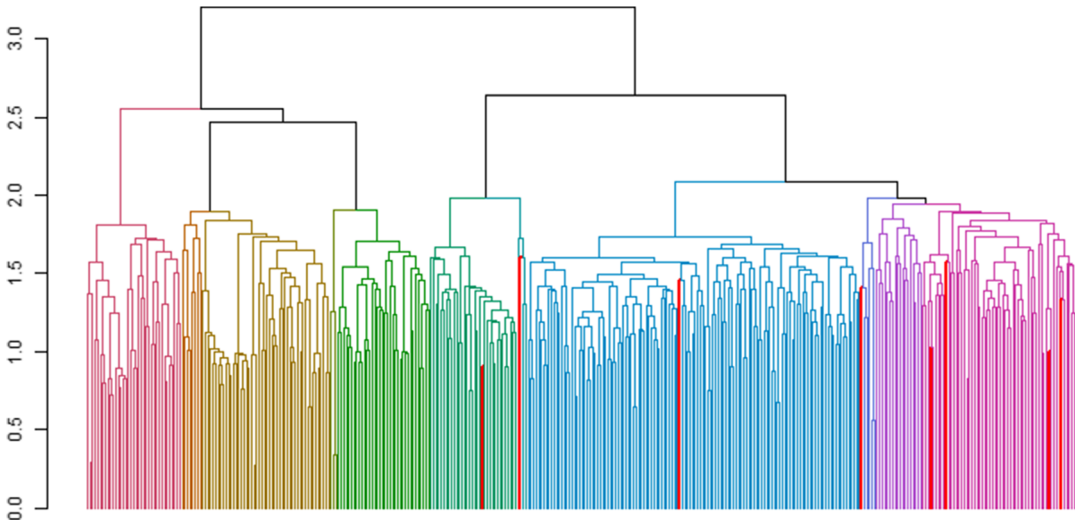
the heatmap (Fig. 2B) and the PCA based on the Stacks output implied that there was an additional small cluster which was only loosely related to the rest of the population, this was not found to be true after SNP calling with TASSEL 5 (Fig. 1A, Fig. 2A). To investigate this phenomenon, the KGD G5 relationship matrix was clustered further and the clusters visualized in dendrograms. Based on the SNPs called with TASSEL 5, the population was found to show a marked difference between seasons (2015/16 vs. 2016/17), but not between populations, with the 8 breeder queens from batch 1 clustering across different clusters (Fig. 3). This was expected due to the continuous flow of genetics from one population to the other.

Visual inspection of the dendrogram suggested the presence of seven distinct clusters, with the queens collected from BBR falling into two clusters (magenta, purple in Fig. 3) that also included four of the eight breeder queens from TPHC. This was expected due to the continuous flow of genetics from one population to the other. Simultaneously, the queens that were known to be sisters clustered together, with the late 2016 queens forming the expected three clusters (red, tan, and green in Fig. 3) and the early 2016 mating yard forming another cluster (turquoise).

However, the relationships calculated based on the Stacks output did not reflect these known connections. Seven similar clusters could be recovered for the production queens sampled in both early and late 2016 (co-assignment rate to clusters ranged from ~60% to ~90%, depending on the cluster), but the breeder queens showed a very distinct separation from the rest of the queens (Fig. 4).

The eight breeder queens in early 2016 had been part of a pilot project that included 78 samples consisting of queen replicates based on drone pools and the corresponding queen, whereas the rest of the samples had been part of libraries of 180 samples each. This subset had been deliberately included in the data set to emulate a real-life situation in dealing with commercial samples. The different library size resulted in higher read depth per sample for the breeder queens. This was found to be a lot more pronounced with Stacks than with TASSEL 5 (Fig. 5). A similar pattern was found for the callrate (proportion of non-missing genotypes; Fig. 6), suggesting that Stacks had called a faction of SNPs exclusively in the high-depth breeder queens.

