

BLUEPRINT proposed data submission schemas  
Draft v0.1

José María Fernández

September 13, 2012



# Contents

<b>1</b>	<b>Tabular format of input files</b>	<b>5</b>
1.1	Regulatory regions - Metadata File . . . . .	5
1.2	Regulatory regions - Primary Analysis File . . . . .	6
1.3	Protein-DNA interaction - Metadata File . . . . .	6
1.4	Methylation - Secondary Analysis File . . . . .	7
1.5	Copy Number Somatic Mutations - Secondary Analysis File . . . . .	8
1.6	Methylation - Metadata File . . . . .	8
1.7	Analyzed Sample Data File . . . . .	9
1.8	Donor Data File . . . . .	10
1.9	Expression - Metadata File . . . . .	10
1.10	Expression - Gene File . . . . .	11
1.11	Regulatory regions - Secondary Analysis File . . . . .	12
1.12	Simple Somatic Mutations - Primary Analysis File . . . . .	13
1.13	Exon Junction - Metadata File . . . . .	15
1.14	Structural Somatic Mutations - Secondary Analysis File . . . . .	15
1.15	Simple Germline Variations - Metadata File . . . . .	17
1.16	Specimen Data File . . . . .	17
1.17	Exon Junction - Primary Analysis File . . . . .	18
1.18	Donor Family History . . . . .	20
1.19	Copy Number Somatic Mutations - Primary Analysis File . . . . .	20
1.20	Structural Somatic Mutations - Primary Analysis File . . . . .	21
1.21	Protein-DNA interaction - Secondary Analysis File . . . . .	23
1.22	Simple Germline Variations - Primary Analysis File . . . . .	23
1.23	Simple Somatic Mutations - Metadata File . . . . .	25
1.24	Methylation - Primary Analysis File . . . . .	25
1.25	Simple Somatic Mutations - Secondary Analysis File . . . . .	26
1.26	Copy Number Somatic Mutations - Metadata File . . . . .	27
1.27	Protein-DNA interaction - Primary Analysis File . . . . .	28
1.28	Structural Somatic Mutations - Metadata File . . . . .	29
<b>A</b>	<b>CV Tables</b>	<b>31</b>
A.1	CV Table appendix_B10.tsv . . . . .	31
A.2	CV Table appendix_B5.tsv . . . . .	32
A.3	CV Table appendix_B12.tsv . . . . .	33
A.4	CV Table appendix_B6.tsv . . . . .	34
A.5	CV Table meth_p__chromosome_strand.tsv . . . . .	35
A.6	CV Table meth_p__validation_status.tsv . . . . .	35
A.7	CV Table sp__analyzed_sample_type.tsv . . . . .	35
A.8	CV Table dr__donor_sex.tsv . . . . .	35
A.9	CV Table dr__donor_vital_status.tsv . . . . .	35
A.10	CV Table exp_g__gene_strand.tsv . . . . .	35
A.11	CV Table exp_g__is_annotated.tsv . . . . .	35
A.12	CV Table exp_g__validation_status.tsv . . . . .	36
A.13	CV Table ssm_p__mutation_type.tsv . . . . .	36

A.14	CV	Table ssm_p__chromosome_strand.tsv . . . . .	36
A.15	CV	Table ssm_p__refsnp_strand.tsv . . . . .	36
A.16	CV	Table ssm_p__is_annotated.tsv . . . . .	36
A.17	CV	Table ssm_p__validation_status.tsv . . . . .	36
A.18	CV	Table jcn_m__seq_coverage.tsv . . . . .	36
A.19	CV	Table specimen__specimen_type.tsv . . . . .	37
A.20	CV	Table specimen__specimen_processing.tsv . . . . .	37
A.21	CV	Table specimen__specimen_storage.tsv . . . . .	37
A.22	CV	Table specimen__specimen_available.tsv . . . . .	37
A.23	CV	Table jcn_p__gene_strand.tsv . . . . .	37
A.24	CV	Table jcn_p__exon1_strand.tsv . . . . .	38
A.25	CV	Table jcn_p__exon2_strand.tsv . . . . .	38
A.26	CV	Table jcn_p__is_fusion_gene.tsv . . . . .	38
A.27	CV	Table jcn_p__is_novel_splice_form.tsv . . . . .	38
A.28	CV	Table jcn_p__junction_type.tsv . . . . .	38
A.29	CV	Table jcn_p__validation_status.tsv . . . . .	38
A.30	CV	Table dr_family_history__relationship_type.tsv . . . . .	38
A.31	CV	Table dr_family_history__relationship_sex.tsv . . . . .	39
A.32	CV	Table cns_m_p__mutation_type.tsv . . . . .	39
A.33	CV	Table cns_m_p__is_annotated.tsv . . . . .	39
A.34	CV	Table cns_m_p__validation_status.tsv . . . . .	39
A.35	CV	Table appendix_B9.tsv . . . . .	39
A.36	CV	Table stsm_p__chr_from_strand.tsv . . . . .	40
A.37	CV	Table stsm_p__chr_to_strand.tsv . . . . .	40
A.38	CV	Table stsm_p__evidence.tsv . . . . .	40
A.39	CV	Table stsm_p__zygosity.tsv . . . . .	40
A.40	CV	Table stsm_p__validation_status.tsv . . . . .	40
A.41	CV	Table sg_v_p__variation_type.tsv . . . . .	40
A.42	CV	Table sg_v_p__chromosome_strand.tsv . . . . .	41
A.43	CV	Table sg_v_p__refsnp_strand.tsv . . . . .	41
A.44	CV	Table sg_v_p__is_annotated.tsv . . . . .	41
A.45	CV	Table sg_v_p__validation_status.tsv . . . . .	41
A.46	CV	Table appendix_B7.tsv . . . . .	42

# Chapter 1

## Tabular format of input files

### 1.1 Regulatory regions - Metadata File

Regulatory regions [rreg] – Metadata File [m]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
matched_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed control sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix_B10.tsv)
platform	VARCHAR(512)	R	Platform or technology used in detecting the methylation (See Table appendix_B5.tsv)
experimental_protocol	VARCHAR(512)	O	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
variation_calling_algorithm	VARCHAR(512)	R	Name of variation calling algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	O	Names of other analysis algorithms. Separate multiple algorithms by commas.
raw_data_repository	VARCHAR(128)	O	Public repository where raw data is submitted (See Table appendix_B12.tsv)
raw_data_accession	VARCHAR(128)	O	Accession and URL for referencing the raw data at the public repository
note	TEXT	O	Optional field to leave notes

