



肿瘤事业部 Control 集合建立及更新方案

撰写人		日 期	
/Draft :	刘博	/Date :	2020.09
审核人		日 期	
/Review :		/Date :	
批准人		日 期	
/Approve :		/Date :	
文件密级/	<input type="checkbox"/> 普通/Unclassified <input type="checkbox"/> 秘密/Secret		
Classification :	<input checked="" type="checkbox"/> 机密/Highly Secret <input type="checkbox"/> 绝密/Top Secret		

目录

1	目的 (Objectives)	1
2	适用范围 (Scopes)	1
3	职责 (Responsibilities)	1
4	术语和定义 (Terms and definitions)	1
5	数据来源和数量 (Data Sources)	2
6	集合构建 (Data Analysis)	2
6.1.	Control 集合构建流程获取	2
6.2.	数据分析流程	2
6.3.	集合构建结果	4
6.4.	集合使用方法	5
7	相关记录 (Related Records)	6
8	参考文献资料 (References)	6
9	附录 (Appendix)	6

1 目的 (Objectives)

经过生信部门的内部讨论，确定了肿瘤事业部 A0 版本的 Control 集合构建标准，制订本方案是为保证明确记录肿瘤事业部内部信息分析流程的 Control 集合构建流程，并在事业部内部形成共识，规范化事业部内部各个产品的 Control 集合的构建标准，并建立统一的规范，形成标准化管理。同时也为后期可能的数据追溯，Control 集合数据集更新建立参考文档，指导后期相关方案的落实。

2 适用范围 (Scopes)

本方案适用于肿瘤事业部内部所有需要进行 Control 集合构建的生信分析流程，如无特殊说明，肿瘤事业部内部所有 Control 集合均基于本方案进行构建，后期如果在使用过程中发现了本方案的缺陷或不足，则应该对本方案方法进行优化及迭代说明，并对所有基于本方案生成的 Control 集合于指定时间内进行同步更新。

3 职责 (Responsibilities)

3.1 生信部负责本方案的整体撰写及第一版脚本的梳理准备工作。

3.2 生信部负责人负责本方案的审核工作。

3.3 生信部各个产品生信分析流程负责人员负责本方案的执行，并生成相应的 Control 集合应用于各个产品对应的生产流程。生信所有相关人员负责使用及了解本方案的使用方法 & 原理，如果因发现缺陷不足需要对本方案进行升级，则应通知生信部门全员 (bgi_tumor@genomics.cn) 知晓相关调整，从而确保各个产品的 Control 集合可以得到及时更新并保持构建逻辑的一致性。

4 术语和定义 (Terms and definitions)

Control 集合：基于华大基因（简称“华大”，下同）接收的历史临床样本的白细胞数据构建临床收样患者在人群层面表现出来的多态性位点（人群高发，大概率不致病）和检测体系的假阳性位点（散发低频变异，体系偏好性导致），用于后期对患者变异结果进行过滤从而剔除一些和患者的肿瘤发生发展无关的变异集合。

ErrorBaseLine：基于历史检出样本，构建位点假阳性错误波动范围文件，各个产品由于上游在实验、测序平台的差异，因此会存在不同的错误偏好性。不同的产品应该独立的进行本集合的构建。

Vcf：变异检测结果文件格式，用于记录样本的变异检出结果。

批注 [刘博(Bo1)]: 和 ErrorBaseLine 似乎是一样的功能，是否有必要保留。

热点变异：基于解读收录整理的一些已知致病变异，这类变异由于已知和肿瘤的发生发展存在相关性，因此需要从 Control 集合中进行过滤，从而避免造成临床解读层面的假阴性。

5 数据来源和数量 (Data Sources)

为了更好的反馈检测系统的偏差因此构建 Control 集合所用的**变异结果(原始 VCF 文件) 应和产品信息分析流程保持一致**，从而避免分析流程中存在的差异对 Control 集合带来偏好性的影响。

构建 Control 集合时，应该优先使用真实的临床样本，同时避免癌种的过渡集中（癌种的过渡集中会导致该癌种中的热点变异被富集，被错误的识别为人群多态性位点，尤其是一些暂未被广泛研究的隐藏热点变异被过滤）。

