# illumina

## Infinium® iSelect® Custom Genotyping Assays

Guidelines for using the Illumina Assay Design Tool (ADT) to create and order custom assays.

## Introduction

The Illumina Infinium Assay is a powerful and widely used assay for highly multiplexed genotyping. Along with a broad portfolio of fixed content BeadArray<sup>™</sup> products, Illumina offers custom and semi-custom genotyping panels deployed on various BeadChip formats. Researchers can design their own content panels with the assistance of the online Assay Design Tool (ADT) and Illumina scientists. The Infinium iSelect HD and HTS Custom BeadChips can be completely customized to target between 3,072 and 1M custom bead types (e.g., SNPs and indels) per sample for any species. The BeadChips come in a 24-sample HD format for 3,072–90k attempted bead types (ABTs) or a 24-sample HTS format for >90k ABTs.

Researchers create custom panels by selecting and submitting a list of requested loci to Illumina. To ensure successful assay development, this list is evaluated within ADT to select an initial panel including desired assays that are predicted to have a high likelihood of success. Metrics returned by ADT provide predicted success information, Infinium validation status, and minor allele frequencies from published studies.

This technical note describes the process of designing, evaluating, and ordering a custom panel of markers. Researchers must submit preliminary design files describing the custom panel using one of four different file types corresponding to the different methods for preliminary evaluation of custom SNP loci: Gene, Region, Identity, and Sequence. Existing Design and Score File are additional input formats if SNPs have been previously scored by ADT. This information is evaluated by ADT and results are used to refine the panel definition. Each of the file type format options for ADT input is described with examples. Template files can be downloaded from the Mylllumina website<sup>1</sup>, or by contacting Illumina Technical Support<sup>2</sup>.

## Infinium Custom (iSelect/iSelect+) and Semi-Custom BeadChips

Illumina offers customers the benefit to augment content on existing custom or commercial BeadChips with newly discovered content from genome-wide association, whole-genome sequencing, and exome sequencing studies. Add-On content provides a significant benefit to customers who want to add newly discovered content to a custom array after the initial design period is completed.

The fully custom Add-On Content option (iSelect+), based on a custom iSelect Base Content BeadChip (e.g., a researcher's own iSelect BeadChip or an Illumina Consortium iSelect BeadChip), allows optimal customization of content. The maximum number of Add-On Content that can be supplemented depends on both the format of the BeadChip and the amount of base content already present. As a reference, the 24-HD BeadChip can support up to 90k total ABTs of Base + Add-On Content, and the 24-sample HTS BeadChip can support up to 700k total ABTs.

## Bead Types to SNPs

### Bead Types and Assay Design

Depending on the type of SNP or marker being assayed, the Infinium Assays use one of two probe (or bead type) designs: Infinium I or Infinium II. The Infinium II probe design, which stops at the base before the SNP of interest, uses only one probe per locus (i.e., one probe for both alleles). This probe design is suitable for the majority of loci in most organisms (e.g., approximately 84% of known SNPs in the human genome). Infinium I probe design is required for relatively less-common A/T and C/G SNPs and uses two probes (or bead types) per SNP to determine relative intensity ratio of the two possible target alleles.

For ordering and pricing, custom BeadChip products are defined by their number of ABTs, not number of loci assayed. Only if a project is exclusively limited to loci using Infinium II designs (A/G, A/C, T/G, T/C SNPs; and all indels) are the number of markers equal to the number of bead types. The more markers ordered that require Infinium I design, the fewer total markers that can be attempted on a custom BeadChip.

• For example, a researcher designing a custom iSelect BeadChip containing 9,500 Infinium II assays and 500 Infinium I assays would require a BeadChip with a total of 10,500 ABTs: 9,500 Infinium II bead types (one probe per marker) plus 1,000 Infinium I bead types (two probes per marker).

## Bead Type Success Rate

Customized Infinium BeadChip manufacturing achieves at least 80% of the attempted assays. However, Infinium iSelect and semi-custom BeadChips commonly achieve greater than 80% conversion. Please contact Illumina if your custom project has specific conversion requirements.

"Must-have" markers can be duplicated in the final design file to increase its likelihood of converting onto the final array. This significantly increases an assay's ability to avoid a random exclusion from the final BeadChip product. If an assay is to be replicated, it requires a unique assay name in the input file submission.

