

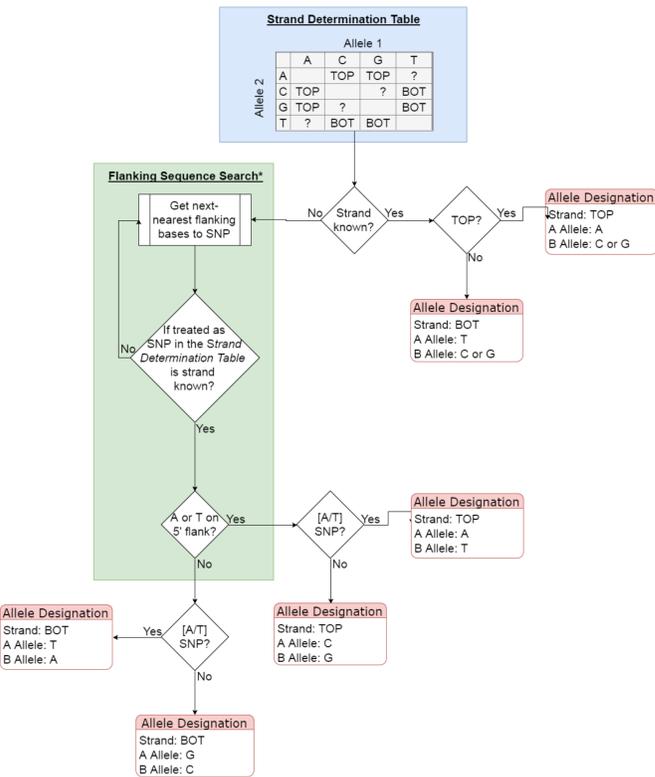
New Output Formats for the Applied Biosystems™ Axiom™ Genotyping Arrays

Joseph M Foster*, Alessandro Davassi, Ali Pirani, Shantanu Kaushikkar, 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

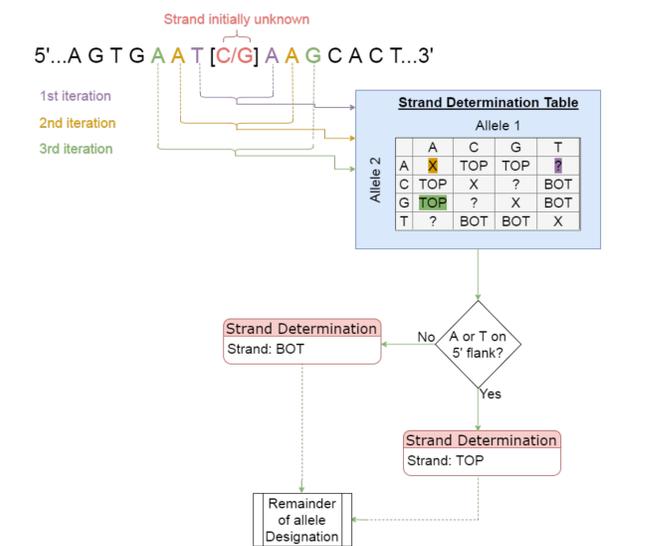


Figure 2 shows the strand determination process for [A/T] and [G/C] SNPs. Starting with the nearest pair of bases either side of the SNP, the *Strand Determination Table* is checked to see if a strand determination can be made. For each time a strand determination cannot be made the algorithm increments 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, we have developed the Applied Biosystems™ Axiom™ Long Format Export (AxLE) Tool¹; a companion application to the Applied Biosystems™ Axiom™ Analysis Suite software¹. 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 compatible with data from alternative arrays. 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