样本量应该尽可能的多，从而更准确的反应人群的分布情况，减少偶然性带来的影响。建议样本量在 1000+以上。如果转产初期样本量少，建议在样本量分别达到 100 例、500 例、1000 例时，对 Control 集分别进行定期更新。同时由于样本量较少的情况下，为确保变异不会因为偶然性出现非多态性位点的富集，因此在样本量较少（500 例以下时），建议人群频率使用 10%进行过滤。

6 集合构建 (Data Analysis)

6.1. 构建流程获取

为了便于整个生信分析流程的管理，所有 Control 集合构建流程均保存在华大 Gitlab，生信工具集仓库中（bioinfokit: <https://gitlab.genomics.cn/liubo4/bioinfokit/-/tree/master/>），如有权限问题，联系（liubo4@genomics.cn）

首先下载 git 仓库 `git clone ssh://git@gitlab.genomics.cn:2200/liubo4/bioinfokit.git`

下载完成后目录 `bioinfokit/02.toolkit/04.Create_ControlSite` 中对应的即为 Control 集合构建的流程。使用方式参考下文及示例 shell: `02.toolkit/04.Create_ControlSite/work.demo.sh`

下载完成后目录 `02.toolkit/04.Create_ErrorBaseLine` 中对应的即为 ErrorBaseLine 集合构建流程，使用方式参考下文及示例 shell: `02.toolkit/04.Create_ErrorBaseLine/demo.sh`

6.2. Control 集合构建 - 基于血细胞检测结果

Control 集合构建参考如下命令：

```
perl Create_Control_Database.pl -l demoInput.vcf.list -h Create_Control_Database.HotMutList
```

-c 0.05



肿瘤事业部 Control 集合建立及更新方案

文件编号：R-BIN-050

版本号：A0

第3页 共4页

其中：

Create_Control_Database.pl 为 Control 集合构建主流程

demoInput.vcf.list 为所有用于构建 Control 集合的 vcf 文件 list，**注意由于 vcf 格式不一致，需要确定 vcf 文件中 AF 列是否存在，若不存在或 AF 所在列不一致，需要对流程进行调整**

输入的 list 文件示例如下：

```
(base) b2c_v03_pipeline@t1-jdpc-24-4[Wed Oct 14] /jdfst1/B2C_CM_P1/PipeAdmin/04_Pipeline/bin/afcol1/02_toolkit/04_Create_ControlGis
$ head demoInput.vcf.list
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM503_ljyanan_2051085371_2061085371_somatic.sn.vcf
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM503_lizhichang_1963201047_1963201047_Analyze/sn/DM503_lizhichang_1963201047_somatic.sn.vcf
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM503_lulongying_2056955282_2066955282_Analyze/sn/DM503_lulongying_2056955282_somatic.sn.vcf
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM503_zhengshuneng_2066955282_2066955282_Analyze/sn/DM503_zhengshuneng_2066955282_somatic.sn.vcf
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM503_zhengshuneng_2066955282_2066955282_Analyze/sn/DM503_zhengshuneng_2066955282_somatic.sn.vcf
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM503_zhouyang_1963557157_1963557157_Analyze/sn/DM503_zhouyang_1963557157_somatic.sn.vcf
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM502_zhuchenglong_1963557157_1963557157_Analyze/sn/DM502_zhuchenglong_1963557157_somatic.sn.vcf
/jdfst1/B2C_CM_P1/Clinical_Product/pipeline/pancancer/20200915/pancancer688_DM502_zhuchenglong_1963557157_1963557157_Analyze/sn/DM502_zhuchenglong_1963557157_somatic.sn.vcf
```