#### **Beadpool and Product Shelf Life**

Illumina guarantees the life of custom iSelect HD and HTS beadpools for 12 months from the date of manufacture. Illumina should be notified, upon purchase (prior to manufacturing), if a custom project is part of a Consortium (shared by multiple institutions) or is forecasted to require support of greater than 1 year.

## **Overview of Workflow**

Custom assay designs begin with ADT and include the following steps outlined in three simple stages.

#### Stage 1. Create an Input File

1. Create Input file with custom markers using one of the Illumina templates downloaded from Mylllumina.

#### Stage 2. Preliminary Scoring

- 1. On the Mylllumina welcome page, select Custom Products | Infinium iSelect Genotyping.
- 2. Under the Preliminary Files Tab, click **Start a new design**.
- 3. Select iSelect, iSelect+, or Commercial+.
- 4. Select the File Type.
- Species defaults to Human for Identity, Gene, and Region files. If submitting a Sequence, Existing Design, or Score file select the desired Species (or Other if not listed).
- 6. Enter the **Description** and **Comments** (both options are userdefined and optional).
- 7. Click **Browse** to navigate to the input file created using one of the previously downloaded Illumina templates in Stage 1.
- 8. Within a few minutes or hours (typical) or up to 2 days, ADT scoring will be completed, and an email notification will be sent. At that point, the file status of the file being scored will be updated from **Processing** to **Success** under the Preliminary Files Tab. In some instances, the Status may appear as **Failed** or **Exception**. Make sure that the input file is correctly formatted and that its size in within the specified limit of 700k markers.
- Download and review the preliminary score results. Edit the list based on error and warning codes, design scores, and research requirements.

#### Stage 3. Final Validation

- 1. Under the Final Files Tab, click Upload Final File.
- 2. Select iSelect, iSelect+, or Commercial+.
- Select the File Type, Species, enter the Description and Comments (optional), and Browse to the edited and finalized score file.
- 4. After the Final Score File is uploaded and successfully validated, Status will update to Content Success. In some instances, the Status may appear as Content Warning (markers with warnings are included in design), Failed (uploaded file incorrectly formatted), or Content Failure (markers with errors are included in design, or attempted design outside the scope for iSelect).

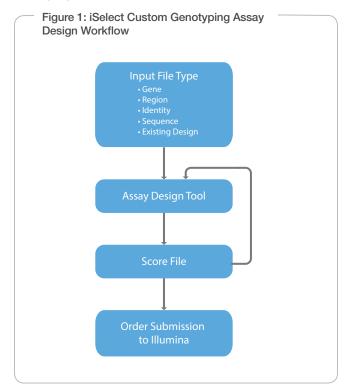
#### Stage 4. Ordering

1. Place an order for the product by selecting **Order** under the Final Files tab in ADT, and follow the instructions on Mylllumina.

## **Detailed Workflow**

#### **Preliminary Input Files**

To submit preliminary input files, researchers can use different file types corresponding to the different methods for preliminary evaluation of custom SNP loci: Gene, Region, Identity, Sequence, and Existing Design. After preliminary evaluation, ADT generates a Score file that can be used as an input file in subsequent rounds of evaluation or ordering (Figure 1).



At this time, ADT returns only human sequences from Gene, Region, and Identity input files. It is important to note that ADT only supports one build of the human genome at a time. Illumina keeps the supported version of the human genome current and gives users at least two weeks notice before switching to a new version. ADT returns the version and build of the current reference database in score results.

ADT supports full duplicate and repeat region checking for assays where full genome reference support is available (e.g., human, mouse, rat, and bovine). When designs are created for other species supported by ADT (e.g., those found in the dropdown menu of the ADT portal in Mylllumina), varying levels of support for duplicate and repeat region checking will be given, depending on the availability of public sequence data. If ADT does not list the species in the Mylllumina ADT portal, no duplicate and repeat region checking will be performed and customers should ensure their assays target unique sequences. The ADT output **Final Score** file includes the 50 bp Illumina design probe sequences so that researchers can cross-check against their own reference data to confirm that these probes hybridize to unique locations in the genome.