As KGD omits SNPs which do not appear in both individuals for each pairwise comparison, it



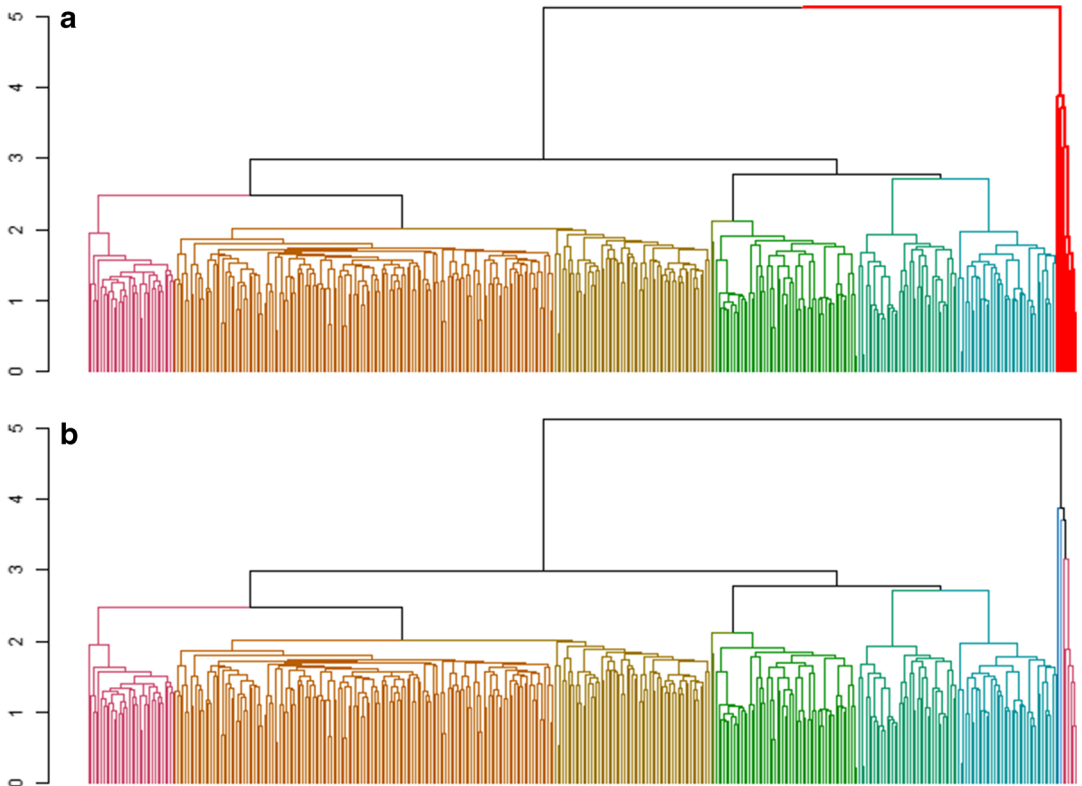
**Figure 3.** Dendrogram of relationships within a KGD relationship matrix of 388 honeybee queens calculated based on SNPs called with TASSEL5. The first branch marks the differentiation between the samples collected in early 2016 and late 2017. Queens from the queen breeder, BBR, all fall into the two magenta, and purple clusters. 8 breeder queens (in red, bold) had been sequenced to a much higher depth in a pilot project with fewer samples, but clustered across different clusters within the other queens

does not introduce imputation bias. However, the presence of a large proportion of missing genotypes in a specific subset of individuals, as found here, skews the depiction of relationships in downstream analysis. Stacks found a high proportion of SNPs in most of the breeder queens, but not in the rest of the clusters, resulting in a separate “breeder queen cluster” (see Fig. 7B). This does not reflect the suspected relationships between breeder queens and production queens which should have led to the breeder queen clustering with groups of production queens that are their daughters or sisters.

