
吉康医学 10X 单细胞转录组常见问题及解读

SINGLE CELL SEQUENCING QUESTION AND ANSWER



目录

一、生产环节 (纯化, 建库与测序等)	3
1.fastQC 结果中 read1 的质量和 GC 含量图分布为什么出现异常? 30bp 后没有 oligdT 的峰值?	3
2.影响 cDNA 浓度的主要因素有哪些?	3
3.什么是 a template switch oligo (TSO)?	3
4.单细胞转录组 3' 试验与 5' 试验的区别有哪些?	4
5.V2 与 V3 试剂主要的区别是什么, 两种试剂结果是否可以整合分析?	4
6.上机细胞数计算方法?	5
7.QC 报告中, read2 的 GC distribution over all sequences 出现双峰的原因有哪些, 是否会影响分析结果?	5
8.分子生物学中, Mb, kb, bp 分别代表什么意思, 它们之间怎样换算?	5
9.QC 报告中, 为什么 read2 的 duplication 水平高低不定? 对分析结果有无影响?	6
二、标准分析环节	6
1.为什么样本中线粒体基因表达较高?	6
2.为什么样本中核糖体蛋白基因表达较高?	6
3.在 Cell Ranger 分析与 cloupe 中, 是如何计算“Log2 Fold Change”的?	7
4.什么是测序饱和度?	7
5.影响测序深度的因素?	8
6.10x Cell Ranger 依据 barcode 鉴定 cell 的原则?	8
7.差异基因展示中, 小提琴图怎么看?	8
8.Cell Ranger 网页版报告, 第二页差异基因的筛选规则是什么?	9
9.如何理解 Graph-based 这种聚类方法?	9

10.如何理解 K-Means 这种聚类方法? 9

11.10x 的 UMI 序列有多少种, 如何认为 UMI 序列是有效的? 9

12.Cell Ranger 的比对规则是什么? 9

13.我们交付的结果, 是否对细胞进行了过滤? 10

14.10x 单细胞转录组 cloup 软件如何使用? 10

三、售后及个性化分析环节 10

1.3'转录组与 5'转录组是否可以整合分析? 10

2.K-means 聚类, 如何选择最优 K 值? 10

3.是否需要将 UMI 转化为 TPM, RPKM 或 FPKM? 11

4.Cell Ranger 与 Seurat 软件分析的区别? 11

5.细胞类型定义的方法有哪些? 11

6.一些文章中对单细胞数据去除了细胞周期基因, 我们交付客户的结果是否进行该项内容, 是否有必要去除? 12

7.Reads Mapped Confidently to Transcriptome < 60%怎么办? 12

8.为什么相同的数据 Cell Ranger 分析后, 得到的 t-SNE 图会不一样? 13

9.Cell Ranger 可以去除多细胞捕获吗? 13

10.Reads Mapped to Genome 明显低于常规水平? 13

11.相同的物种, 但参考基因组版本不一致, 可以一起整合分析吗? 13

12.单细胞分析 Cell Ranger 可以分析多少细胞? 13

13.影响基因检出的因素有哪些? 13

14.影响捕获细胞数目的因素有哪些? 14

15.拟时分析有什么要求? 14

一、生产环节 (纯化, 建库与测序等)

1.fastQC 结果中 read1 的质量和 GC 含量图分布为什么出现异常? 30bp 后没有 oligdT 的峰值?

答: 1) 由于10X特殊的实验方式, 测序结果read1从起始1bp至16bp为barcode序列, 紧接着是10bp (或12bp) 的UMI序列, 再接着是poly(dT)VN, 由于序列的特殊结构, 导致26bp (或28bp) 之后的碱基测序质量和GC分布抖动, 这是正常现象, 而且10x Cell Ranger在分析时read1只用到前26bp (或前28bp) 的序列信息, 后面序列的信息不会对结果产生影响。