Input files may be created or edited with a text editor or spreadsheet program. However, before submitting them to ADT, files must be saved in a comma-separated values (\*.csv) format. The examples provided in this document show files created in Microsoft Excel. Blank lines are not permitted in the data fields. The following formatting specifications must be met.

- Comma-separated values with a \*.csv file extension. Because the input file format is comma-delimited, no commas may be used within the values.
- Each file type includes specific column headings for the data, as described in Table 1.
- File contains no more than 700k markers or indels. If the number of markers exceeds this limit, the file must be split into smaller files with a maximum of 700k lines.

#### Gene List

The Gene List file type provides a method for returning designs on all markers within a gene and in the regions upstream and downstream from a gene in the supported build of the human genome (Table 2). A Gene List requires input using RefSeq NM accession IDs (preferred) or HUGO identifiers. ADT maps these accession numbers to the human genome to identify gene regions and return all SNPs in those regions. The size of upstream and downstream bases is customizable via the Gene List input file format (Figure 2). Markers in overlapping gene regions will be listed in the Score output file only one time, but will be annotated as being present in both regions in the design output Region\_Description field. The 700k marker limit applies to the Gene List file; a general guideline is that a Gene List with up to 600 genes and a range of 10,000 bases upstream and downstream will be within the marker limit.

#### Figure 2: Gene List Format Example

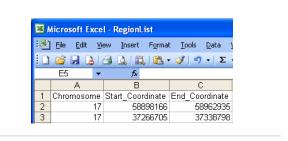
× 1	Microsoft Excel - GeneList							
:	Eile	<u>E</u> dit	⊻iev	v <u>I</u> nser	t F <u>o</u> rmat	<u>T</u> ools	<u>D</u> ata	Window
1	2		16	3 🖪   1	🖏   🛍 -	🦪   🌖	- Σ	- 2↓   10
	F8		•	fs.	r			
		Α			В		С	
1	Gene	Nam	е	Bases	Upstream	Base	s_Dowr	nstream
2	GenelD:1073		500		D		500	
3	GenelD:11261		500		D		500	
4	GenelD:6387		500		D	500		
5	NM_020134.2		500		0		500	
6	NM_182685.1		500		0		500	
7	CHRNA1			500		0		500

Example of properly formatted entries in a Gene List, suitable for upload through Mylllumina.

#### **Region List**

A Region List file contains a list of regions in the human genome identified by physical chromosome and coordinates. ADT will search and evaluate from among cataloged markers in a current Illuminainternal version of dbSNP. This internal database does not contain multinucleotide repeats (MNPs), microsatellites (simple sequence repeats), or markers with ambiguous or multiple locations. Markers in overlapping regions will be listed in the Score output file only one time, but will be annotated as being present in both regions in the Region\_Description field. Because ADT limits output to 700k markers, submitting fewer than 60 Mb of regions per file is recommended. Figure 3 provides an example of a properly formatted Region List, suitable for upload via Mylllumina.

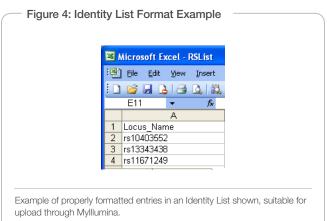




Example of properly formatted entries in a Region file, suitable for upload through Mylllumina.

#### **Identity List**

Known markers described in the current version of dbSNP for the human reference genome can be requested specifically using the Identity List. A current internal version of dbSNP is the source for rs marker and flanking sequence data. The column heading **Locus\_ Name** is the only input field in an Identity List. Figure 4 provides an example of a properly formatted Identity List, suitable for upload via Mylllumina. An error output file will be generated to indicate any merged SNP IDs or unsupported molecule (e.g., RNA) or marker types (e.g., tri-allelic markers).



#### **Sequence List**

The SequenceList allows researchers to evaluate markers from private databases or other sources, including any species. The Locus\_Name field is used to name sequences for easy identification. Locus\_Name entries contained in this file must not begin with "rs" because that prefix designates rs-IDs in the local dbSNP database and will trigger a database search.