## 1.2 Regulatory regions - Primary Analysis File

Regulatory regions [rreg] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	TEXT	R	Unique identifier for the analysis performed for a particular group of samples
analyzed_sample_id	TEXT	R	Unique identifier for the analyzed sample
regulatory_region_id	VARCHAR(128)	R	Unique identifier for the identified regulatory region
chromosome	VARCHAR(64)	R	Name of the chromosome containing the methylation (See Table <i>appendix_B6.tsv</i> )
chromosome_start	INTEGER	R	Start position of the methylation on the chromosome
chromosome_end	INTEGER	R	End position of the methylation on the chromosome
chromosome_strand	INTEGER	O	Chromosome strand (See Table <i>meth_p__chromosome_strand.tsv</i> )
normalized_read_count	FLOAT(5,2)	R	Normalized count of sequencing reads if analyzed by sequencing platforms
raw_read_count	INTEGER	R	Raw count of sequencing reads if analyzed by sequencing platforms
validation_status	VARCHAR(64)	R	Indicate if the methylation has been validated (See Table <i>meth_p__validation_status.tsv</i> )
validation_platform	VARCHAR(512)	O	Platform or technology used in validation (See Table <i>appendix_B5.tsv</i> )
note	TEXT	O	Optional field to leave notes

## 1.3 Protein-DNA interaction - Metadata File

Protein-DNA [pdna] – Metadata File [m]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
matched_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed control sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix_B10.tsv)
platform	VARCHAR(512)	R	Platform or technology used in detecting the methylation (See Table appendix_B5.tsv)
experimental_protocol	VARCHAR(512)	O	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
variation_calling_algorithm	VARCHAR(512)	R	Name of variation calling algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	O	Names of other analysis algorithms. Separate multiple algorithms by commas.
raw_data_repository	VARCHAR(128)	O	Public repository where raw data is submitted (See Table appendix_B12.tsv)
raw_data_accession	VARCHAR(128)	O	Accession and URL for referencing the raw data at the public repository
NSC	FLOAT(5,2)	O	Normalized strand-cross correlation
RSC	FLOAT(5,2)	O	Relative strand-cross correlation
note	TEXT	O	Optional field to leave notes

## 1.4 Methylation - Secondary Analysis File

Methylation [meth] – Secondary Analysis File [s]

Name	Type	R/O	Description / Values
analysis_id	TEXT	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	TEXT	R	Unique identifier for the analyzed sample
methyalted_fragment_id	TEXT	R	Unique identifier for the methylation
gene_affected	VARCHAR(128)	R	Gene(s) containing the methylation. Use Ensembl gene id. Separate multiple genes with vertical bars in the form of geneA—geneB—geneC. If no gene is affected, use -888 (not applicable).
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation
note	TEXT	0	Optional field to leave notes

## 1.5 Copy Number Somatic Mutations - Secondary Analysis File

Copy Number Somatic Mutations [cnsm] – Secondary Analysis File [s]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
mutation_id	VARCHAR(128)	R	Unique identifier for the mutation
gene_affected	VARCHAR(512)	R	Gene(s) containing the mutation/variation. Use Ensembl gene id. Separate multiple genes with vertical bars in the form of 'geneA—geneB—geneC'. If no gene is affected, use -888 (not applicable).
transcript_affected	TEXT	0	Transcript(s) containing the mutation/variation. Use Ensembl transcript id. Separate multiple transcripts from the same gene with commas, and separate transcripts from different genes with vertical bars. e.g.: 'transcriptA1,transcriptA2—transcriptB1—transcriptC1,transcriptC2,transcriptC3'. If no transcript is affected, use -888 (not applicable).
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation
note	TEXT	0	Optional field to leave notes

## 1.6 Methylation - Metadata File

Methylation [meth] – Metadata File [m]



Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
matched_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed control sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix_B10.tsv)
platform	VARCHAR(512)	R	Platform or technology used in detecting the methylation (See Table appendix_B5.tsv)
experimental_protocol	VARCHAR(512)	0	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
variation_calling_algorithm	VARCHAR(512)	R	Name of variation calling algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	0	Names of other analysis algorithms. Separate multiple algorithms by commas.
raw_data_repository	VARCHAR(128)	0	Public repository where raw data is submitted (See Table appendix_B12.tsv)
raw_data_accession	VARCHAR(128)	0	Accession and URL for referencing the raw data at the public repository
note	TEXT	0	Optional field to leave notes

## 1.7 Analyzed Sample Data File

Analyzed Sample Data File [sample] (required)

This submission file describes an analyzed sample on which molecular characterization was performed. It includes both control samples and tumour samples.

Name	Type	R/O	Description / Values
analyzed_sample_id	VARCHAR(64)	R	<i>Unique identifier</i> for the sample assigned by data provider
specimen_id	VARCHAR(64)	R	<i>Unique identifier</i> for the specimen assigned by data provider. The corresponding specimen id <b>must</b> appear in the specimen data submission file
analyzed_sample_type	VARCHAR(128)	R	Controlled vocabulary description of sample type (See Table <i>sp_analyzed_sample_type.tsv</i> )
analyzed_sample_type_other	VARCHAR(64)	0	Free text description of site of sample if "other" was specified in <i>sample_type</i> field
analyzed_sample_interval	INTEGER	0	Interval from specimen acquisition to sample use in an analytic procedure (e.g. DNA extraction), in days
analyzed_sample_notes	TEXT	0	Freetext notes about sample allowed

## 1.8 Donor Data File

Donor Data File [donor] (required)

This submission file describes a donor from which one or more specimens were obtained.

Name	Type	R/O	Description / Values
donor_id	VARCHAR(64)	R	<i>Unique identifier</i> for the donor; assigned by data provider.
donor_sex	VARCHAR(128)	R	Donor biological sex. <i>"Other" has been removed from the controlled vocabulary due to identifiability concerns. (See Table dr_donor_sex.tsv)</i>
donor_region_of_residence	VARCHAR(64)	R	Country, and optionally state or province code, but not city. <i>ISO3166-1-alpha-2 or ISO3166-2 codes, eg: "CA" or "CA-ON"</i>
donor_vital_status	VARCHAR(128)	R	Donor's last known vital status (See Table <i>dr_donor_vital_status.tsv</i> )
donor_age_at_diagnosis	INTEGER	R	Age at primary diagnosis <i>Use "90" for patients <math>\geq 90</math></i>
donor_diagnosis_do	VARCHAR(64)	R	Disease Ontology code <i>Disease                      Ontology                      code</i> ( <a href="http://diseaseontology.sourceforge.net/">http://diseaseontology.sourceforge.net/</a> )
donor_notes	TEXT	0	Free text notes concerning donor <i>Any additional non-identifying information can be included here.</i>

## 1.9 Expression - Metadata File

Expression [exp] – Metadata File [m]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix.B10.tsv)
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation
platform	VARCHAR(512)	R	Platform or technology used in detecting the expression (See Table appendix.B5.tsv)
experimental_protocol	VARCHAR(512)	O	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
normalization_algorithm	VARCHAR(512)	R	Name of normalization algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	O	Names of other analysis algorithms. Separate multiple algorithms by commas.
seq_coverage	FLOAT(5,2)	O	Sequence coverage if analyzed by sequencing platforms
raw_data_repository	VARCHAR(128)	O	Public repository where raw data is submitted (#) (See Table appendix.B12.tsv)
raw_data_accession	VARCHAR(128)	O	Accession and URL for referencing the raw data at the public repository
note	TEXT	O	Optional field to leave notes