2) 理论上reads1序列30bp后有oligdT序列, FastQC的ATCG分布图呈现dT峰值; 而我们测得有些reads1不是这种情况, 30bp后dT峰值不明显, 序列没有oligdT结构, 这是由于:

a. 若30bp后测序质量很高, fastQC有dT峰值, 序列有oligdT结构; 若测序质量差, fastQC没有dT峰值, 序列没有有oligdT结构。

b. 若平台为 NovoSeq 测序, 包 lane 则 fastQC 没有 dT 峰值, 序列没有 oligdT 结构; 不包 lane 则有 dT 峰值和 oligdT 结构; 若测序平台为 X-ten, 不论是否包 lane, reads1 的 fastQC 都呈现 dT 峰值、序列有 oligdT 结构。此问题咨询了 10x 官方, 因在分析时只用 Reads1 的前 26bp, 所以这对分析是没有影响的。

2.影响 cDNA 浓度的主要因素有哪些?

答: 1) 样本的类型会影响扩增后的cDNA浓度, 不同细胞的基因表达水平不一, 另外有些细胞容易粘黏或贴管壁。

2) 细胞捕获效率, 直接影响细胞捕获效率。比如: 细胞活性与浓度、细胞状态、组织解离、FACS等。样本上芯片前的存放时间, 实验操作是否严谨, 活性浓度计数是否准确, 及PCR循环数 (增加cDNA的PCR扩增循环或样本index的循环数可能会导致文库duplicates比例升高) 等。

3.什么是 a template switch oligo (TSO)?

答: TSO (template switch oligo)是一种寡核苷酸, 在逆转录过程中, 它与逆转录酶添加的未模板化的C核苷酸杂交。TSO为全长cDNA添加了一个通用的5'序列, 用于下游cDNA扩增。

TSO 在单细胞 3'试验和单细胞 5'试验中不同。在 3'试验中, polyd(T)序列是凝胶珠 oligo 的一部分(凝胶珠 oligo 也包含 10x 条码、UMI 和部分 Illumina Read 1 序列), TSO 在 RT 引物中提供。5'试验中, RT 引物中含有 polyd(T), TSO 是凝胶珠寡聚物的一部分。

Single Cell 3' assay after reverse transcription:



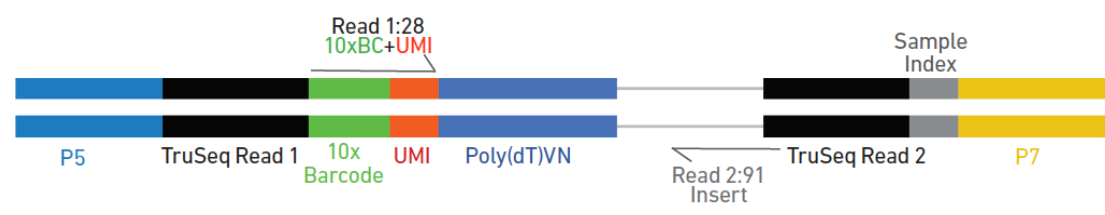
Single Cell 5' assay after reverse transcription:



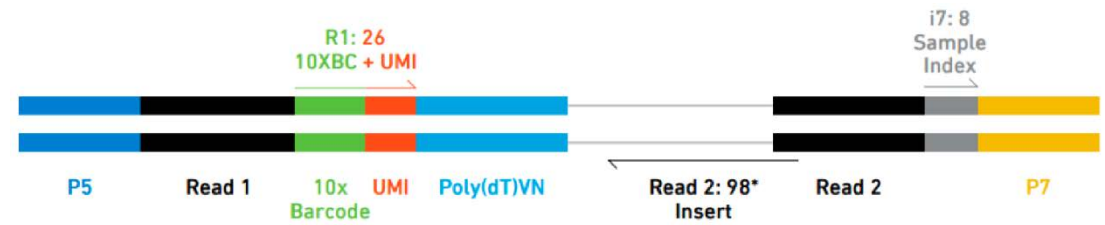
4.单细胞转录组 3'试验与 5'试验的区别有哪些？

答：方法相似，主要为 10X 文库捕获 polyadenylated 尾的转录本时不同。都是使用 ploydT 引物反转录，3'试验 polydT 引物位于 gel bead oligo（凝胶珠寡头上）而 5'试验 polydT 供给是在 RT 引物里。

