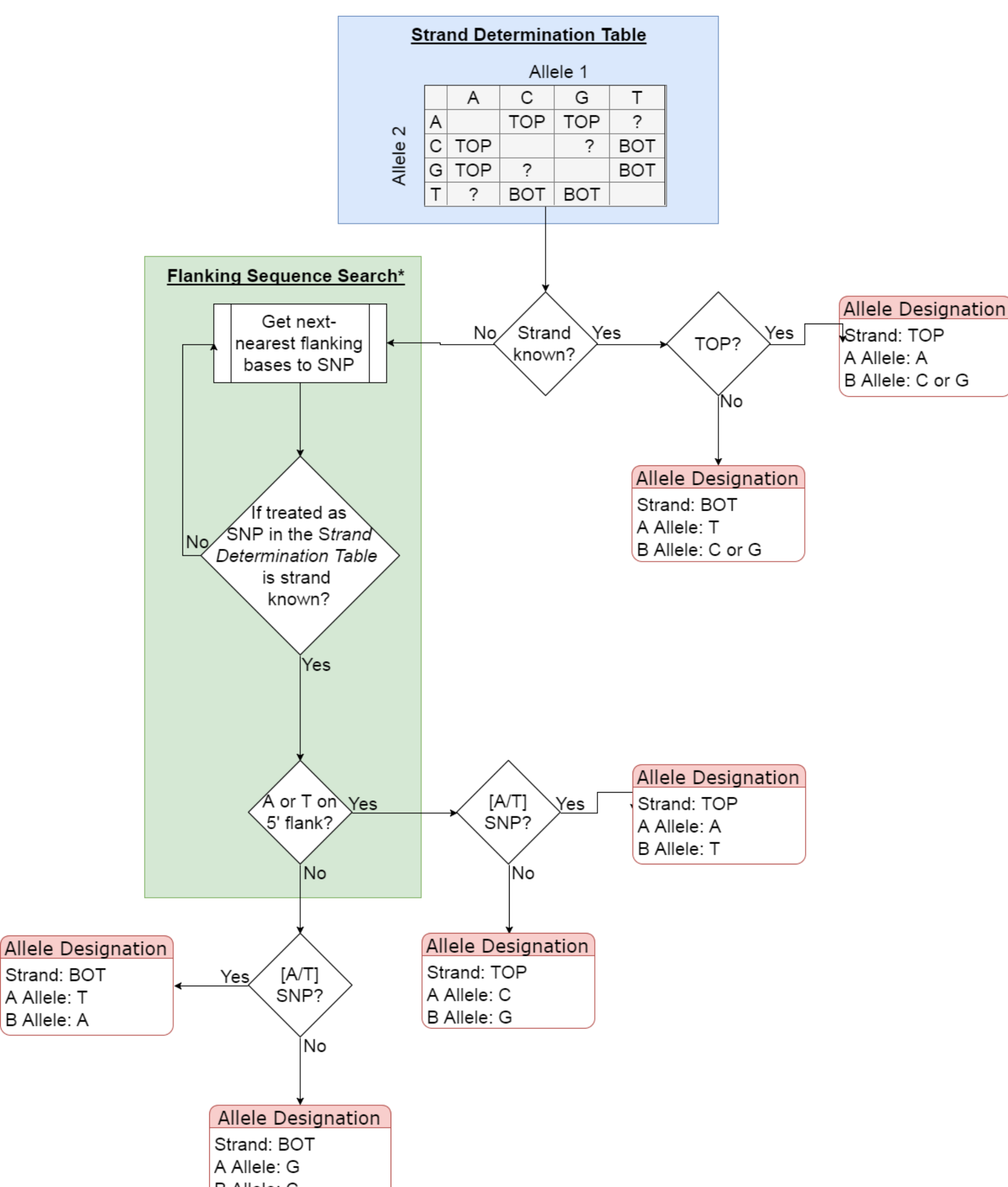# New output formats for Axiom genotyping arrays

## Ali Pirani, Joseph M Foster, Alessandro Davassi, Brant Wong, Mohini A Patil, and Luis Jevons
## Thermo Fisher Scientific, 3450 Central Expressway, Santa Clara, CA 95051

## INTRODUCTION

The high throughput agricultural genotyping landscape encompasses a broad range of applications and technical platforms. One of the major challenges of adopting a new platform or performing meta-analyses is data format congruity. Biallelic genotypes are recorded in one of three ways; "AA", "AB" and "BB" call codes, "0", "1", and "2" numeric call codes and base calls "A", "T", "G" or "C". For call codes and numeric call codes, the A and B alleles must be designated. Historically, two formats have dominated the designation of variant alleles; "Forward" and "TOP". For bi-allelic SNPs this can create a situation where the "A" allele designated by one format differs from the other.