## 1.10 Expression - Gene File

Expression [exp] – Gene File [g]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
gene_stable_id	VARCHAR(64)	R	For annotated gene, use Ensembl gene ID. Otherwise, use assemblyBuild_chr_start_end where assemblyBuild is hg18 or hg19.
gene_chromosome	VARCHAR(64)	R	Name of the chromosome containing the expressed gene/region interrogated (See Table appendix.B6.tsv)
gene_strand	INTEGER	R	Strand of the chromosome containing the expressed gene/region (See Table exp_g_gene_strand.tsv)
gene_start	INTEGER	R	Start position of the gene on the chromosome
gene_end	INTEGER	R	End position of the gene on the chromosome
normalized_read_count	FLOAT(5,2)	R	Normalized count of sequencing reads if analyzed by sequencing platforms
raw_read_count	INTEGER	R	Raw count of sequencing reads if analyzed by sequencing platforms
normalized_expression_level	FLOAT(5,2)	0	Normalized value of expression level if analyzed by microarray platforms
fold_change	FLOAT(5,2)	0	Expressed fold change if differential expression is measured
reference_sample	VARCHAR(64)	0	ID of the reference analyzed sample if differential expression is measured
quality_score	INTEGER	0	Quality score for the expression call
probability	FLOAT(3,2)	0	Probability of the expression call
is_annotated	VARCHAR(64)	0	Indicate if the expressed fragment is annotated in Ensembl (See Table exp_g_is_annotated.tsv)
validation_status	VARCHAR(64)	R	Indicate if the expressed fragment has been validated (See Table exp_g_validation_status.tsv)
validation_platform	VARCHAR(512)	0	Platform or technology used in validation (See Table appendix.B5.tsv)
probeset_id	VARCHAR(128)	0	ID of the probeset used in microarray
note	TEXT	0	Optional field to leave notes

## 1.11 Regulatory regions - Secondary Analysis File

Regulatory regions [rreg] – Secondary Analysis File [s]

Name	Type	R/O	Description / Values
analysis_id	TEXT	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	TEXT	R	Unique identifier for the analyzed sample
regulatory_region_id	TEXT	R	Unique identifier for the identified regulatory region
gene_affected	VARCHAR(128)	R	Gene(s) containing the methylation. Use Ensembl gene id. Separate multiple genes with vertical bars in the form of geneA—geneB—geneC. If no gene is affected, use -888 (not applicable).
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation
note	TEXT	O	Optional field to leave notes

## 1.12 Simple Somatic Mutations - Primary Analysis File

Simple Somatic Mutations [ssm] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
mutation_id	VARCHAR(128)	R	Unique identifier for the mutation
mutation_type	VARCHAR(128)	R	Type of mutation (See Table <i>ssm_p_mutation_type.tsv</i> )
chromosome	VARCHAR(64)	R	Name of the chromosome containing the mutation/variation (See Table <i>appendix.B6.tsv</i> )
chromosome_start	INTEGER	R	Start position of the mutation/variation on the chromosome
chromosome_end	INTEGER	R	End position of the mutation/variation on the chromosome
chromosome_strand	INTEGER	R	Chromosome strand (See Table <i>ssm_p_chromosome_strand.tsv</i> )
refsnp_allele	VARCHAR(512)	R	RefSNP alleles from dbSNP (use a dash for each missing base) e.g.: A/T, —/AAA
refsnp_strand	INTEGER	0	Strand of RefSNP allele (See Table <i>ssm_p_refsnp_strand.tsv</i> )
reference_genome_allele	VARCHAR(512)	R	Allele in the reference genome (use a dash for each missing base)
control_genotype	VARCHAR(512)	R	Genotype of the control sample (use a dash for each missing base)
tumour_genotype	VARCHAR(512)	R	Genotype of the tumour sample (use a dash for each missing base)
mutation	VARCHAR(1024)	R	Mutation, e.g. C > G
expressed_allele	VARCHAR(512)	0	The expressed allele(s) as revealed by RNA-seq, etc.
quality_score	INTEGER	0	Average quality score for the mutation/variation call
probability	FLOAT(3,2)	0	Probability of the mutation/variation call
read_count	FLOAT(5,2)	0	Average number of times the bases are covered by raw reads
is_annotated	VARCHAR(64)	0	Indicate if the mutation/variation is annotated in dbSNP (See Table <i>ssm_p_is_annotated.tsv</i> )
validation_status	VARCHAR(64)	R	Indicate if the mutation/variation has been validated (See Table <i>ssm_p_validation_status.tsv</i> )
validation_platform	VARCHAR(512)	0	Platform or technology used in validation (See Table <i>appendix.B5.tsv</i> )
xref_ensembl_var_id	VARCHAR(128)	0 <sup>14</sup>	Cross-reference: Ensembl Variation ID from Ensembl Variation database. e.g.: rs12345; ENSSNP53189
note	TEXT	0	Optional field to leave notes

## 1.13 Exon Junction - Metadata File

Exon Junction [jcn] – Metadata File [m]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
sample_id	TEXT	R	Unique identifier for the analyzed sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (#) (See Table <i>appendix.B10.tsv</i> )
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation
platform	VARCHAR(512)	R	Platform or technology used in detecting the expression (See Table <i>appendix.B5.tsv</i> )
experimental_protocol	VARCHAR(512)	O	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
normalization_algorithm	VARCHAR(512)	R	Name of normalization algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	O	Names of other analysis algorithms. Separate multiple algorithms by commas.
seq_coverage	FLOAT(5,2)	O	Sequence coverage if analyzed by sequencing platforms (See Table <i>jcn_m_seq_coverage.tsv</i> )
raw_data_repository	VARCHAR(128)	R	Public repository where raw data is submitted (#) (See Table <i>appendix.B12.tsv</i> )
raw_data_accession	VARCHAR(128)	R	Accession and URL for referencing the raw data at the public repository
note	TEXT	O	Optional field to leave notes

## 1.14 Structural Somatic Mutations - Secondary Analysis File

Structural Somatic Mutations [stsm] – Secondary Analysis File [s]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular group of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
sv_id	TEXT	R	Unique identifier for variant
placement	INTEGER	R	Ordering of breakpoint pairs within a single structural change event
bkpt_from_context	VARCHAR(64)	0	Contextual description of the first break location (Exonic, Intronic, Intergenic)
gene_affected_by_bkpt_from	VARCHAR(512)	0	"Gene(s) affected by the breakpoints. Use Ensembl gene id. Separate multiple genes with vertical bars in the form of geneA—geneB—geneC. If both breakpoints affect genes, then use ""—" to separate them. If no gene is affected, use -888 (not applicable)."
transcript_affected_by_bkpt_from	TEXT	0	Transcript(s) affected by the breakpoints. Use Ensembl transcript id. Separate multiple transcripts from the same gene with commas, and separate transcripts from different genes with vertical bars. e.g.: transcriptA1,transcriptA2—transcriptB1—transcriptC1
bkpt_to_context	VARCHAR(64)	0	Contextual description of the second break location (Exonic, Intronic, Intergenic)
gene_affected_by_bkpt_to	VARCHAR(512)	0	"Gene(s) affected by the breakpoints. Use Ensembl gene id. Separate multiple genes with vertical bars in the form of geneA—geneB—geneC. If both breakpoints affect genes, then use ""—" to separate them. If no gene is affected, use -888 (not applicable)."
transcript_affected_by_bkpt_to	TEXT	0	Transcript(s) affected by the breakpoints. Use Ensembl transcript id. Separate multiple transcripts from the same gene with commas, and separate transcripts from different genes with vertical bars. e.g.: transcriptA1,transcriptA2—transcriptB1—transcriptC1
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation
note	TEXT	0	Optional field to leave notes