Single Cell 3' v3 Gene Expression Library:



Single Cell 3' v2 Gene Expression Library:



5.V2 与 V3 试剂主要的区别是什么，两种试剂结果是否可以整合分析？

- 答：1) V2试剂的UMI序列为10bp, barcode库约75万；V3试剂的UMI序列为12bp, barode库约360万；
- 2) V3试剂基因检出有所增多，但成本更高，建库价格更高。
- 3) 两种试剂结果可以进行 Cell Ranger 的 aggr 整合（在信息人员分析整合两种试剂的数据时，需注意添加参数），也可以用 Seurat 软件整合。虽然软件矫正批次效应问题，采用同一种试剂的效果会更好。

6.上机细胞数计算方法?

答: 如10X 外出实验信息登记表中有细胞浓度 (cell/ml) : 1.06×10^6 ; cell/NF水: 15.5/18.3 水, 则上机细胞数为: $1.06 \times 10^6 \times 15.5 \div 103=16430$ 个细胞。个细胞。

NF 为 Nuclease-Free Water, 根据客户预期捕获细胞数目, 细胞浓度, 对应 Cell Suspension Volume Table, 获得需要 NF 水和细胞悬液的体积, 之后上机捕获细胞。详情参考单细胞板块外出试验操作。

7.QC 报告中, read2 的 GC distribution over all sequences 出现双峰的原因有哪些, 是否会影 响分析结果?

答: 有可能是纯化过程不够纯, 存在测序接头污染 (比较难做到特别纯); 也有可能是文库自身特点。若存在测序接头污染情况, 在Cell Ranger分析中, STAR比对时, 接头序列会被排除在外, 不会对细胞聚类分型结果有任何影响。

对上述问题, 10X 官方回复:

I don't see any major red flags in those FAST QC reports. I would go ahead with running Cell Ranger. Since Cell Ranger uses the STAR aligner underneath to align the reads many contaminants will not align, and thus be excluded from the analysis. The criteria for a read to be counted towards the expression of a gene is also quite strict. There is more on what Cell Ranger considers confidently mapped reads here:

<https://support.10xgenomics.com/single-cell-gene-expression/software/pipelines/latest/algorithms/overview>

8.分子生物学中, Mb, kb, bp 分别代表什么意思, 它们之间怎样换算?

答: 用来表示DNA长度的单位, 也就是常说的碱基数 (bp), G表示碱基数 (Gbase, Gb), 而不是计算机的储存单位 (gigabyte, GB)。

Mb=1000kb=1000000bp; G=1000M。

与测序策略有关:

PE150 (2*150), 即双端测序, 每条read长150bp。150bp X 2端X read数= 数据量。

双端测序, 一个RNA片段, 即fragment, 也叫read, 会测出来2条序列。

SE50 (1*50), 即单端测序, 每条read长50bp。50bp X 1端X read数= 数据量。测20M read, 则50bp X 1端X 20M read = 1000M = 1G。

另外, 测序下机得到的原始数据 (Sequenced Reads, 测序读段) 通常是压缩的 fastq 格式,

文件大小不是测序数据量。测序数据量需要通过 FastQC 才能得到。

9. QC 报告中，为什么 read2 的 duplication 水平高低不定？对分析结果有无影响？

答：对于普通转录组，duplication水平可能会反映数据质量，可能与PCR扩增有关。但这一指标不适用于10x单细胞转录组的数据。原因如下：

- 1) 10X单细胞转录组仅测序3'（或5'）端的150bp，不同的的转录本可能对应相同的一个基因。且当某一个基因表达量很高时，对应得UMI计数会更高，而细胞数量又很大，所以duplication水平会很高。
- 2) 10X 单细胞转录组基因的表达水平是以 UMI 进行计数，相同序列的 reads 因为 UMI 一样只会计数一次，所以不会对结果产生显著的影响。

二、标准分析环节

1.为什么样本中线粒体基因表达较高？

答：线粒体基因在多数细胞中均有表达，其表达水平与细胞类型，状态有关。

高表达水平可能原因：1) 样本质量差，较多的细胞处于凋亡或溶解状态，如果只是单个或较少cluster的细胞含有差异上调的线粒体基因和较低的总UMI count，这个cluster可能为死的或凋亡的细胞。