RSID sequences from dbSNP within compatible adjacent polymorphisms can be adjusted to be designable by replacing an adjacent polymorphism with a known major allele in the population,

	Table	1:	Score	Output	File	Column	Headers
--	-------	----	-------	--------	------	--------	---------

Column	Description			
Locus_Name	rsID or customer's unique name			
Sequence	The bracketed site identified by the Locus_Name with $\geq$ 50 bases of flanking sequence			
Genome_Build_Version	Genome build. Contact Technical Support <sup>2</sup> for the currently supported build.			
Chromosome	Chromosome on which the marker is located. Must be a valid chromosome for the species being analyzed.			
Coordinate	Chromosomal coordinate of marker			
Source	Source of the sequence and annotation data			
Source_Version	Source version number			
Sequence_Orientation	Must contain one of the following three values: <b>forward, reverse</b> , or <b>unknown</b> . A score file resulting from a sequence input and field <b>Sequence_Orientation</b> is customer-supplied and not validated.			
Region_Description	Description of the region of interest			
Final_Score	Final scores are based on a proprietary algorithm and can range from 0 to 1 with higher values reflecting their "probe-ability," or likelihood of success for a custom assay design. A final score of 1.1 indicates that the SNP has been specifically validated by Illumina as a successful design for Infinium.			
Failure_Codes	If applicable, reasons why a successful assay at this marker locus is unlikely (see Table 8)			
Validation_Class	Numerical representation of validation_bin (see Table 9)			
Validation_Bin	Manner in which designed assays have been validated (see Table 9)			
MAF_(Population)	Minor allele frequency from the largest peer-reviewed study conducted in the indicated population, the			
Chr_Count_(Population)	study size in terms of number of chromosomes, and the study type. Data are retrieved from dbSNP for each of the following human genome populations. Score output resulting from Sequence input files will be blank for these fields.			
Study_Name_(Population)				
	Caucasian, African, African-American, Han Chinese, Japanese, Unknown			
Normalization_Bin	A, B, or C			
Bead_Types/Assay	1 for Infinium II or 2 for Infinium I			
Assay_Type	Infinium I or Infinium II			
ILMN_ID	Unique identifier assigned by ADT for the designed assay			
Gene_ID	Gene ID number from NCBI			
Gene_symbol	HUGO identifier			
Accession	RefSeq Accession number			
Location	Structural location of the marker: intron, coding, flanking_5UTR, flanking_3UTR, 5UTR, 3UTR, UTR			
Probe Sequences	Illumina's 50 bp design probe sequences			
Location_relative_to_gene	If the marker does not fall within an exon, the value is the actual base pair distance from gene start. The absolute value of this number is the distance to the closest transcript. The negative sign is a formatting symbol and is not meant to imply strand or direction. If the SNP is within an exon, two values separated by a '/ ' are gvien, which represent distances to the exon-intron boundaries. Information in this column can be used to identify potential splice site variants.			
Coding_status	<b>NONSYN</b> or <b>SYNON</b> . If the marker falls within an exon, this field notes a synonymous or nonsynonymous amino acid change.			
Amino_acid_change	If the marker falls within an exon, this field notes the actual change to the amino acid, followed by the GenBank protein sequence used in numbering the change			
ld_with_mouse	Ratio of identical bases within 60 bp of flanking sequence compared to mouse sequence that has been aligned with the homologous human sequence and covers the marker in question			
Phast_conservation	Metric used by the UCSC Genome Browser to identify highly conserved markers among species			
Design_Date	Date of design			

Column	Description
Gene_Name	Customer-supplied gene name. Can be a RefSeq accession ID or HUGO gene symbol
Bases_Upstream	Number of bases to search upstream of the gene starting coordinate
Bases_Downstream	Number of bases to search downstream of the gene starting coordinate

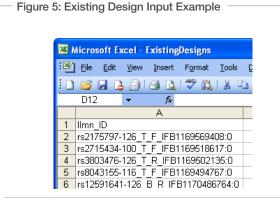
if the adjacent polymorphism is greater than 10 bases from the target SNP (20 bases recommended). To submit such designs as a sequence file, add a prefix (e.g., **adj\_rs1234**).