## 1.15 Simple Germline Variations - Metadata File

Simple Germline Variations [sgv] – Metadata File [m]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular group of samples
control_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed control sample
matched_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed matched sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix_B10.tsv)
platform	VARCHAR(512)	R	Platform or technology used in detecting the mutation/variation (See Table appendix_B5.tsv)
experimental_protocol	VARCHAR(512)	O	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
variation_calling_algorithm	VARCHAR(512)	R	Name of variation calling algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	O	Names of other analysis algorithms. Separate multiple algorithms by commas.
seq_coverage	FLOAT(5,2)	O	Sequence coverage if analyzed by sequencing platforms
raw_data_repository	VARCHAR(512)	O	Public repository where raw data is submitted (#) (See Table appendix_B12.tsv)
raw_data_accession	VARCHAR(512)	O	Accession and URL for referencing the raw data at the public repository
note	TEXT	O	Optional field to leave notes

## 1.16 Specimen Data File

Specimen Data File [specimen] (required)

This submission file describes a specimen from which one or more samples were derived. Use additional rows for more than one specimen from the same patient. If more than one specimen was extracted during the same procedure, each gets a distinct ID.

Name	Type	R/O	Description / Values
donor_id	VARCHAR(64)	R	Unique identifier for the donor; assigned by data provider. It must be coded, and correspond to a donor ID listed in the donor data file.
specimen_id	VARCHAR(64)	R	Unique identifier for the specimen assigned by data provider.
specimen_type	VARCHAR(128)	R	Controlled vocabulary description of specimen type. (See Table <i>specimen__specimen_type.tsv</i> )
specimen_type_other	VARCHAR(64)	R	Free text description of site of specimen if "normal control (other)" or "tumour (other)" was specified in specimen_type field.
specimen_processing	VARCHAR(128)	R	Description of technique used to process specimen (See Table <i>specimen__specimen_processing.tsv</i> )
specimen_processing_other	VARCHAR(64)	R	If "other" specified for specimen_processing, may indicate technique here.
specimen_storage	VARCHAR(128)	R	Description of how specimen was stored. For specimens that were extracted freshly or immediately cultured, answer (1) "NA". (See Table <i>specimen__specimen_storage.tsv</i> )
specimen_storage_other	VARCHAR(64)	R	If "other" specified for specimen_storage, may indicate technique here.
specimen_biobank	VARCHAR(64)	R	If the specimen was obtained from a biobank, provide the biobank name here
specimen_biobank_id	VARCHAR(64)	R	If the specimen was obtained from a biobank, provide the biobank accession number here.
specimen_available	VARCHAR(128)	R	Whether additional tissue is available for followup studies. (See Table <i>specimen__specimen_available.tsv</i> )
specimen_notes	TEXT	0	Free text notes allowed <i>Any additional non-identifying information can be included here.</i>

## 1.17 Exon Junction - Primary Analysis File

Exon Junction [jcn] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular group of samples
sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
junction_id	VARCHAR(256)	R	For known exons, use exonID1_exonID2 where exonID1 and exonID2 are Ensembl IDs of the 5' and 3' exons, respectively. For novel or putative exons, use assemblyBuild_exon1chr_exon1end_exon2chr_exon2start where assemblyBuild is hg18 or hg19; exon1chr and exon2chr are the chromosomes of the 5' and 3' exons, respectively; exon1end is the end position of the 5' exon; exon2start is the start position of the 3' exon.
gene_stable_id	VARCHAR(64)	R	Stable ID of the gene containing the 5' exon at the junction. For annotated gene, use Ensembl gene ID. For putative and novel gene, use assemblyBuild_chr_start_end where assemblyBuild can be hg18 or hg19.
gene_chromosome	VARCHAR(64)	R	Name of the chromosome containing the above gene. <i>(See Table appendix.B6.tsv)</i>
gene_strand	INTEGER	R	Strand of the chromosome <i>(See Table jcn_p_gene_strand.tsv)</i>
gene_start	INTEGER	R	Start position of the entire gene on the chromosome as annotated in Ensembl
gene_end	INTEGER	R	End position of the entire gene on the chromosome as annotated in Ensembl
second_gene_stable_id	VARCHAR(64)	O	In the case of a fusion gene, provide the Stable ID of the gene containing the 3' exon at the junction. For annotated genes, use Ensembl gene ID. For putative and novel genes, use assemblyBuild_chr_start_end where assemblyBuild can be hg18 or hg19.
exon1_chromosome	VARCHAR(64)	R	Name of the chromosome containing the 5' exon (#) <i>(See Table appendix.B6.tsv)</i>
exon1_number_bases	INTEGER	R	Number of bases from 5' exon
exon1_end	INTEGER	R	End position of the 5' exon on the chromosome
exon1_strand	INTEGER	O	Chromosome strand of the 5' exon <i>(See Table jcn_p_exon1_strand.tsv)</i>
exon2_chromosome	VARCHAR(64)	R	Name of the chromosome containing the 3' exon (#) <i>(See Table appendix.B6.tsv)</i>
exon2_number_bases	INTEGER	R	Number of bases from 3' exon
exon2_start	INTEGER	R	Start position of the 3' exon on the chromosome
exon2_strand	INTEGER	O	Chromosome strand of the 3' exon <i>(See Table jcn_p_exon2_strand.tsv)</i>
is_fusion_gene	VARCHAR(16)	O	Indicate if the function is the result of a fusion gene <i>(See Table jcn_p_is_fusion_gene.tsv)</i>
is_novel_splice_form	VARCHAR(16)	O	Indicate if the splice form is novel <i>(See Table jcn_p_is_novel_splice_form.tsv)</i>
junction_seq	TEXT	O	Provide junction sequence if either is_fusion_gene or is_novel_splice_form is true
junction_type	VARCHAR(64)	O	Type of junction <i>(See Table jcn_p_junction_type.tsv)</i>

## 1.18 Donor Family History

Donor Family History [family] (optional)

This file describes the family history of the donor.