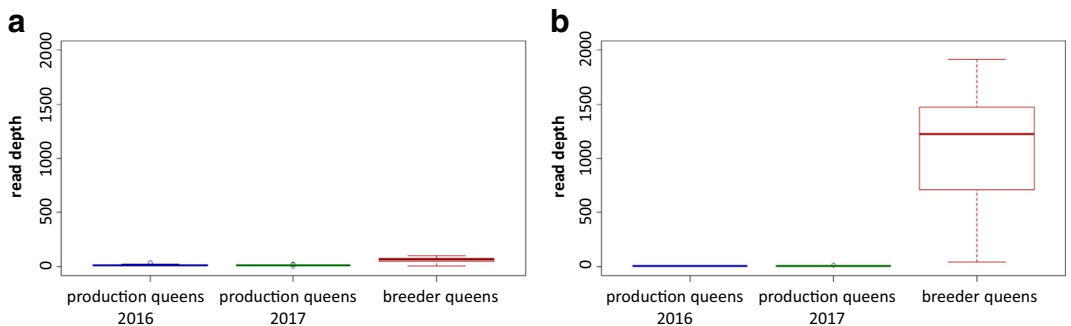
Although suspected relationships could be confirmed with the use clustering tools on the basis the results of TASSEL 5 and KGD, the resulting relatedness measures did not adequately reflect expected relationships. Due to the honeybee sex determination mechanism and the resulting haplodiploid population, the expected relationships of related honeybee queens can range from 0.25 to 0.75 as depicted in Table 3.

In this study, the relationships between queens based on SNP calling with TASSEL 5 were found to range from  $-0.15$  to  $1.15$ . In

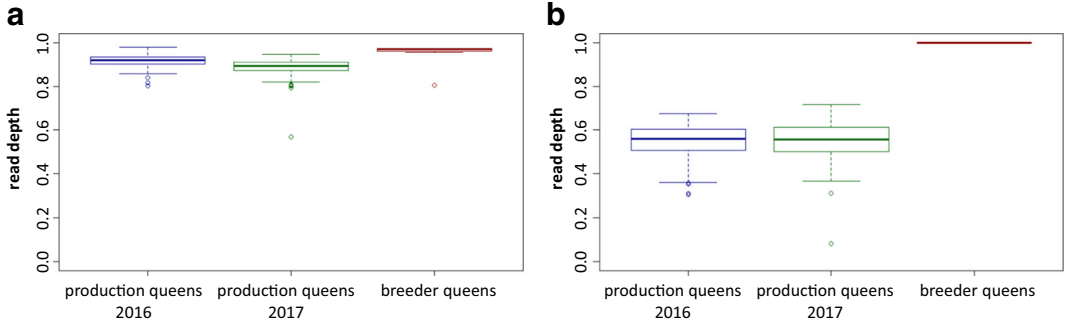
a pilot study regarding GBS in honeybees, the mean relatedness between replicate genotyped samples based on aliquots of the same DNA extract was found to be 1.01, suggesting that (at least within the same batch), GBS is consistent. However, the relatedness measures between replicates ranged from 0.83 to 1.12, indicating that there is a certain level of uncertainty attached to the relatedness measures based on GBS data. While a 12–17% loss in accuracy in the assignment of relationships is less problematic in species with a basic relatedness pattern between of unrelated = 0, sibling or parent = 0.5, self, or twin = 1, it creates some issues with honeybee relationships where 17% uncertainty can be enough to suggest a completely different relationship between sisters (e.g. suggesting that they are super sisters when they are full sisters or the other way around). Based on these results it seems unlikely that GBS will yield SNPs with a sufficient resolution to allow calculation of an accurate genomic relationship matrix that could serve as full substitution of a pedigree in genomic selection.



**Figure 4.** Dendrogram of relationships within a KGD relationship matrix of 388 honeybee queens calculated based on SNPs called with Stacks. The first branch marks the differentiation between the breeder queens and production queens. The breeder queens (in red, bold (A)) had been sequenced to a much higher depth in a pilot project with fewer samples and clustered together outside the bulk to the other individuals as a result



**Figure 5.** Read depth per genotyped sample group (2016 production queens, 2017 production queen, breeder queens) in GBS data from 388 honeybee queens as found using two different SNP calling pipelines, TASSEL 5 (a) and Stacks (b). The mean read depth of the breeder queen samples sequenced in a much smaller sequencing run is higher with both SNP calling pipelines, but this is much more pronounced with Stacks (mean read depth for breeder queens = 1400) than with TASSEL 5 (mean read depth for breeder queens = 75)