2) 样本的独特性，如肿瘤或肝脏组织代谢旺盛，线粒体基因表达水平较高。

2.为什么样本中核糖体蛋白基因表达较高？

答：单细胞 3'文库包含比对到核糖体蛋白转录本(Rps, Rpl)的 reads。根据细胞类型不同，一般规律为：

Cell Type	Fraction of reads mapping to ribosomal proteins
Barnyard Cells(1:1 HEK293T:3T3)	25-35%
Neuronal Cells	15-20%
PBMCs	35-40%
Isolated Pan T Cells	40-45%

For specific gene information,our PBMC and isolated Pan T Cell datasets are available on the Support website([4K PBMCs and Pan T Cells](#))

在数据分析过程中解决核糖体基因的表达问题，可以采用以下方法：

1) Exclude PCs (Principal Components) that correlate with ribosomal protein genes. This amounts to subtracting effects that are presumably technical.

2) Only include in the PC analysis genes that are "highly variable". We do this by selecting the genes with the highest dispersion across the dataset and performing PCA on those genes only. The Seurat package uses a similar approach.

相关链接:

<https://kb.10xgenomics.com/hc/en-us/articles/218169723-What-fraction-of-reads-map-to-ribosomal-proteins>

3.在 Cell Ranger 分析与 cloupe 中, 是如何计算“Log2 Fold Change”的?

答: “Log2 Fold Change”是归一化平均基因UMI计数, 每个cluster或组中, 相对于所有其他cluster或组的比值。

计算公式: 在Log2 Fold Change中引入了一个pseudocount

$\text{'log2_fold_change':np.log2}((1+\text{gene_sums_a})/(1+\text{size_factor_a}))$

$\text{np.log2}((1+\text{gene_sums_b})/(1+\text{size_factor_b}))$

相关链接:

<https://kb.10xgenomics.com/hc/en-us/articles/360007388751-How-is-Log2-Fold-Change-calculated->

4.什么是测序饱和度?

答: 这个词量化了来自已经检出的UMI的fraction of reads的比例。更具体地说, 这是可靠映射的部分, 有效的细胞barcode, 有效的UMIreads是非唯一的(匹配现有的细胞barcode, UMI, 基因组合)。反映文库的复杂程度。

计算公式: $\text{Sequencing Saturation} = 1 - (\text{n_deduped_reads} / \text{n_reads})$

n_deduped_reads = Number of unique (valid cell-barcode, valid UMI, gene) combinations among confidently mapped reads.

n_reads = Total number of confidently mapped, valid cell-barcode, valid UMI reads.

注意, 该分数的分子是 n_deduped_reads , 而不是定义中提到的非唯一读取。 n_deduped_reads 是惟一性的程度, 而不是复制/饱和的程度。因此, 我们使用 $1 - (\text{n_deduped_reads} / \text{n_reads})$ 来测量饱和度。

相关链接:

<https://kb.10xgenomics.com/hc/en-us/articles/115003646912-How-is-sequencing-saturation-calculated->

5.影响测序深度的因素?

答：测序饱和度是在实验中对文库复杂度进行测序的分数。测序饱和度的反比可以解释为一个新的read所能找到的新转录本的数量。如果测序饱和度为50%，则意味着每2个新的reads将检测到1个新的UMI count (unique transcript)。相比之下，90%的测序饱和度意味着需要10个新的read才能获得一个新的UMI计数。

测序饱和度取决于文库的复杂度和测序深度。不同的细胞类型有不同数量的RNA，因此最终文库中不同转录本的总数也不同(也称为文库复杂性)。随着测序深度的增加，可以检测到更多的基因，但根据细胞类型的不同，在不同的测序深度，这一过程会达到饱和。

测序深度也影响测序饱和度；通常，测序read越多，可以检测到的额外的独特转录本就越多。但受到库复杂性的限制。

相关链接：

<https://kb.10xgenomics.com/hc/en-us/articles/115005062366-What-is-sequencing-saturation->

6.10x Cell Ranger 依据 barcode 鉴定 cell 的原则?