Name	Type	R/O	Description / Values
donor_id	TEXT	R	Unique identifier for the donor; assigned by data provider. It must be coded, and correspond to a donor ID listed in the donor data file.
relationship_type	VARCHAR(128)	R	Relationship to the donor (See Table <i>dr_family_history__relationship_type.tsv</i> )
relationship_type_other	TEXT	R	If "other" answered in previous column, specify the relationship type here
relationship_sex	VARCHAR(128)	R	Biological sex of related individual (See Table <i>dr_family_history__relationship_sex.tsv</i> )
relationship_age	INTEGER	R	Age of relative at primary diagnosis (years) Use 90 for ages $\geq 90$ years.
relationship_diagnosis_do	TEXT	R	Disease Ontology code for the relative's diagnosis status
relationship_diagnosis	TEXT	R	Diagnosis (disease or healthy status) <i>e.g. "breast cancer"</i>

## 1.19 Copy Number Somatic Mutations - Primary Analysis File

Copy Number Somatic Mutations [cnsm] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
mutation_id	VARCHAR(128)	R	Unique identifier for the mutation
mutation_type	VARCHAR(128)	R	Type of mutation (See Table <i>cnsmp__mutation_type.tsv</i> )
chromosome	VARCHAR(64)	R	Name of the chromosome containing the mutation/variation (#) (See Table <i>appendix.B6.tsv</i> )
chromosome_start	INTEGER	R	Start position of the mutation/variation on the chromosome
chromosome_end	INTEGER	R	End position of the mutation/variation on the chromosome
chromosome_start_range	INTEGER	0	Number of bases around chromosome_start that may contain the start position
chromosome_end_range	INTEGER	0	Number of bases around chromosome_end that may contain the end position
start_probe_id	VARCHAR(128)	0	Probe id containing the chromosome_start if array platform was used
end_probe_id	VARCHAR(128)	0	Probe id containing the chromosome_end if array platform was used
copy_number	FLOAT(3,2)	0	DNA copy number estimated
segment_mean	FLOAT(5,2)	0	Mean LRR per segment
segment_median	FLOAT(5,2)	0	Median LRR per segment
quality_score	FLOAT(5,2)	0	Quality score for the mutation/variation call
probability	FLOAT(3,2)	0	Probability of the mutation/variation call
is_annotated	VARCHAR(64)	0	Indicate if the mutation/variation is annotated in the Database of Genomic Variants (See Table <i>cnsmp__is_annotated.tsv</i> )
validation_status	VARCHAR(64)	R	Indicate if the mutation/variation has been validated (See Table <i>cnsmp__validation_status.tsv</i> )
validation_platform	VARCHAR(512)	0	Platform or technology used in validation (See Table <i>appendix.B5.tsv</i> )
note	TEXT	0	Optional field to leave notes

## 1.20 Structural Somatic Mutations - Primary Analysis File

Structural Somatic Mutations [stsm] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
sv_id	VARCHAR(32)	R	Unique identifier for variant
placement	INTEGER	R	Ordering of breakpoint pairs within a single structural mutation/variation event
annotation	VARCHAR(256)	R	Annotation describing sequence mutation/variation based on breakpoint pairs
interpreted_annotation	VARCHAR(256)	O	HGVS nomenclature for description of sequence mutation/variation. <i>e.g.: chr3:g.1234567-2345678inv</i>
variant_type	VARCHAR(128)	R	Type of mutation/variation (See Table appendix_B9.tsv)
chr_from	VARCHAR(64)	R	Name of the donor chromosome containing the mutation/variation (See Table appendix_B6.tsv)
chr_from_bkpt	INTEGER	R	Breakpoint position of the mutation/variation on the donor chromosome
chr_from_strand	INTEGER	R	Donor chromosome strand (See Table stsm_p_chr_from_strand.tsv)
chr_from_range	INTEGER	O	Number of bases around chr_from_bkpt that may contain the real breakpoint
chr_from_flanking_seq	VARCHAR(512)	O	Flanking sequences that are 200bp upstream and 200bp downstream to the chr_from_bkpt position.
chr_to	VARCHAR(64)	R	Name of the acceptor chromosome containing the mutation/variation (See Table appendix_B6.tsv)
chr_to_bkpt	INTEGER	R	Breakpoint position of the mutation/variation on the acceptor chromosome
chr_to_strand	INTEGER	R	Acceptor chromosome strand (See Table stsm_p_chr_to_strand.tsv)
chr_to_range	INTEGER	O	Number of bases around chr_to_bkpt that may contain the real breakpoint
chr_to_flanking_seq	VARCHAR(512)	O	Flanking sequences that are 200bp upstream and 200bp downstream to the chr_to_bkpt position
microhomology_sequence	TEXT	O	If a microhomology is inserted, provide sequence
non_templated_sequence	TEXT	O	If non-templated DNA is inserted, provide sequence
evidence	VARCHAR(128)	O	Evidence supporting a structural mutation/variation (See Table stsm_p_evidence.tsv)
quality_score	INTEGER	O	Quality score for the mutation/variation call
probability	FLOAT(3,2)	O <sup>22</sup>	Probability of the mutation/variation call
zygosity	VARCHAR(64)	O	Zygosity (See Table stsm_p_zygosity.tsv)
validation_status	VARCHAR(64)	R	Indicate if the mutation/variation has been validated (See Table stsm_p_validation_status.tsv)

## 1.21 Protein-DNA interaction - Secondary Analysis File

Protein-DNA [pdna] – Secondary Analysis File [s]

Name	Type	R/O	Description / Values
analysis_id	TEXT	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	TEXT	R	Unique identifier for the analyzed sample
protein_dna_interaction_id	TEXT	R	Unique identifier for the protein-DNA interaction
gene_affected	VARCHAR(128)	R	Gene(s) containing the methylation. Use Ensembl gene id. Separate multiple genes with vertical bars in the form of geneA—geneB—geneC. If no gene is affected, use -888 (not applicable).
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation
note	TEXT	0	Optional field to leave notes

## 1.22 Simple Germline Variations - Primary Analysis File

Simple Germline Variations [sgv] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	TEXT	R	Unique identifier for the analysis performed for a particular set of samples
control_sample_id	TEXT	R	Unique identifier for the analyzed control sample
variation_id	VARCHAR(128)	R	Unique identifier for the variation
variation_type	VARCHAR(64)	R	Type of variation (See Table <i>sgv_p__variation_type.tsv</i> )
chromosome	VARCHAR(64)	R	Name of the chromosome containing the mutation/variation (See Table <i>appendix.B6.tsv</i> )
chromosome_start	INTEGER	R	Start position of the mutation/variation on the chromosome
chromosome_end	INTEGER	R	End position of the mutation/variation on the chromosome
chromosome_strand	INTEGER	R	Chromosome strand (See Table <i>sgv_p__chromosome_strand.tsv</i> )
refsnp_allele	VARCHAR(512)	R	RefSNP alleles from dbSNP (use a dash for each missing base) <i>e.g.: A/T, —/AAA</i>
refsnp_strand	INTEGER	O	Strand of RefSNP allele (See Table <i>sgv_p__refsnp_strand.tsv</i> )
reference_genome_allele	VARCHAR(512)	R	Allele in the reference genome (use a dash for each missing base)
control_genotype	VARCHAR(512)	R	Genotype of the control sample (use a dash for each missing base)
tumour_genotype	VARCHAR(512)	R	Genotype of the tumour sample (use a dash for each missing base)
expressed_allele	VARCHAR(512)	O	The expressed allele(s) as revealed by RNA-seq, etc.
quality_score	INTEGER	O	Average quality score for the mutation/variation call
probability	FLOAT(3,2)	O	Probability of the mutation/variation call
read_count	FLOAT(5,2)	O	Average number of times the bases are covered by raw reads
is_annotated	VARCHAR(64)	O	Indicate if the mutation/variation is annotated in dbSNP (See Table <i>sgv_p__is_annotated.tsv</i> )
validation_status	VARCHAR(64)	R	Indicate if the mutation/variation has been validated (See Table <i>sgv_p__validation_status.tsv</i> )
validation_platform	VARCHAR(512)	O	Platform or technology used in validation (See Table <i>appendix.B5.tsv</i> )
xref_ensembl_var_id	VARCHAR(128)	O	Cross-reference: Ensembl Variation ID in Ensembl Variation database. <i>e.g.: rs12345; ENSNP53189</i>
note	TEXT	O	Optional field to leave notes