**Figure 6.** Callrate (proportion of non-missing genotypes) per genotyped sample group (2016 production queens, 2017 production queen, breeder queens) in GBS data from 388 honeybee queens as found using two different SNP calling pipelines, TASSEL 5 (a) and Stacks (b). The mean callrate of the breeder queen samples sequenced in a much smaller sequencing run is higher with both SNP calling pipelines, but this is much more pronounced with Stacks than with TASSEL 5

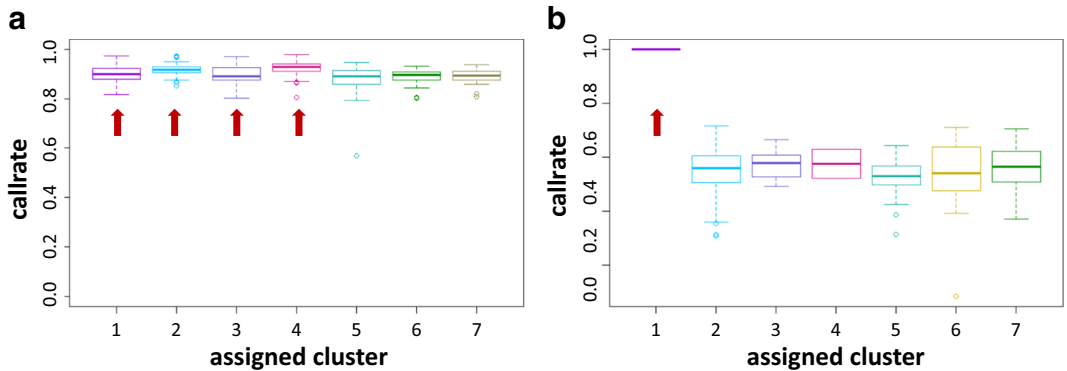
**3.4. Implications**

Practical issues in implementation of genotyping to service a breeding program include:

- 1) Fluctuations in genotyping conditions, from changing sample numbers resulting in different library sizes, to different providers.
- 2) sampling in the field.
- 3) the high costs of DNA extraction and data analysis (including costs of computation).

In respect of the bias aspect across libraries, TASSEL 5 and KGD performed better than the Stacks and KGD in recovering suspected relationships, even though Stacks discovered significantly more SNPs. This indicates that SNP calling with TASSEL 5 is a robust way to analyse data from batches with differing sample sizes for genotyping (batch size in this study varied from 78 samples/lane to 180 samples/lane) that can be expected at different times of the year.

The sampling issue observed when adult drones were collected (as described in Petersen



**Figure 7.** Callrate (proportion of non-missing genotypes) per assigned cluster (= 7) in GBS data from 388 honeybee queens as found using two different SNP calling pipelines, TASSEL 5 (a) and Stacks (b). *k* was assigned after visual inspection of a dendrogram depicting clustered measures of queen relatedness. A small number of breeder queens had been sequenced in a smaller batch at higher depth. While with TASSEL 5 these breeder queens could be assigned to four different clusters (see red arrows in A), with Stacks this failed, and all breeder queens were assigned to their own cluster (red arrow in B)

**Table III.** Relationships between a honeybee queen and her relatives. Since drones have clonal sperm and arise from unfertilised eggs, they cannot be considered sires in an animal breeding context (while still biological fathers). Super sisters are the result of the same mating (Dam x drone), full sisters are the result of two matings with brother drones, half sisters are the result matings with two unrelated drones.

| To    | Daughter | Son | Dam | Father | Super sister | Full sister | Half sister |
|-------|----------|-----|-----|--------|--------------|-------------|-------------|
| Queen | 0.5      | 0.5 | 0.5 | 0.5    | 0.75         | 0.5         | 0.25        |

et al. 2017) should be satisfactorily resolved through drone larval sampling which guarantees sampling integrity since these must be the progeny of the resident queen. In addition, the sampling of larvae was found to be easier to perform in the field, with lower risks of stings and less disturbance of the sampled colony.