### Figure 1. TOP/BOT format allele designation



The flow diagram in *Figure 1* describes the process for allele designation of bi-allelic SNPs for the TOP/BOT format. Initially strand determination is done by the *Strand Determination Table*. For [A/C] and [A/G] SNPs the strand is defined as TOP, for [T/C] and [T/G] SNPs the strand is defined as BOT. Where strand determination is possible in this manner, allele designation follows such that any A or T base is considered the A allele and the other base that constitutes the SNP is designated as the B allele. For [A/T] and [G/C] SNPs the strand is unknown and the flanking sequence is used to determine strand by the *Flanking Sequence Search* process (more detailed view in *Figure 2*). For SNPs determined to be on the TOP strand, the A or C base is designated as the A allele and the T or G base as the B allele for [A/T] and [G/C] SNPs respectively. For the SNPs determined to be on the BOT strand, the reverse is true; the T or G base is designated as the A allele and the A or C base as the B allele for [A/T] and [G/C] SNPs respectively.

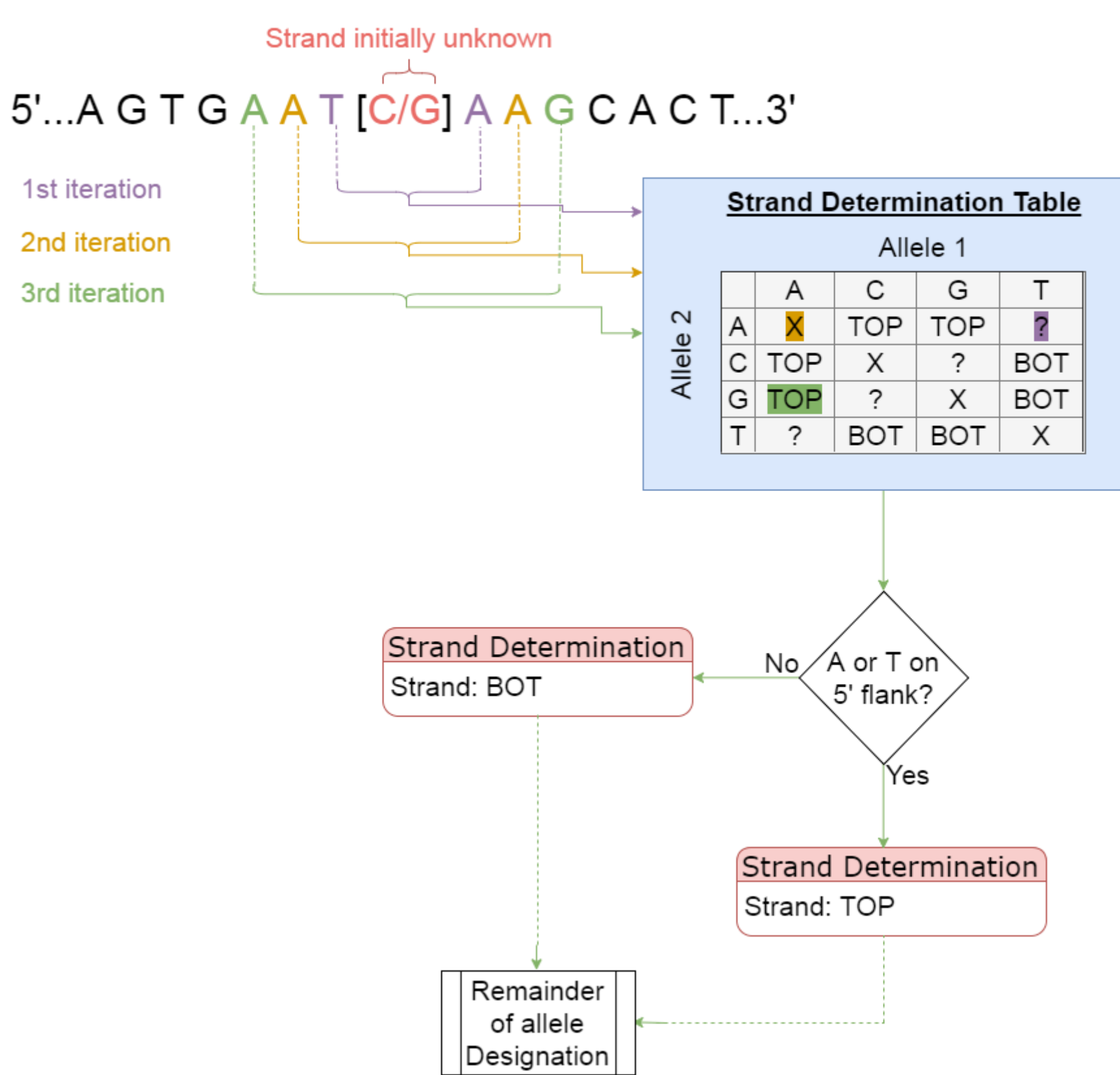### Figure 2. TOP/BOT strand determination for [A/T] and [G/C] SNPs