## 1.23 Simple Somatic Mutations - Metadata File

Simple Somatic Mutations - Metadata File

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular group of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
matched_sample_id	VARCHAR(64)	R	Unique identifier for the control sample matched to the analyzed sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix.B10.tsv)
platform	VARCHAR(512)	R	Platform or technology used in detecting the mutation/variation (See Table appendix.B5.tsv)
experimental_protocol	VARCHAR(512)	O	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
variation_calling_algorithm	VARCHAR(512)	R	Name of variation calling algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	O	Names of other analysis algorithms. Separate multiple algorithms by commas.
seq_coverage	FLOAT(5,2)	O	Sequence coverage if analyzed by sequencing platforms
raw_data_repository	VARCHAR(128)	O	Public repository where raw data is submitted (#) (See Table appendix.B12.tsv)
raw_data_accession	VARCHAR(128)	O	Accession and URL for referencing the raw data at the public repository
note	TEXT	O	Optional field to leave notes

## 1.24 Methylation - Primary Analysis File

Methylation [meth] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	TEXT	R	Unique identifier for the analysis performed for a particular group of samples
analyzed_sample_id	TEXT	R	Unique identifier for the analyzed sample
methyated_fragment_id	VARCHAR(128)	R	Unique identifier for the methylated fragment
chromosome	VARCHAR(64)	R	Name of the chromosome containing the methylation (See Table appendix_B6.tsv)
chromosome_start	INTEGER	R	Start position of the methylation on the chromosome
chromosome_end	INTEGER	R	End position of the methylation on the chromosome
chromosome_strand	INTEGER	0	Chromosome strand (See Table meth_p__chromosome_strand.tsv)
beta_value_methylation	FLOAT(5,2)	0	Methylation Beta value for interrogated site
beta_value_hydroxymethylation	FLOAT(5,2)	0	Hydroxymethylation Beta value for interrogated site
quality_score_methylation	INTEGER	0	Quality score for the methylation call
quality_score_hydroxymethylation	INTEGER	0	Quality score for the hydroxymethylation call
probability	FLOAT(3,2)	0	Probability of the methylation call
validation_status	VARCHAR(64)	R	Indicate if the methylation has been validated (See Table meth_p__validation_status.tsv)
validation_platform	VARCHAR(512)	0	Platform or technology used in validation (See Table appendix_B5.tsv)
note	TEXT	0	Optional field to leave notes

## 1.25 Simple Somatic Mutations - Secondary Analysis File

Simple Somatic Mutations [ssm] – Secondary Analysis File [s]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
mutation_id	VARCHAR(128)	R	Unique identifier for the mutation
consequence_type	VARCHAR(512)	R	Functional consequence of the SNP. (See Table appendix_B7.tsv)
aa_mutation	VARCHAR(512)	0	Changes at amino acid level. Indicate the reference aa, position and mutation aa. <i>e.g.: P234W</i>
cds_mutation	VARCHAR(512)	0	Changes in coding sequence. Indicate position, reference base and mutation base. <i>e.g.: 12324T&gt;G</i>
protein_domain_affected	VARCHAR(128)	0	Protein domain containing the mutation/variation. Use Pfam accession.
gene_affected	VARCHAR(512)	0	Gene(s) containing the mutation/variation.
transcript_affected	TEXT	0	Transcript(s) containing the mutation/variation. Use Ensembl transcript id.
gene_build_version	INTEGER	R	Version of Ensembl gene build used for annotation <i>e.g.: 61</i>
note	TEXT	0	Optional field to leave notes

## 1.26 Copy Number Somatic Mutations - Metadata File

Copy Number Somatic Mutations [cnsm] – Metadata File [m]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular group of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
matched_sample_id	VARCHAR(64)	R	Unique identifier for the control sample matched to the analyzed sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix_B10.tsv)
platform	VARCHAR(512)	R	Platform or technology used in detecting the mutation/variation (See Table appendix_B5.tsv)
experimental_protocol	VARCHAR(512)	0	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	0	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	0	Name of alignment algorithm and URL to written protocol
variation_calling_algorithm	VARCHAR(512)	0	Name of variation calling algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	0	Names of other analysis algorithms. Separate multiple algorithms by commas.
seq_coverage	FLOAT(3,2)	0	Sequence coverage if analyzed by sequencing platforms
raw_data_repository	VARCHAR(128)	0	Public repository where raw data is submitted (See Table appendix_B12.tsv)
raw_data_accession	VARCHAR(128)	0	Accession and URL for referencing the raw data at the public repository
note	TEXT	0	Optional field to leave notes

## 1.27 Protein-DNA interaction - Primary Analysis File

Protein-DNA [pdna] – Primary Analysis File [p]

Name	Type	R/O	Description / Values
analysis_id	TEXT	R	Unique identifier for the analysis performed for a particular group of samples
analyzed_sample_id	TEXT	R	Unique identifier for the analyzed sample
protein_dna_interaction_id	VARCHAR(128)	R	Unique identifier for the protein-DNA interaction
protein_stable_id	VARCHAR(128)	R	Ensembl Protein stable id of the interacting protein
chromosome	VARCHAR(64)	R	Name of the chromosome containing the methylation (See Table <i>appendix.B6.tsv</i> )
chromosome_start	INTEGER	R	Start position of the methylation on the chromosome
chromosome_end	INTEGER	R	End position of the methylation on the chromosome
chromosome_strand	INTEGER	O	Chromosome strand (See Table <i>meth_p_chromosome_strand.tsv</i> )
idr	FLOAT(5,2)	R	Irreproducible discovery rate
fdr	FLOAT(5,2)	O	False discovery rate
rank_type	VARCHAR(64)	O	Kind of used ranking
rank_value	FLOAT(5,2)	O	Rank value
validation_status	VARCHAR(64)	R	Indicate if the methylation has been validated (See Table <i>meth_p_validation_status.tsv</i> )
validation_platform	VARCHAR(512)	O	Platform or technology used in validation (See Table <i>appendix.B5.tsv</i> )
note	TEXT	O	Optional field to leave notes