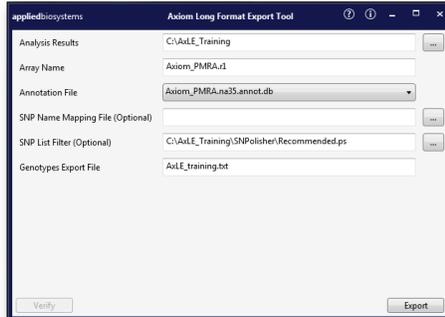


Figure 3 shows a screen shot of the AxLE tool. After executing a Best Practices Workflow using Axiom Analysis Suite, the AxLE tool can be accessed via the "External Tools" menu of Axiom Analysis Suite software. The following steps are taken to generate Long format:

- "Analysis Results": Select the appropriate Axiom Analysis Suite results folder of the analysis to be converted
- The "Array Name" and "Annotation File" will automatically populate once an analysis results folder has been selected
- Optionally a "SNP Name Mapping File" can be selected. This gives the end user to specify a file of Axiom SNP names to user-defined SNP names
- Optionally the analysis results can be filtered by a "SNP List File" to return only those SNP in the exported file
- A name must be assigned to the long format output file
- Click the "Export" button to complete

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 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: SNPClassifier² conversion type (category)

Figure 4. AxLE Example Output

1	A	B	C	D	E	F	G	H	I	J
1	[Header]									
2	Version	1.19.02.7.0								
3	Processing Date	3/6/17 10:14								
4	Content	Axiom_PMRA.r1								
5	Num SNPs	920636								
6	Total SNPs	920636								
7	Num Samples	190								
8	Total Samples	190								
9	[Data]									
10	SNP Name	Sample ID	Allele1 - Forward	Allele2 - Forward	Allele1 - TOP	Allele2 - TOP	Allele1 - AB	Allele2 - AB	Confidence	SNP Classification
11	AXFX-SP-000001	AS50778-4310250-123017-864_A01.CEL	C	G	C	G	A	B	0.00001	PolyHighResolution
12	AXFX-SP-000001	AS50778-4310250-123017-864_A02.CEL	G	C	C	A	A	A	0.00002	PolyHighResolution
13	AXFX-SP-000001	AS50778-4310250-123017-864_A03.CEL	C	G	C	G	B	B	0.00002	PolyHighResolution
14	AXFX-SP-000001	AS50778-4310250-123017-864_A04.CEL	G	C	C	C	A	A	0.00001	PolyHighResolution
15	AXFX-SP-000001	AS50778-4310250-123017-864_A05.CEL	G	C	C	G	A	B	0.00001	PolyHighResolution
16	AXFX-SP-000001	AS50778-4310250-123017-864_A06.CEL	C	G	C	G	B	B	0.00001	PolyHighResolution
17	AXFX-SP-000001	AS50778-4310250-123017-864_A07.CEL	G	C	C	G	A	B	0	PolyHighResolution
18	AXFX-SP-000001	AS50778-4310250-123017-864_A08.CEL	G	C	C	A	A	A	0.00002	PolyHighResolution
19	AXFX-SP-000001	AS50778-4310250-123017-864_A09.CEL	G	C	C	G	A	B	0	PolyHighResolution
20	AXFX-SP-000001	AS50778-4310250-123017-864_A10.CEL	C	G	C	G	B	B	0.00002	PolyHighResolution
21	AXFX-SP-000001	AS50778-4310250-123017-864_A11.CEL	G	C	C	G	A	B	0.00001	PolyHighResolution
22	AXFX-SP-000001	AS50778-4310250-123017-864_A12.CEL	G	C	C	G	A	B	0.00001	PolyHighResolution
23	AXFX-SP-000001	AS50778-4310250-123017-864_A13.CEL	G	C	C	G	A	B	0.00001	PolyHighResolution
24	AXFX-SP-000001	AS50778-4310250-123017-864_B01.CEL	G	C	C	G	A	B	0.00001	PolyHighResolution
25	AXFX-SP-000001	AS50778-4310250-123017-864_B02.CEL	G	C	C	G	A	B	0.00001	PolyHighResolution
26	AXFX-SP-000001	AS50778-4310250-123017-864_B03.CEL	G	C	C	A	A	A	0.00002	PolyHighResolution
27	AXFX-SP-000001	AS50778-4310250-123017-864_B04.CEL	C	G	C	G	B	B	0.00006	PolyHighResolution
28	AXFX-SP-000001	AS50778-4310250-123017-864_B05.CEL	C	G	C	G	B	B	0.00001	PolyHighResolution
29	AXFX-SP-000001	AS50778-4310250-123017-864_B06.CEL	G	C	C	A	A	A	0.00002	PolyHighResolution