本流程构建所用 vcf 文件格式示意如下：

```
##contig=<ID=chr3,length=198022430>
##contig=<ID=chr4,length=191154276>
##contig=<ID=chr5,length=109515266>
##contig=<ID=chr6,length=171115067>
##contig=<ID=chr7,length=159138663>
##contig=<ID=chr8,length=146364022>
##contig=<ID=chr9,length=141213431>
##contig=<ID=chr10,length=135534747>
##contig=<ID=chr11,length=135005116>
##contig=<ID=chr12,length=133851885>
##contig=<ID=chr13,length=115169878>
##contig=<ID=chr14,length=107349540>
##contig=<ID=chr15,length=102531392>
##contig=<ID=chr16,length=9054175>
##contig=<ID=chr17,length=8119216>
##contig=<ID=chr18,length=78877248>
##contig=<ID=chr19,length=59128983>
##contig=<ID=chr20,length=63025520>
##contig=<ID=chr21,length=48129695>
##contig=<ID=chr22,length=51304566>
##contig=<ID=chrX,length=155270560>
##contig=<ID=chrY,length=59373566>
##contig=<ID=chrM,length=16571>
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT c a e control
chr1 2488976 T C BAYES=-3.472;STRAND_FS=-0.906;POS_KS=-0.467;FILT_MQ_FS=0.000;FILT_BQ_FS=-0.807;MQ_FS=0.793;146,0,0;ctrlDP4=208,27,0,1;DUPLICATION=3,0,0,0 GT:AD:AF:DP 0/1:898,6:0.606:947 0/1:235,1:0.004:236
chr1 2489954 A G BAYES=-5.998;STRAND_FS=-0.436;POS_KS=-0.724;FILT_MQ_FS=0.000;FILT_BQ_FS=-0.836;MQ_FS=0.453;805,8,0;ctrlDP4=99,208,1,0;DUPLICATION=3,0,0,0 GT:AD:AF:DP 0/1:1268,8:0.006:1279 0/1:307,1:0.003:308
chr1 2492101 C G BAYES=-2.206;STRAND_FS=0.000;POS_KS=-0.020;FILT_MQ_FS=0.000;FILT_BQ_FS=0.000;MQ_FS=0.059;1,5;ctrlDP4=232,135,0,0;DUPLICATION=3,0,0,0 GT:AD:AF:DP 0/1:1319,6:0.005:1326 0/1:367,0:0.000:367
chr1 2493028 T C BAYES=-4.409;STRAND_FS=0.000;POS_KS=-0.222;FILT_MQ_FS=0.000;FILT_BQ_FS=-0.663;MQ_FS=0.925;134,7,0;ctrlDP4=305,30,1,0;DUPLICATION=6,0,0,0 GT:AD:AF:DP 0/1:1059,7:0.007:1075 0/1:335,1:0.003:339
chr1 2493950 A T BAYES=-3.155;STRAND_FS=0.000;POS_KS=-0.029;FILT_MQ_FS=0.000;FILT_BQ_FS=0.000;MQ_FS=0.00
```

Create_Control_Database.HotMutList 为解读提供的热点变异文件列表，如无更新可参考流程自带的 list 文件，文件格式如下：

```
(base) b2c_v03_pipeline@t1-jdpc-24-4[Wed Oct 14] /jdfst1/B2C_CM_P1/PipeAdmin/04_Pipeline/bin/afcol1
$ head Create_Control_Database.HotMutList
#Chr Stop Ref Call Gene_Symbol Transcript cHGVS pHGVS
chr1 115252202 G A NRAS NM_002524.4 c.438C>T p.(=)
chr1 115252202 G C NRAS NM_002524.4 c.438C>G p.(=)
chr1 115252202 G T NRAS NM_002524.4 c.438C>A p.(=)
chr1 115252203 G A NRAS NM_002524.4 c.437C>T p.A146V
chr1 115252203 G C NRAS NM_002524.4 c.437C>G p.A146G
chr1 115252203 G T NRAS NM_002524.4 c.437C>A p.A146D
chr1 115252204 C A NRAS NM_002524.4 c.436G>T p.A146S
chr1 115252204 C T NRAS NM_002524.4 c.436G>A p.A146T
chr1 115252204 C G NRAS NM_002524.4 c.436G>C p.A146P
```

-c 0.05 为指定人群频率过滤阈值，人群发生频率低于该阈值的变异则会在 Control 集合中进行剔除。目前默认人群频率阈值为 5% 以上的会在临床检测中进行过滤。