## 1.28 Structural Somatic Mutations - Metadata File

Structural Somatic Mutations [stsm] – Metadata File [m]

Name	Type	R/O	Description / Values
analysis_id	VARCHAR(64)	R	Unique identifier for the analysis performed for a particular set of samples
analyzed_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed sample
matched_sample_id	VARCHAR(64)	R	Unique identifier for the analyzed control sample
assembly_version	VARCHAR(64)	R	Version of reference genome assembly (See Table appendix_B10.tsv)
platform	VARCHAR(512)	R	Platform or technology used in detecting the mutation/variation (See Table appendix_B5.tsv)
experimental_protocol	VARCHAR(512)	O	Name of experimental protocol and URL to written protocol
base_calling_algorithm	VARCHAR(512)	R	Name of base calling algorithm and URL to written protocol
alignment_algorithm	VARCHAR(512)	R	Name of alignment algorithm and URL to written protocol
variation_calling_algorithm	VARCHAR(512)	R	Name of variation calling algorithm and URL to written protocol
other_analysis_algorithm	VARCHAR(512)	O	Names of other analysis algorithms. Separate multiple algorithms by Commas.
seq_coverage	FLOAT(5,2)	O	Sequence coverage if analyzed by sequencing platforms
raw_data_repository	VARCHAR(128)	O	Public repository where raw data is submitted (See Table appendix_B12.tsv)
raw_data_accession	VARCHAR(128)	O	Accession and URL for referencing the raw data at the public repository
note	TEXT	O	Optional field to leave notes

# Appendix A

## CV Tables

### A.1 CV Table appendix\_B10.tsv

Key	Description
1	GRCh37
2	NCBI36
3	GRCh37.p1
4	GRCh37.p2
5	GRCh37.p3
6	GRCh37.p4
7	GRCh37.p5

## A.2 CV Table appendix\_B5.tsv

Key	Description
1	PCR
2	qPCR
3	capillary sequencing
4	SOLiD sequencing
5	Illumina GA sequencing
6	454 sequencing
7	Helicos sequencing
8	Affymetrix Genome-Wide Human SNP Array 6.0
9	Affymetrix Genome-Wide Human SNP Array 5.0
10	Affymetrix Mapping 100K Array Set
11	Affymetrix Mapping 500K Array Set
12	Affymetrix Mapping 10K 2.0 Array Set
13	Affymetrix EMET Plus Premier Pack
14	Agilent Whole Human Genome Oligo Microarray Kit
15	Agilent Human Genome 244A
16	Agilent Human Genome 105A
17	Agilent Human CNV Association 2x105K
18	Agilent Human Genome 44K
19	Agilent Human CGH 1x1M
20	Agilent Human CGH 2x400K
21	Agilent Human CGH 4x180K
22	Agilent Human CGH 8x60K
23	Agilent Human CNV 2x400K
24	Agilent Human miRNA Microarray Kit (v2)
25	Agilent Human CpG Island Microarray Kit
26	Agilent Human Promoter ChIP-on-chip Microarray Set
27	Agilent Human SpliceArray
28	Illumina human1m-duo
29	Illumina human660w-quad
30	Illumina humancytosnp-12
31	Illumina human510s-duo
32	Illumina humanmethylation27
33	Illumina goldengate methylation
34	Illumina HumanHT-12 v4.0 beadchip
35	Illumina HumanWG-6 v3.0 beadchip
36	Illumina HumanRef-8 v3.0 beadchip
37	Illumina microRNA Expression Profiling Panel
38	Illumina humanht-16
39	Illumina humanht-17
40	Nimblegen Human CGH 3x720 Whole-Genome v3.0 Array
41	Nimblegen Human CGH 2.1M Whole-Genome v2.0D Array
42	Nimblegen Gene Expression 385K
43	Nimblegen Gene Expression 4x72K
44	Nimblegen Gene Expression 12x135K
45	Nimblegen Human Methylation 2.1M Whole-Genome sets
46	Nimblegen Human Methylation 385K Whole-Genome sets
47	Nimblegen CGS
48	Illumina Human1M OmniQuad chip
49	PCR and capillary sequencing
50	Custom-designed gene expression array
51	Affymetrix HT Human Genome U133A Array Plate Set
52	Agilent 244K Custom Gene Expression G4502A-07-1
53	Agilent 244K Custom Gene Expression G4502A-07-2
54	Agilent 244K Custom Gene Expression G4502A-07-3
55	Agilent Human Genome CGH Custom Microarray 2x415K
56	Affymetrix Human U133 Plus PM
57	Affymetrix Human U133 Plus 2.0
58	Affymetrix Human Exon 1.0 ST
59	Almac Human CRC



### A.3 CV Table appendix\_B12.tsv

Key	Description
1	EGA
2	dbSNP
3	TCGA
4	CGHub
5	GEO

## A.4 CV Table appendix\_B6.tsv

Key	Description
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	21
22	22
23	X
24	Y
25	MT
26	c5.H2
27	c6.COX
28	c6.QBL
29	NT_113870
30	NT_113871
31	NT_113872
32	NT_113874
33	NT_113878
34	NT_113880
35	NT_113881
36	NT_113884
37	NT_113885
38	NT_113886
39	NT_113888
40	NT_113889
41	NT_113890
42	NT_113898
43	NT_113899
44	NT_113901
45	NT_113902
46	NT_113903
47	NT_113906
48	NT_113908
49	NT_113909
50	NT_113910
51	NT_113911
52	NT_113912
53	NT_113915
54	NT_113916
55	NT_113917
56	NT_113923
57	NT_113924
58	NT_113925
59	NT_113926
60	NT_113927

## A.5 CV Table meth\_p\_\_chromosome\_strand.tsv

Key	Description
1	1
-1	1

## A.6 CV Table meth\_p\_\_validation\_status.tsv

Key	Description
1	validated
2	not tested
3	not valid

## A.7 CV Table sp\_\_analyzed\_sample\_type.tsv

Key	Description
1	Normal blood
2	Leukemic blood
3	Normal control adjacent to primary
4	Normal control from non-tumour site
5	Control from cell line derived from normal tissue
6	Normal mouse host
7	Primary tumour
8	Mouse xenograft derived from tumour
9	Cell line derived from tumour
10	Cell line derived from xenograft
11	Other (specify)

## A.8 CV Table dr\_\_donor\_sex.tsv

Key	Description
1	male
2	female

## A.9 CV Table dr\_\_donor\_vital\_status.tsv

Key	Description
1	alive
2	deceased

## A.10 CV Table exp\_g\_\_gene\_strand.tsv

Key	Description
1	1
-1	-1

## A.11 CV Table exp\_g\_\_is\_annotated.tsv

Key	Description
1	annotated
2	not annotated

### A.12 CV Table exp\_g\_\_validation\_status.tsv

Key	Description
1	validated
2	not tested
3	not valid

### A.13 CV Table ssm\_p\_\_mutation\_type.tsv

Key	Description
1	single base substitution
2	insertion of <=200bp
3	deletion of <=200bp
4	multiple base substitution (>=2bp and <=200bp)