The cost of extraction and analysis remains high, at around \$US18 per sample. Given that full management of a breeding program will also require determination of the allele status in the sex-determination locus *csd* in selection candidate queens, approaches that reduce the cost of genotyping will continue to be important. The cost of GBS data analysis using Stacks was found to be significantly higher than with the use of TASSEL 5. The full SNP calling pipeline using TASSEL 5 and BWA could be performed on a machine with 8 cores and around 100 GB of RAM in just over a day, whereas Stacks required more powerful hardware as well as significantly more time, with a single batch being analyzed in around 6 days. Since costs of computation can be substantial where thousands of individuals have to be analyzed, TASSEL 5 was found to be preferable over Stacks in this regard.

#### 4. CONCLUSION

Pools of drone larvae provide a reliable source of high-quality DNA for the genotyping of their mothers. While the protocol presented here still includes an individual sample preparation component, most of the process can be automated, opening the possibility of establishing a high-throughput system for a commercial genotyping pipeline.

GBS followed by Tassel 5 for SNP calling and KGD R software for derivation of relationships appears to provide an effective way to

construct relationship matrices for honeybee populations. Although these relationship matrices are unlikely to fully reflect the real relationships due to lack of resolution, this enables the assignment of queens to lines or families even in the absence of sophisticated tools to manage mating such as AI, and can serve to verify a suggested pedigree or determine parents out of a pool of candidates. Determination of *csd* alleles alongside this type of genotype will also be required to manage brood viability issues (Hyink et al. 2013). Unfortunately, due to the high levels of diversity in the *csd* gene and the size of the hypervariable regions that will need to be sequenced to allow *csd*-typing, targeted GBS (inclusion of a *csd* amplicon into the GBS protocol) that has been developed for other applications (e.g. Caulton 2018) will not be possible due to the restriction to short sequencing reads. A separate protocol to determine the *csd* alleles carried by individual queens will be necessary.

However, information on the relatedness of individual queens does provides the opportunity to gain more control over the management of genetic diversity within honeybee populations. It might also allow queen breeders generate rough estimates of genetic merit for individual queens that have no known (i.e. recorded pedigree) connection to previously evaluated breeder queens. It would be reasonable to expect that the rate of genetic gain can be increased beyond that currently possible using the described genotyping methodology, because valuable queens (i.e. that express desired phenotypes) can be brought into the evaluation and selection without compromising the integrity of the pedigree. Genetic improvement could thus be possibly extended across separate populations in the context of an industry-wide genetic improvement program.

## ACKNOWLEDGMENTS

GP performed this work as a Callaghan Innovation PhD fellow in the Department of Biochemistry at the University of Otago, in cooperation with AbacusBio Limited and Taylor Pass Honey Company, Blenheim, New Zealand, with additional funding being provided by Taylor Pass Honey Co. This work was carried out in collaboration with the ‘Genomics for Production & Security in a Biological Economy’ C10X1306 New Zealand Ministry of Business, Innovation and Employment program grant to AgResearch and ‘Selecting Future Bees’ UOOX1610 New Zealand Ministry of Business, Innovation and Employment program grant to PKD.

PKD, PF, and GP conceived this research and designed experiments; GP collected samples and performed DNA extractions and data analysis; SC and TVS performed GBS. GP wrote the paper, PKD, PF, KGD, and SC participated in the revisions of it. All authors read and approved the final manuscript.

The authors declare that they have no conflict of interest in relation to the study in this paper.

The sequence data generated during this study are available in the NCBI Sequence Read Archive, as part of the FutureBees NZ BioProject, Accession number PRJNA490757.

**Génotypage par séquençage de l’ADN de mâle mis en commun pour la gestion des reines d’abeilles vivantes (*Apis mellifera*) dans les exploitations apicoles commerciales en Nouvelle-Zélande.**

**élevage / génotypage par séquençage / génotypage / SNPs / reine.**

**Genotyping-by-Sequencing (GBS) von gepoolter Drohnen-DNA für das Management lebender Honigbienenköniginnen in Berufsimkereien in Neuseeland.**

**Zucht / Genotyping-by-Sequencing (GBS) / SNPs / Königin.**

## REFERENCES

Andrews, S. (2010). FastQC: a quality control tool for high throughput sequence data. Available online at: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>