答：Cell Ranger软件3.0及以上版本，依据barcode识别有效细胞cell，分为两步：

第一步，使用cutoff值，识别高RNA含量的细胞。Cell Ranger将期望捕获的细胞数量N（默认3000个细胞）中barcodes的UMI总数由高到低进行排序，取前N个UMI数值的99%分位数为最大估算UMI总数(m)，将UMI数目超过m/10的barcodes用于细胞计数。

第二步，选择一组具有低 UMI 计数的 barcode，表示“空 GEM”分区，建立背景模型。利用平滑算法，找到“空 GEM”集合中的 RNA 图谱；再将第一步中未作为细胞计数的 barcode 中 RNA 图谱与背景模型进行比较，其 RNA 图谱与背景模型存在较大差异的 barcode 用于细胞计数。

7.差异基因展示中，小提琴图怎么看?

答：小提琴图反映了该cluster的细胞中，某一个基因的表达情况，及密度分布。小提琴图的最大宽度取决于给定cluster内基因在细胞内的表达丰度，与cluster内细胞数、cluster间细胞数差异无关；基因在cluster的大多数细胞表达为0，小提琴图的最大宽度在基因丰度为0时实现，这样高表达丰度的基因在小提琴图上的宽度很小（呈一条线）。

相关链接：

<https://github.com/satijalab/seurat/issues/297>

8. Cell Ranger 网页版报告，第二页差异基因的筛选规则是什么？

答：各cluster中差异高表达的基因，UMI count > 1, log2foldchang > 0, p-value ≤ 0.1。

p-value：基于负二项检验，p值表示的是差异显著性的度量，报告中显示的p值是经过多次Benjamini-Hochberg校正的。

9. 如何理解 Graph-based 这种聚类方法？

答：是基于图的聚类算法，通过构建稀疏邻近的图，然后在图中寻找高度连接的Louvain算法。k值是细胞数量值取对数得到；在聚类的过程中，若cluster中没有差异基因，则会进行进一步的层次聚类，直到没有可以合并的cluster。可以理解为系统自动分群。

10. 如何理解 K-Means 这种聚类方法？

答：k-means的算法是将数据构建k个子集，然后计算每个类的中心点，向中心点聚集，直到聚类结果不再变化时停止，可以理解为人为规定手动进行分群。分析流程默认K-means为10，最大可以做到50。

11. 10x 的 UMI 序列有多少种，如何认为 UMI 序列是有效的？

答：10x的UMI是随机序列，V2试剂10bp，V3试剂12bp，因此存在 4^{10} 或 4^{12} 种可能性。以物种人为例，其参考基因组注释已经非常完善，总基因数不足4万，所以10x UMI的种类足以覆盖细胞中全部mRNA。认为有效UMI的规则如下：

- (1) Must not be a homopolymer, e.g. AAAAAAAAAA;
- (2) Must not contain N;
- (3) Must not contain bases with base quality < 10;
- (4) UMIs that are 1 mismatch away from a higher-count UMI are corrected to that UMI if they share a cell barcode and gene.

12. Cell Ranger 的比对规则是什么？

答：比对到基因组：利用STAR比对，比对速度快，灵敏度高，允许基因的可变剪切存在。比对完之后，利用GTF文件将reads溯源回外显子区、内含子区、基因间区。如果一条read的50%以上与外显子有交集，那么就认为它在外显子区；如果不在外显子区，与内含子有交集，那么就认为它在内含子区；与外显子、内含子都没有交集，那么就认为在基因间区。

利用MAPQ辅助判断：如果reads比对到了一个外显子区，同时也比对到了1个或多个的非外显子区，更相信它在外显子区；然后看MAPQ值，值越大越可信，如果MAPQ的值为255的

话, 那么就可以非常确定它比对到了外显子区。

比对到转录组: 如果上面得到的外显子区域reads同时比对上有注释转录本上的外显子, 并且在同一条链上, 那么认为这个reads也比对到了转录组; 如果只比对到单个基因的注释信息, 那么认为它是特异比对到转录组的(uniquely /confidently mapped), 这样的reads才会拿来作接下来的UMI计数。

13.我们交付的结果, 是否对细胞进行了过滤?