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 a 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 SNPClassifier² 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 formats could be disastrous to the prediction of economically important traits. To support this specific use case in dairy cattle, we have developed the Applied Biosystems™ Council on Dairy Cattle Breeding (CDCB) export tool³; a companion application to Axiom Analysis Suite software. Once an analysis has been completed in Axiom Analysis Suite, the CDCB export tool performs three operations. Firstly, the "A/B" allele designations are swapped where the native "Forward" strand annotation differs to "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 catalog bovine array and also custom bovine designs and is freely available to download.

Figure 5. CDCB Export Tool Usage

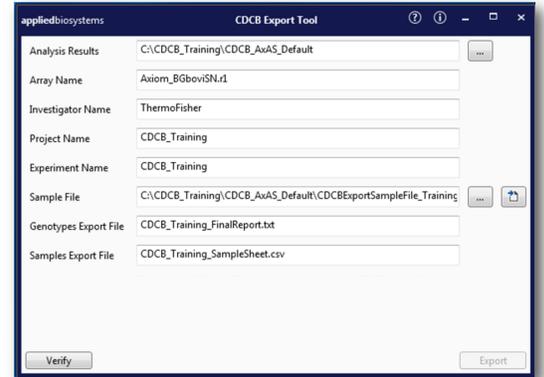


Figure 5 shows a screen shot of the CDCB export tool. After executing a Best Practices Workflow using Axiom Analysis Suite software, the CDCB Export Tool can be accessed via the "External Tools" menu of Axiom Analysis Suite. The following steps are taken to generate CDCB compliant files for upload directly into the CDCB Genetics Evaluations system:

- "Analysis Results" – navigate to the Axiom Analysis Suite results folder to be exported
- "Array Name" will automatically populate
- Complete the "Investigator Name", "Project Name" and "Experiment Name" fields. These will be used to populate the Sample export
- Generate a "Sample File" template with the button and populate it with sample information
- Set the "Genotype Export File" and "Sample Export File" values
- Click Export to generate the output files.

Figure 6. CDCB Export Tool Usage Example Genotyping Output

1	A	B	C	D	E	F	G	H	I
1	[Header]								
2	Axiom Analysis Suite Algorithm	1.19.02.7.0							
3	Processing Date	6/5/17 21:45							
4	Content	Axiom_BGovSN.r1							
5	Num SNPs	46578							
6	Total SNPs	57513							
7	Num Samples	95							
8	Total Samples	95							
9	[Data]								
10	Sample ID	Sample_Plate	Sample_Name	Project	AMP_Plate	Well in AMP Plate	CentriBarcode_A	CentriPosition_A	
11	sample10	550780432660706118202	sample10	CDCB_Training	550780432660706118202	009	SMP4_118202	R15C09	
12	Affx-93058216		AB	AB	AA	AA	AA	AA	AA
13	Hapmap51730-BTA-44937		AB	BB	BB	BB	BB	BB	BB
14	ARS-BFGL-NGS-11068		AB	AB	AA	AA	AA	AA	AA
15	Hapmap51730-BTA-44937		BB	BB	BB	BB	BB	BB	BB
16	ARS-BFGL-NGS-7215		BB	BB	BB	BB	BB	BB	BB
17	BTB-00430000		AB	BB	AB	AB	BB	BB	BB
18	ARS-BFGL-NGS-110683		AB	AB	BB	BB	AA	BB	BB
19	BTA-7865-no-rs		BB	BB	BB	BB	BB	BB	BB
20	ARS-BFGL-NGS-91754		BB	AB	AB	AB	AA	AB	AB
21	BTA-112164-no-rs		AB	BB	BB	BB	AA	BB	BB
22	Hapmap50247-BTA-117369		AA	AA	AA	AA	AA	AA	AA
23	BTB-0125802		AA	AB	AA	AA	AA	AA	AA
24	ARS-BFGL-NGS-70160		BB	AB	AB	AA	AB	BB	AB
25	ARS-BFGL-NGS-4015		BB	AB	BB	AB	AB	BB	BB
26	Hapmap4834-BTA-96124		BB	BB	AB	AB	BB	AB	BB
27	BTB-0042810		AB	AB	AA	AB	BB	BB	AB