- Bienefeld K. (2016) Breeding Success or Genetic Diversity in Honey Bees? *Bee World* **93** (2), 40–44
- Catchen J.M., Hohenlohe P., Bassham S., Amores A., Cresko W. (2013) Stacks: an analysis tool set for population genomics. *Mol. Ecol.* **22** (11), 3124–3140
- Caulton, A. J. (2018). “Development of ‘Targeted’ Genotyping-by-Sequencing in Atlantic Salmon (*Salmo salar*).” PhD thesis. University of Otago.
- Chaline N., Ratnieks F.L.W., Raine N.E., Badcock N.S. (2004) Non-lethal sampling of honey bee, *Apis mellifera*, DNA using wing tips. *Apidologie* **35** (3), 311–318
- Desai S.D., Currie R.W. (2015) Genetic diversity within honey bee colonies affects pathogen load and relative virus levels in honey bees, *Apis mellifera* L. *Behav. Ecol. Sociobiol.* **69** (9), 1527–1541
- Dodds K.G., McEwan J.C., Brauning R., Anderson R.M., van Stijn T.C., Kristjansson T., Clarke S. (2015) Construction of relatedness matrices using genotyping-by-sequencing data. *BMC Genomics* **16** (1047), 1–15
- Eckholm B.J., Huang M.H., Anderson K.E., Mott B.M., DeGrandi-Hoffman G. (2014) Honey bee (*Apis mellifera*) intracolony genetic diversity influences worker nutritional status. *Apidologie* **46** (2), 150–163
- Elshire R.J., Glaubitz J.C., Sun Q., Poland J.A., Kawamoto K., Buckler E.S., Mitchell S.E. (2011) A robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS ONE* **6** (5), e19379
- Galili T. (2015) dendextend: an R package for visualizing, adjusting and comparing trees of hierarchical clustering. *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/btv428>
- Glaubitz J.C., Casstevens T.M., Lu F., Harriman J., Elshire R., Sun Q., Buckler E.S. (2014) TASSEL-GBS: A high capacity genotyping by sequencing analysis pipeline. *PLoS ONE* **9** (2), e90346
- Gregory P.G., Rinderer T.E. (2004) Non-Destructive Sources of DNA Used to Genotype Honey Bee (*Apis Mellifera*) Queens. *Entomologia Experimentalis et Applicata* **111** (3), 173–177
- Henryon M., Liu H., Berg P., Su G., Nielsen H.M., Gebregiwergis G.T., Sorensen A.C. (2019) Pedigree Relationships to Control Inbreeding in Optimum-Contribution Selection Realise More Genetic Gain than Genomic Relationships. *Genetics Selection Evolution* **51** (1). <https://doi.org/10.1186/s12711-019-0475-5>
- Hyink O., Laas F., Dearden P.K. (2013) Genetic Tests for Alleles of Complementary-Sex-Determiner to Support Honeybee Breeding Programmes. *Apidologie* **44** (3), 306–313
- Li H., Durbin R. (2010) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/btp324>
- Mattila H.R., Seeley T.D. (2007) Genetic Diversity in Honey Bee Colonies Enhances Productivity and Fitness. *Science* **317** (5836), 362–364
- Meuwissen T.H.E. (1997) Maximising the response of selection with a predefined rate of inbreeding. *J. Anim. Sci.* **75**, 934–940

- Paris, J. R., Stevens, J. R. and Catchen, J. M. (2017) Lost in Parameter Space: A Road
- Petersen, G. E. L., Fennessy, P. F., van Stijn, T. C., Clarke, S. M. and Dearden, P. K. (2017) Genotyping-by-Sequencing for Genetic Improvement in Honeybees. In Proceedings of the 21st Conference of the Association for the Advancement of Animal Breeding and Genetics, Townsville, Australia: The Association for the Advancement of Animal Breeding and Genetics.
- VanRaden P.M. (2008) Efficient methods to compute genomic predictions. *J. Dairy Sci.* **91**(11), 4414–4423
- Weinstock G.M. et al (2006) Insights into social insects from the genome of the honeybee *Apis mellifera*. *Nature* **443**(7114), 931–949

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.