*Figure 2* shows the strand determination for [A/T] and [G/C] SNPs. Starting with the nearest pair of bases either side of the SNP, the *Strand Determination Table* is checked to see if a strand determination can be made. For each time a strand determination cannot be made the algorithm increments 1 position further away from the SNP in both directions and re-checks the strand determination table until a strand determination would be made. Once this iterative process is complete, if the A or T base is in the 5' flanking sequence the strand is determined as TOP. Conversely, if the A or T allele is in the 3' flanking sequence the strand is determined to be BOT.

## AXIOM LONG FORMAT EXPORT TOOL (AxLE)

To support cross-platform high throughput genotyping analysis, Thermo Fisher Scientific has developed the Axiom Long format Export Tool (AxLE)[1]; a companion application to Axiom™ Analysis Suite[2]. The tool converts Axiom genotype data from native "Forward" format to the "TOP" format based on the polymorphism itself, or the contextual surrounding sequence and designates the A/B allele. The tool also converts the standard Axiom output into a format that is similar to the long format options for other platforms. This makes Axiom genotyping easier to integrate with existing downstream analysis pipelines and large scale meta-analysis of several cross-platform datasets.
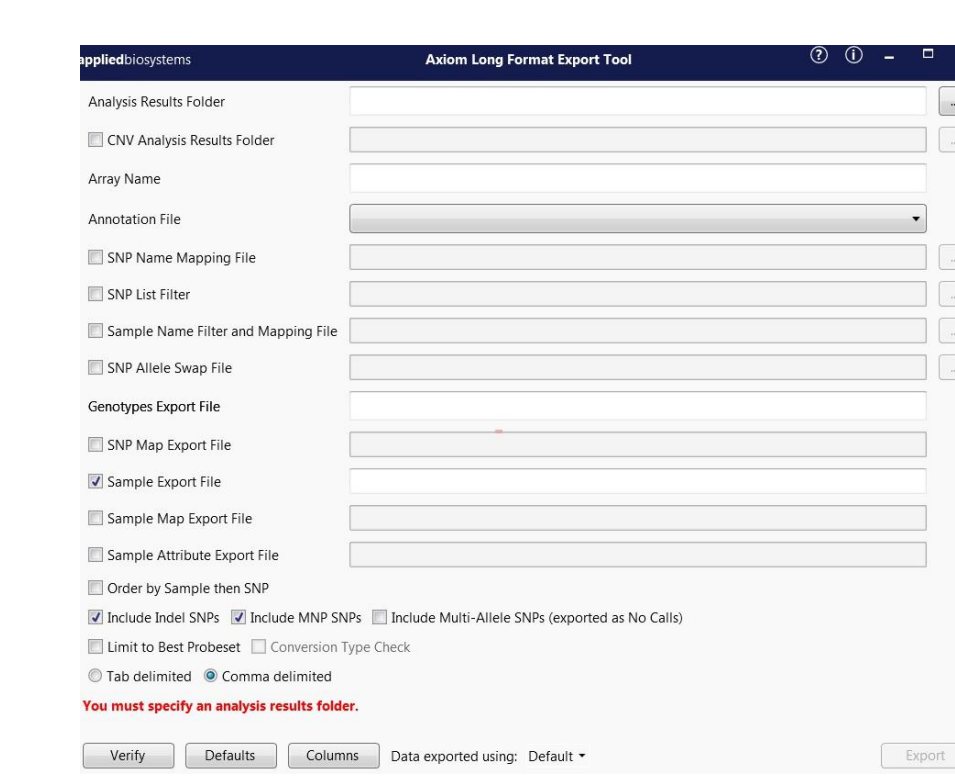
### Figure 3. AxLE Tool Usage