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

1	A	B	C	D	E	F	G	H
1	[Header]							
2	Investigator Name	ThermoFisher						
3	Project Name	CDCB_Training						
4	Experiment Name	CDCB_Training						
5	Date	6/5/2017 9:45 PM						
6	[Manifests]							
7	Axiom_BGovSN.r1							
8	[Data]							
9	Sample_ID	Sample_Plate	Sample_Name	Project	AMP_Plate	Well in AMP Plate	CentriBarcode_A	CentriPosition_A
10	sample11	550780432660706118202	sample11	CDCB_Training	550780432660706118202	009	SMP4_118202	R15C09
11	sample13	550780432660706118202	sample13	CDCB_Training	550780432660706118202	019	SMP4_118202	R15C19
12	sample15	550780432660706118202	sample15	CDCB_Training	550780432660706118202	M09	SMP4_118202	R15C09
13	sample17	550780432660706118202	sample17	CDCB_Training	550780432660706118202	M09	SMP4_118202	R15C09
14	sample19	550780432660706118202	sample19	CDCB_Training	550780432660706118202	M07	SMP4_118202	R15C07
15	sample21	550780432660706118202	sample21	CDCB_Training	550780432660706118202	119	SMP4_118202	R09C19
16	sample24	550780432660706118202	sample24	CDCB_Training	550780432660706118202	A21	SMP4_118202	R01C21
17	sample30	550780432660706118202	sample30	CDCB_Training	550780432660706118202	E01	SMP4_118202	R05C01
18	sample35	550780432660706118202	sample35	CDCB_Training	550780432660706118202	E09	SMP4_118202	R05C09
19	sample36	550780432660706118202	sample36	CDCB_Training	550780432660706118202	K23	SMP4_118202	R11C23
20	sample37	550780432660706118202	sample37	CDCB_Training	550780432660706118202	G13	SMP4_118202	R07C13
21	sample38	550780432660706118202	sample38	CDCB_Training	550780432660706118202	K11	SMP4_118202	R11C11
22	sample33	550780432660706118202	sample33	CDCB_Training	550780432660706118202	G23	SMP4_118202	R15C23
23	sample43	550780432660706118202	sample43	CDCB_Training	550780432660706118202	G03	SMP4_118202	R07C03
24	sample44	550780432660706118202	sample44	CDCB_Training	550780432660706118202	C19	SMP4_118202	R03C19
25	sample45	550780432660706118202	sample45	CDCB_Training	550780432660706118202	E17	SMP4_118202	R05C17
26	sample46	550780432660706118202	sample46	CDCB_Training	550780432660706118202	G11	SMP4_118202	R07C11

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

SOFTWARE REFERENCES

- Axiom Analysis suite: <http://bit.ly/2uqHODo>
- SNPClassifier: <http://bit.ly/2tr8NC4>
- Council on Dairy Cattle Breeding export tool: <http://bit.ly/2soJRLq>



For Research Use Only. Not for use in diagnostic procedures