### A.14 CV Table ssm\_p\_\_chromosome\_strand.tsv

Key	Description
1	1
-1	-1

### A.15 CV Table ssm\_p\_\_refsnp\_strand.tsv

Key	Description
1	1
-1	-1

### A.16 CV Table ssm\_p\_\_is\_annotated.tsv

Key	Description
1	annotated
2	not annotated

### A.17 CV Table ssm\_p\_\_validation\_status.tsv

Key	Description
1	validated
2	not tested
3	not valid

### A.18 CV Table jcn\_m\_\_seq\_coverage.tsv

Key	Description
1	EGA
2	dbSNP

### A.19 CV Table specimen\_\_specimen\_type.tsv

Key	Description
1	primary tumour
2	tumour local recurrence
3	tumour metastasis to local lymph node
4	tumour metastasis to distant location
5	peripheral blood
6	bone marrow
7	lymph node
8	normal control (tissue adjacent to primary)
9	normal control (blood)
10	normal control (other)
11	tumour (other)

### A.20 CV Table specimen\_\_specimen\_processing.tsv

Key	Description
1	cryopreservation in liquid nitrogen (dead tissue)
2	cryopreservation in dry ice (dead tissue)
3	cryopreservation of live cells in liquid nitrogen
4	cryopreservation, other
5	formalin fixed, unbuffered
6	formalin fixed, buffered
7	formalin fixed & paraffin embedded
8	fresh
9	other technique

### A.21 CV Table specimen\_\_specimen\_storage.tsv

Key	Description
1	frozen, liquid nitrogen
2	frozen, -70 freezer
3	frozen, vapor phase
4	RNA later frozen
5	paraffin block
6	cut slide
7	other

### A.22 CV Table specimen\_\_specimen\_available.tsv

Key	Description
1	no
2	yes

### A.23 CV Table jcn\_p\_\_gene\_strand.tsv

Key	Description
1	1
-1	-1

#### A.24 CV Table jcn\_p\_\_exon1\_strand.tsv

Key	Description
1	1
-1	-1

#### A.25 CV Table jcn\_p\_\_exon2\_strand.tsv

Key	Description
1	1
-1	-1

#### A.26 CV Table jcn\_p\_\_is\_fusion\_gene.tsv

Key	Description
1	yes
2	no

#### A.27 CV Table jcn\_p\_\_is\_novel\_splice\_form.tsv

Key	Description
1	yes
2	no

#### A.28 CV Table jcn\_p\_\_junction\_type.tsv

Key	Description
1	canonical
2	non-canonical
3	U12

#### A.29 CV Table jcn\_p\_\_validation\_status.tsv

Key	Description
1	validated
2	not tested
3	not valid

#### A.30 CV Table dr\_family\_history\_\_relationship\_type.tsv

Key	Description
1	sibling
2	parent
3	grandparent
4	uncle/aunt
5	cousin
6	other

### A.31 CV Table dr\_family\_history\_relationship\_sex.tsv

Key	Description
1	male
2	female

### A.32 CV Table cns\_m\_p\_mutation\_type.tsv

Key	Description
1	gain
2	loss
3	copy neutral LOH
4	copy neutral
5	hemizygous del LOH
6	amp LOH

### A.33 CV Table cns\_m\_p\_is\_annotated.tsv

Key	Description
1	annotated
2	not annotated

### A.34 CV Table cns\_m\_p\_validation\_status.tsv

Key	Description
1	validated
2	not tested
3	not valid

### A.35 CV Table appendix\_B9.tsv

Key	Description
1	deletion
2	tandem duplication
3	inversion
4	inverted duplication - head-to-head
5	inverted duplication - tail-to-tail
6	insertion
7	intrachromosomal rearrangement with inverted orientation
8	intrachromosomal rearrangement with non-inverted orientation
9	fold-back inversion
10	complex intrachromosomal rearrangement
11	reciprocal translocation
12	unbalanced translocation
13	interchromosomal insertion
14	interchromosomal rearrangement - unknown type
15	complex interchromosomal rearrangement
16	intrachromosomal amplicon-to-amplicon
17	intrachromosomal amplicon-to-nonamplified dna
18	interchromosomal amplicon-to-amplicon
19	interchromosomal amplicon-to-nonamplified dna
20	extrachromosomal

### A.36 CV Table stsm\_p\_\_chr\_from\_strand.tsv

Key	Description
1	1
-1	-1

### A.37 CV Table stsm\_p\_\_chr\_to\_strand.tsv

Key	Description
1	1
-1	-1

### A.38 CV Table stsm\_p\_\_evidence.tsv

Key	Description
1	Copy number change
2	FISH
3	Flow-sorted chromosome evidence
4	Paired sequence either side of breakpoint
5	Partner breakpoint found
6	PCR product across breakpoint
7	Protein evidence
8	Seen in multiple samples
9	Sequence across breakpoint

### A.39 CV Table stsm\_p\_\_zygosity.tsv

Key	Description
1	homozygous
2	heterozygous
3	hemizygous
4	nullizygous

### A.40 CV Table stsm\_p\_\_validation\_status.tsv

Key	Description
1	validated
2	not tested
3	not valid

### A.41 CV Table sgvs\_p\_\_variation\_type.tsv

Key	Description
1	single base substitution
2	insertion of $\leq 200$ bp
3	deletion of $\leq 200$ bp
4	multiple base substitution ( $\geq 2$ bp and $\leq 200$ bp)



#### A.42 CV Table `sgv_p__chromosome_strand.tsv`

Key	Description
1	1
-1	-1

#### A.43 CV Table `sgv_p__refsnp_strand.tsv`

Key	Description
1	1
-1	-1

#### A.44 CV Table `sgv_p__is_annotated.tsv`

Key	Description
1	annotated
2	not annotated

#### A.45 CV Table `sgv_p__validation_status.tsv`

Key	Description
1	validated
2	not tested
3	not valid

## A.46 CV Table appendix\_B7.tsv

Key	Description
1	3prime_utr
2	5prime_utr
3	upstream
4	downstream
5	essential_splice_site,3prime_utr
6	essential_splice_site,5prime_utr
7	essential_splice_site,intronic
8	essential_splice_site,non_synonymous_coding
9	essential_splice_site,stop_lost
10	essential_splice_site,synonymous_coding
11	frameshift_coding
12	frameshift_coding,splice_site
13	intergenic
14	intronic
15	non_synonymous_coding
16	non_synonymous_coding,splice_site
17	splice_site,3prime_utr
18	splice_site,5prime_utr
19	splice_site,intronic
20	splice_site,synonymous_coding
21	stop_gained
22	stop_gained,splice_site
23	stop_lost
24	stop_lost,splice_site
25	synonymous_coding
26	utr
27	splice_site
28	noncoding_rna
29	complex_indel
30	regulatory_region
31	inframe_indel
32	start_lost
33	ambiguous
34	complex_substitution