*Figure 3* shows a screen shot of the AxLE tool. After executing a Best Practices Workflow in Axiom Analysis Suite, the AxLE tool can be accessed via the "External Tools" menu of Axiom Analysis Suite. Some of the following steps are taken to generate Long format:

1. "Analysis Results": Select the appropriate results folder of the analysis to be converted to Long format.
2. The "Array Name" and "Annotation File" will automatically be populated.
3. Optionally a "SNP Name Mapping File" can be selected. If selected it will replace the native probeset_ids with alternate SNP identifiers in the Long format output.
4. Optionally the analysis results can be filtered by a "SNP List" to return only those SNPs in the Long format output.
5. A name must assigned to the Long format output file.
6. Review the other options and click the "Export" button.

AxLE output consists of:
- A descriptive header
- SNP Name: the default is probeset_id e.g. AX-123456789, but these can be modified to a user-defined value by using a "SNP Name Mapping File"
- Sample ID: the default is the CEL file name
- Allele 1/2 – Forward: Base call relative to Forward Strand
- Allele 1/2 – Top: Base call normalized to TOP strand
- Allele 1/2 - A/B: Axiom designated A/B allele call
- Confidence: AxiomGT1 algorithm confidence score for this genotype assignment
- SNP Classification: SNPolisher[3] conversion type (category)

### Figure 4. AxLE Example Output



The table in *Figure 4* demonstrates the output of the Axiom Long format Export tool (AxLE). Each row represents a genotype call for a single SNP in single sample, described by both the Axiom native Forward format and the "TOP" format. In addition, the genotype call confidence as determined by the AxiomGT1 algorithm and the SNP classification by SNPolisher[3] is reported.

## COUNCIL ON DAIRY CATTLE BREEDING (CDCB) EXPORT TOOL

A clear requirement for the standardisation of allele designation is in the downstream application of genotyping data to genetic evaluation systems where mixing of the formats could be disastrous to the prediction of economically important traits. To support this specific use case in dairy cattle, Thermo Fisher Scientific has developed the CDCB (Council on Dairy Cattle Breeding) export tool[4]. Once an analysis has been completed in Axiom Analysis Suite[2], the CDCB export tool[4] performs three operations. Firstly, the "A/B" allele designations are swapped where the native "Forward" strand annotation differs from "TOP" based on a predefined list of affected markers. Secondly, the native SNP identifiers are mapped to the CDCB approved SNP identifiers. This occurs when a SNP has previously been submitted to the CDCB as part of another supported array and that name takes priority. Finally, it formats the calls and generates a sample sheet to enable direct upload to the Council on Dairy Cattle Breeding website. The tool is capable of consuming data from any Axiom catalogue bovine array and also custom bovine designs and is freely available to download.