集合构建结果

Control 集合构建结束后会生成如下两个文件：

```
(base) b2c_rd3_pipeadmin@bjtj-login-24-4[Wed Oct 14] /jdfstj1/B2C_COM_P1/PipeAdmin/04.Pipeline/bioinfotoolkit/02.toolkit/04.Create_ControlSite
$ ll
total 6806
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 11393 Oct 14 09:14 Create_Control_Database.HotMutList
-rwxr-xr-x 1 b2c_rd3_pipeadmin b2c_rd3 7610 Oct 14 09:14 Create_Control_Database.pl
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 3408 Oct 14 09:14 demoInput.vcf.List
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 4292344 Oct 14 09:14 demoInput.vcf.List.DetailInfo
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 2652880 Oct 14 09:14 demoInput.vcf.List.PopCutoff_0.95.countInfo
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 102 Oct 14 09:15 demoInput.vcf.List.PopCutoff_0.05.countInfo
(base) b2c_rd3_pipeadmin@bjtj-login-24-4[Wed Oct 14] /jdfstj1/B2C_COM_P1/PipeAdmin/04.Pipeline/bioinfotoolkit/02.toolkit/04.Create_ControlSite
```

其中后缀 *DetailInfo 的文件记录了这些变异在各个样本中的变异检出频率等相关细节，用于可能需要的结果核查，数据回溯等。示例如下：

```
$ head demoInput.vcf.List.DetailInfo
chr1 201981873 A T 0 0 0.315789473684211 0.315789473684211 0 0 6 6 19
chr1 201981873 A G 0 0 0.368421052631579 0.368421052631579 0 0 7 7 19
chr1 201981873 A C 0 0 0.105263157894737 0.105263157894737 0 0 2 2 19
chr14 81610900 T C 0 0 0.263157894736842 0.263157894736842 0 0 5 5 19
chr1 39325334 G C 0 0 0.263157894736842 0.263157894736842 0 0 5 5 19
chr19 36216049 T C 0 0 0.263157894736842 0.263157894736842 0 0 5 5 19
chr1 51436071 T C 0 0 0.210526315789474 0.210526315789474 0 0 4 4 19
chr19 52693305 A C 0 0 0.105263157894737 0.105263157894737 0 0 2 2 19
chr17 37639305 G A 0 0 0.157894736842105 0.157894736842105 0 0 3 3 19
```

其中后缀 *PopCutoff_0.05.countInfo 的文件，为最终生成的 Control 集合文件，结果示例如下：

```
(base) b2c_rd3_pipeadmin@bjtj-login-24-4[Wed Oct 14] /jdfstj1/B2C_COM_P1/PipeAdmin/04.Pipeline/bioinfotoolkit/02.toolkit/04.Create_ControlSite
$ head demoInput.vcf.List.PopCutoff_0.05.countInfo
chr Site REF ALT >25%SamPopFreq <5%SamPopFreq LocalPopFreq >25%SamNum 5%_25%SamNum <5%SamNum altSampleNum totalSampleNum
chr1 201981873 A T 0 0 0.315789473684211 0.315789473684211 0 0 6 6 19
chr1 201981873 A G 0 0 0.368421052631579 0.368421052631579 0 0 7 7 19
chr1 201981873 A C 0 0 0.105263157894737 0.105263157894737 0 0 2 2 19
chr14 81610900 T C 0 0 0.263157894736842 0.263157894736842 0 0 5 5 19
chr1 39325334 G C 0 0 0.263157894736842 0.263157894736842 0 0 5 5 19
chr19 36216049 T C 0 0 0.263157894736842 0.263157894736842 0 0 5 5 19
chr1 51436071 T C 0 0 0.210526315789474 0.210526315789474 0 0 4 4 19
chr19 52693305 A C 0 0 0.105263157894737 0.105263157894737 0 0 2 2 19
chr17 37639305 G A 0 0 0.157894736842105 0.157894736842105 0 0 3 3 19
```

其中前 4 列为变异信息分别为染色体、位置、Ref、Alt；

其中 5-8 列为不同变异频率阈值下的人群频率，(>25%SamPopFreq, 表示该变异检出频率大于 25%的患者在所有患者中的占比，LocalPopFreq 为不进行频率过滤的人群占比，)