To specify a SNP, the customer should put brackets around a polymorphic locus in the submitted sequence, and separate the two alleles with a forward slash (TGC[A/C]CCG). Similarly, to specify an indel, a forward slash should be used between a single minus sign (indicating the deletion) and bases representing the insertion (e.g., TGC[-/AT]CCG). A minimum of 50 bp of sequence on either side of the SNP is ideal for evaluating both strands for the best design. ADT will also accept IUPAC codes for degenerate bases in the flanking sequence and avoid placement of probes over these polymorphisms adjacent to the targeted SNP. If the Lowercase\_Weighting checkbox on the Mylllumina submission form is checked (or file header value in an email submission indicates TRUE for the Lowercase\_Weighting option), lowercase nucleotides will have penalized final scores reflecting suboptimal probe placement over these lowercase regions. Because lowercasing in public databases is not a standard way to indicate repetitive or duplicated regions, Illumina recommends clearing the Lowercase\_Weighting checkbox by default.

The columns and description information shown in Table 3 must be provided in the Sequence List.

#### **Using Previous Assay Designs**

Illumina has created a method for conveniently ordering the same assays that were designed and used on a previous iSelect or Illumina Commercial product. As shown in Figure 5, an ExistingDesign file only contains a list of the Illumina IDs (Ilmn\_Ids) from the original design or beadpool manifest.



Example of properly formatted entries in an Existing Design List, suitable for upload through Mylllumina.

#### **Score Output File**

ADT preliminary file submission results are returned as a Score file for review and revision, or for input into a Final file submission at the time of purchase.

The Score file header section includes additional summary information, such as the total number of markers in the file. This is further broken down into numbers of markers in each of three Normalization\_Bins: A, B, and C. Bin C assays include all Infinium Il designs requiring a single bead type. Bin A and B assays are Infinium I designs in the red and green channels, respectively, and are classified into one of these two bins based on the color channel required to detect the target alleles across the two bead types used in the Infinium I assay. For a final order, if there are any loci in a given optimal normalization bin, there must be at least 100 loci in that bin to ensure normalization of the intensity data from the scanner. If Bin A or Bin B needs to be supplemented to reach 100, the customer should submit additional A/T or C/G SNPs to ADT for scoring or change the Force Infinium I column in the input file (for an Infinium II Bin C SNP) to increase representation in Bin A and/or Bin B. ADT will report whether each of these SNPs are in Bin A or Bin B in the output score file. The appropriate assays can then be added to the original design as needed.

The Number of Bead Types value is important for ordering, because iSelect BeadChip orders are priced based on the number of ABTs. The number of bead types may be different from the number of assays because Infinium I assays require two bead types per marker and Infinium II assays only require one bead type per marker.

Following the Score file header section, detailed information for each marker is listed in the data section. All Score columns are described in Table 1. Important performance values are also presented for each SNP. The **Final\_Score** indicates the expected success for designing a given assay, and may be supplemented with **Failure\_Codes** for further information (Table 4). Validation status is also indicated to provide even greater confidence in design success. To assist researchers in ordering the most applicable markers for their studies, minor allele frequencies (MAFs) in several populations are provided for SNPs when available from dbSNP. MAF from the largest study is reported, and is qualified based on peer-reviewed publication, study design and study size, and verified results.

#### **Filtering and Selecting Custom Lists**

In addition to being an output file format, Score files can be used as input files to ADT. Thus, users can create a filtered or edited output file (with markers removed or added) for iterative ADT analysis during final SNP selection. Markers identified using more than one input search

Column	Description	
Locus_Name	Customer-supplied unique locus name (cannot begin with rs or cg). Name must not contain any of the following characters: " $\% / = \?$ , @;', or space. Duplicate probe names are not allowed. Underscores are discouraged and have the risk of interfering with the assignment of Illumina IDs (ILMN_IDs) in the Illumina database.	
Target_Type	Must be SNP or INDEL (case-insensitive)	
Sequence	Limited to 10,000 bases. May only contain one bracketed SNP or indel. Output will be $\leq$ 122 bases per line.	
Chromosome	Chromosome on which the marker is located. Must be a valid chromosome for the species being analyzed. Enter <b>0</b> if unknown. Contig numbers are not accepted and may result in a corrupt manifest if passed through ADT.	
Coordinate Chromosome coordinate of marker. Enter <b>0</b> if unknown.		
Genome_Build_Version	Version number supplied by customer. Otherwise, enter 0.	
Source Identify the source of the sequence and annotation data. Must be completed. Enter <b>unknowr</b> information is available.		
Source_Version	Source version number. Enter 0 if unknown.	
Sequence_Orientation	Must be either <b>Forward</b> , <b>Reverse</b> , or <b>Unknown</b> (case-insensitive). For consistency and cross-platform data compatibility, Forward and Reverse strand should be based upon the public database definition of strand, if available.	
Plus_Minus Must be either <b>Plus</b> , <b>Minus</b> , or <b>Unknown</b> (case-insensitive). For consistency and cross-platform data compatibility, Plus and Minus strands should be based upon the public database definition of strand, if		
Force_Infinium_I	Must be <b>TRUE</b> or <b>FALSE</b> ; enables the ability to force an Infinium I design, which may be needed to fulfill a normalization bin.	