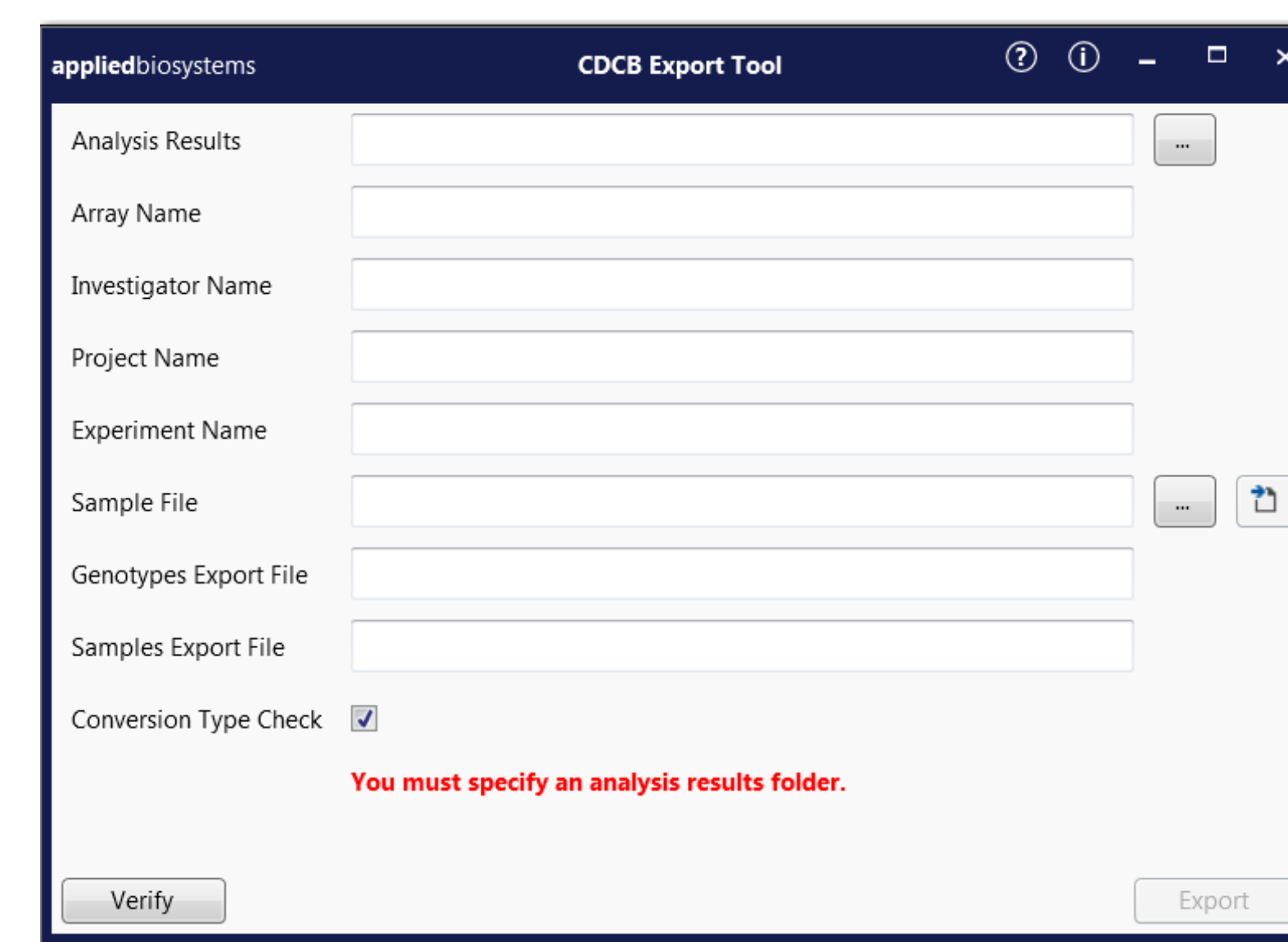### Figure 5. CDCB Export Tool Usage



*Figure 5* shows a screen shot of the CDCB export tool. After executing a Best Practices Workflow using Axiom Analysis Suite, the CDCB Export Tool can be accessed via the "External Tools" menu. The following steps are taken to generate CDCB compliant files for upload directly into the CDCB Genetics Evaluations system:
1. "Analysis Results" – navigate to the Axiom results folder to be exported.
2. "Array Name" will be automatically populated.
3. Complete the "Investigator Name", "Project Name" and "Experiment Name" fields. These will be used to populate the Sample export.
4. Generate a "Sample File" template with the ⬚ button and populate it with your sample information.
5. Set the "Genotypes Export File" and "Samples Export File" values.
6. Click Export to generate the output files.

### Figure 6. CDCB Export Tool Usage Example Genotyping Output



The table in *Figure 6* demonstrates the genotyping output of the CDCB Export tool. Meta data describing the processing, array, total SNPs on the array, reported SNPs and samples resides at the top, followed by a table of A/B genotyping calls in TOP format. SNP identifiers, where already present in the CDCB database prior to a new Axiom array being added are reported with the original SNP name. Where an Axiom array contains novel SNPs the SNP identifier native to that array is used. A mapping file between native SNP ID and CDCB SNP ID is provided with the array library files.

### Figure 7. CDCB Export Tool Usage Example Sample Output



The table in *Figure 7* demonstrates the sample output of the CDCB Export tool.

## SOFTWARE REFERENCES

[1] Axiom Long format Export tool (AxLE): *https://bit.ly/2R5JYr0*
[2] Axiom Analysis Suite: *http://bit.ly/2uqHODo*
[3] SNPolisher: *http://bit.ly/2tr8NC4*
[4] Council on Dairy Cattle Breeding export tool: *http://bit.ly/2soJRLq*

**Thermo Fisher SCIENTIFIC**