其中 9-12 列为不同变异频率阈值下的变异患者数目。

其中第 13 列为样本总数。

注：流程中的人群频率过滤仅针对 “>25%SamPopFreq” 和 “<5%SamPopFreq” 两个变异频率筛选后的人群频率进行过滤。

其中 “>25%SamPopFreq” 大于 5%，则认为该位点为人群多态性位点；

其中 “<5%SamPopFreq” 大于 5%，则认为该位点存在系统性偏好错误。

后续需要基于流程使用位置确定是否需要**对变异信息进行左对齐或 3' 对齐**。

6.3. ErrorBaseLine 集合构建 - 基于 case 样本

构建数据集合原理：

- 使用天津本地分析流程对所有临床样本进行变异检测分析。通过分析获得变异检出结果文件原始（vcf），该文件获取方式，应和临床实际分析模式保持一致，避免分析过程等其他差异导致错误集合和临床真实错误错误集合出现偏差。
- 对文件进行整理和汇总，获得组织样本人群检出频率大于 10%的所有变异位点及相关检出频率信息；
- 针对保留下来的所有变异，统计检出的突变频率，并基于所有检出频率结果计算出位点的检测阈值， $cutoff = \text{均值} + 3 * \text{标准差}$ 。

ErrorBaseLine 集合构建方法：

使用：仓库脚本

bioinfotoolkit/02.toolkit/04.Create_ErrorBaseLine/Create_Control_ErrorBaseLine.pl

整理历史临床样本，对应组织结构检测得到的 vcf 文件对应的文件路径 list 作为 ErrorBaseLine 集合构建脚本的输入文件。流程运行后会生成对应的 ErrorBaseLine.vcf，为原始的过滤集合，后续需要基于流程使用位置确定是否需要**对变异信息进行左对齐或 3' 对齐**。

```
(base) b2c_rd3_pipeadmin@jlogin-24-4[Mon Nov 15] /jdfstj1/B2C_COM_P1/PipeAdmin/04.Pipeline/bioinfotoolkit/02.toolkit/04.Create_ErrorBaseLine
$ perl Create_Control_ErrorBaseLine.pl -l demoInput.vcf.list
(base) b2c_rd3_pipeadmin@jlogin-24-4[Mon Nov 15] /jdfstj1/B2C_COM_P1/PipeAdmin/04.Pipeline/bioinfotoolkit/02.toolkit/04.Create_ErrorBaseLine
$ ll
total 4577
-rwxr-xr-x 1 b2c_rd3_pipeadmin b2c_rd3 6141 Nov 15 09:48 Create_Control_ErrorBaseLine.pl
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 2148 Nov 12 18:07 demoInput.vcf.list
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 2827441 Nov 15 09:48 demoInput.vcf.list.PopCutoff_0.1.ErrorBaseLine.Detail.tsv
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 1837701 Nov 15 09:48 demoInput.vcf.list.PopCutoff_0.1.ErrorBaseLine.vcf
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 191 Nov 15 09:35 demo.sh
-rw-r--r-- 1 b2c_rd3_pipeadmin b2c_rd3 11393 Nov 12 17:17 HotMutList
```

批注 [刘博(Bo2)]: 目前未检出的样本未纳入方差及均值的计算；是否有必要考虑

6.4. 集合使用方法

完成*PopCutoff_0.05.countInfo 文件的获取后，需要在原有的生信流程中添加 Control 集合过滤步骤，将过滤后人群频率大于 5%的变异位点进行剔除，仅保留剩余的变异位点用于后期的遗传解读。基于染色体、位置、Ref、Alt 作为变异信息的唯一识别码进行识别并过滤。

完成* ErrorBaseLine.vcf 文件的获取后，需要在原有的生信流程中添加 ErrorBaseLine 过滤步骤，将样本中检出同时频率低于 ErrorBaseLine 确定频率阈值的相关变异进行剔除，仅保

留不位于 ErrorBaseLine 集合，或检出频率大于 ErrorBaseLine 集合频率阈值的变异用于后期的遗传解读。

7 相关记录 (Related Records)

无

8 参考文献资料 (References)

无

9 附录 (Appendix)

附表 1 主要标准作业指导书一览表

附表 1 主要标准作业指导书一览表

操作	文件名称	文件编号	版本号

-----终止符-----