## Table 3: Sequence List Column Descriptions

method (e.g., Gene, Region, Identity, Sequence, or ExistingDesign) can be combined as one Score file and resubmitted to ADT as an input file for evaluation as a single list.

Illumina recommends applying the following criteria for discriminating marker lists to create a successful product that achieves the scientific aims of the experiment and has the highest probability of generating meaningful results.

- 1. Remove markers that cannot be ordered (error codes in the 101–110 range).
- 2. Consider research requirements (e.g., tags, spacing, or MAF).
- 3. When appropriate, favor Infinium-validated markers, because they have successfully converted to functional assays.
- 4. Use two-hit validated markers (Table 5) based on the Validation\_ Bin field. Higher Final\_Scores are preferred.
- 5. Do not be hindered by proximity limits for Infinium assays. Markers can typically be within 11 to 50 bases without interfering with probe hybridization to the targeted SNP. Changing IUPAC codes to the major allele in the target population for adjacent polymorphisms beyond 11 bases from the targeted SNP will allow ADT to consider the flanking sequence in sequencing for design. This is recommended only when the initial design is failed for both sides of the SNP.
- Submit a number of bead types and corresponding SNPs to ADT that are at least 20% over the targeted number for final design. It is easier to remove SNPs from the design file than to

go back and add in more SNPs that have not yet been evaluated by ADT. For example, Human OmniExpress+ supports 1,000 to 50,000 ABTs; therefore, Illumina recommends starting the selection process by submitting at least 1,200 to 60,000 SNPs, respectively, for consideration in the final design.

#### **Final Score File**

The Final Score file is created automatically by ADT during Final scoring. However, the "header" section, created by ADT during preliminary scoring, must be deleted before this file can be submitted for Final Scoring. Make sure that the body of the scored file retains **Locus\_Name** as the first line of the CSV.

## Ordering

To place an order, the edited score file from preliminary scoring is submitted for Final Scoring. If the list passes without errors, the product will automatically be available for ordering. Pricing will be displayed before the order is entered. Alternatively, the customer can request a quote from Illumina before initiating the design process.

After a file is successfully uploaded under the **Final Files** tab (see *Overview of Workflow*), an order can be placed by selecting **Order**. The order finalization and confirmation page opens to request the number of samples requested for this project to **Get/Update Pricing** and enable **Add to Cart**. After **Add to Cart** is selected, a final review page will populate with the option to **UPDATE** or **CHECKOUT**. When **CHECKOUT** is selected, the order enters into Illumina's manufacturing process.

#### Table 4: List of ADT Failure Codes

101	Flanking sequence is too short					
102	Polymorphism or sequence formatting error. Possible causes:					
	Check polymorphism format: SNP => $[X/Y]$ , INDEL => $[-/XYZ]$ , CpG => $[CG]$					
	More than one set of brackets in sequence					
	Missing brackets around polymorphism					
	SNP alleles not separated by a "/"					
	Spaces found in submitted sequence.					
103	TOP/BOT strand cannot be determined. Possible cause:					
	Low sequence complexity					
104	Variant is not appropriate for Illumina platform. Possible causes:					
	Variant is not bi-allelic					
	Contains characters other than A,G,C, or T					
	Non-DNA molecule type not supported (e.g., RNA)					
106	Degenerate nucleotide(s) in assay design region (e.g., W, R, S, N)					
107	SNP sequence not found					
109	Indels not supported for Infinium I					
110	Locus name duplicated in the base commercial content					
Warnings	designable)					
301	Polymorphism in duplicated/repetitive region					
302	Melting temperature (Tm) outside assay limits					
304	There are known SNPs within the probe region. See Underlying_SNP column for details.					
311	SNP Inappropriate for Infinium I assay type					
340	Another polymorphism in this list is $\leq$ 60 nucleotides away					
360	Low score warning					
399	Multiple contributing issues					
601	Potentially nonspecific against the genome					
602	Probe sequence is duplicated in the base commercial content					
603	Probe sequence is duplicated in the custom content					

#### **Reordering an Existing Beadpool or Product**

To place an order based on an existing product or custom beadpool, click **My Chips/Reorder**. Reorders have a minimum order requirement of 288 samples, and must be ordered in multiples of 48. Only bead pools in inventory within acceptable QC metrics are available for reordering. Illumina Customer Service will receive a notification to confirm that inventory is available before processing the online order. If a design file has been used to place an BeadChip order, that design file is never removed.

#### **Completing an Order**

Once all desired orders are added to the shopping cart, click **Continue to Checkout**. The shopping cart is automatically displayed after clicking **Add to Cart**; otherwise click **View Cart**. During the checkout process, shipping and payment information is required. An initial shipping date and partial shipment request can be entered before final submission. After clicking **Submit**, the order will be sent to Illumina Customer Service for review. If needed, customers can modify an order by emailing orders@illumina.com during a 48 hour review period immediately following submission.

#### **Shipping Schedule**

Shipping schedules can be defined after an order is placed through Mylllumina. Otherwise, orders may ship complete or according to a default schedule. The first shipment of custom BeadChips will typically arrive within 8–12 weeks after order confirmation.

#### **Ordering by Email**

To place an order by email, submit the Final Score output file and a purchase order to orders@illumina.com. Illumina will send a confirmation and contact you if essential information is missing or if it contains unorderable designs.

#### **Table 5: Validation Status Descriptions**

Validation_Bin	Validation_Class	Description
NonValidated	1	Locus has been seen in only one study or population. Even if it has a high design score, there is a chance that it is monomorphic.
OneKGenomeValidated	100	Locus has been sequenced in the 1000 Genomes Project.
TwoHitValidated	110	Both alleles have been seen in two independent studies and populations.
HapMapValidated	120	Locus has been genotyped by the HapMap Project.
TwoHit_OneKGenomeValidated	200	Both alleles have been seen in two independent studies and populations. Locus has been sequenced in the 1000 Genomes Project.
HapMap_OneKGenomeValidated	210	Locus has been genotyped by the HapMap Project. Locus has been sequenced in the 1000 Genomes Project.
TwoHit_HapMapValidated	220	Locus has been genotyped by the HapMap Project. Both alleles have been seen in two independent studies and populations.
TwoHit_HapMap_OneKGenomeValidated	300	Locus has been genotyped by the HapMap Project. Both alleles have been seen in two independent studies and populations. Locus has been sequenced in the 1000 Genomes Project.
InfiniumValidated	910	Design has been previously designed and successfully generated polymorphic results using the Infinium assay.

## Summary

Custom Infinium DNA Analysis products by Illumina allow researchers to create assays tailored to their specific needs for focused, highthroughput genotyping or fine-mapping of candidate disease association regions. The ADT is a simple and powerful tool for evaluating individual loci and creating the successful custom BeadChips for genotyping. By following the guidelines in this technical note, researchers can make sure that their orders are designed and placed quickly and easily.

#### References

- 1. https://my.illumina.com/Custom/Index (requires login)
- 2. To contact Technical Support, send email to techsupport@illumina.com or call 1.800.809.4566.

## Additional Information

Visit www.illumina.com or contact us to learn more about iSelect Custom Infinium products from Illumina.

Illumina • 1.800.809.4566 toll-free (U.S.) • +1.858.202.4566 tel • techsupport@illumina.com • www.illumina.com

#### FOR RESEARCH USE ONLY

© 2007-2014 Illumina, Inc. All rights reserved.

Illumina, BeadArray, Infinium, iSelect, the pumpkin orange color, and the Genetic Energy streaming bases design are trademarks or registered trademarks of Illumina, Inc. All other brands and names contained herein are the property of their respective owners. Pub. No. 370-2007-021 Current as of 03 June 2014

illumina