答: 过滤了, 人和小鼠模式生物, 针对以下内容进行了过滤:

- 1) 血红蛋白基因占比情况;
- 2) 线粒体基因占比情况;
- 3) barcode中过高的UMI与基因数。

对于其他物种, 若客户提供血红蛋白基因或线粒体基因列表, 也可以进行过滤。

14.10x 单细胞转录组 cloup 软件如何使用?

答: 单细胞研究板块提供cloup使用说明文档pdf, 10X官网也有教学视频, 相关链接:
<https://support.10xgenomics.com/single-cell-gene-expression/software/visualization/latest/what-is-loupe-cell-browser>

三、售后及个性化分析环节

1.3'转录组与 5'转录组是否可以整合分析?

答: 理论上是可以的, 但需要注意aggr不能矫正不同试剂带来的影响。对单细胞RNA-seq数据中的系统效应或批量效应进行归一化和校正是目前比较活跃的研究领域, 目前10x 还没有一个具体的建议。从文献中可以看出, R中有许多包, 比如Seurat、scran和scone, 试图解决这些问题。

相关链接:

<https://kb.10xgenomics.com/hc/en-us/articles/115003145272-Can-I-combine-gene-expression-data-from-3-and-5-assay-chemistries->

2.K-means 聚类, 如何选择最优 K 值?

答: 推荐使用手肘法(SSE)与轮廓系数法(Silhouette Coefficient)结合, 判断最优K值。

手肘法的核心指标是SSE(sum of the squared errors, 误差平方和), 随着聚类数的增大, 样本

划分会更加精细，每个簇的聚合程度会逐渐提高，因此误差平方和SSE会逐渐变小。但当K小于真实聚类数时，由于K的增大会大幅增加每个簇的聚合程度，所以SSE的下降幅度会很大；当到达真实聚类数时，再增加所得到的聚合程度回报会迅速变小，所以SSE的下降幅度会骤减，随着值的继续增大而趋于平缓，也就是说SSE和的关系图会形成一个手肘的形状，而肘部对应的值就是数据的真实聚类数。轮廓系数法的核心指标是轮廓系数(Silhouette Coefficient)，该方法结合了聚类的凝聚度(Cohesion)和分离度(Separation)，用于评估聚类的效果。该值处于-1~1之间，值越大，表示聚类效果越好。

轮廓系数法确定出的最优值不一定是最优的，需要根据SSE 进行辅助选取，以上内容仅供参考，请结合实际情况，合理选择最优K值。

相关链接：

https://blog.csdn.net/qq_15738501/article/details/79036255

3.是否需要将 UMI 转化为 TPM, RPKM 或 FPKM?

答：不需要。10x 也不建议如此。Cell Ranger分析使用唯一标识UMI用于基因表达水平的定量，后续分析（包括Cell Ranger及其他分析软件）也均基于UMI。

此外，10X单细胞与传统转录组测序不同，在传统的RNA-seq数据中，完整的转录本被片段化，然后cDNA合成、末端修复和adapter连接等。在这个流程中，从长文本中抽取片的概率要高于从短文本中抽取片的概率。因此，根据reads长度(例如TPM、RPKM、FPKM)对读取计数进行规范化是有意义的。而在10x单细胞3'或5'的实验中，这种基因长度偏差并不存在。所以，不建议根据基因长度对UMI计数进行标准化。

相关链接：

<https://kb.10xgenomics.com/hc/en-us/articles/115003684783-How-to-calculate-TPM-RPKM-or-FPKM-instead-of-counts->

4.Cell Ranger 与 Seurat 软件分析的区别?

答：二者侧重不同，Cell Ranger侧重于拆库定量，可以得到cell-gene矩阵，初步的细胞聚类分型结果及cloud可视化软件。Seurat分析基于Cell Ranger结果进行，在参数调整，差异基因展示，细胞周期分析等方面更具优势，可进行数据挖掘。但Seurat结果的可视化操作不佳，需要导入cloud中，而cloud只可通过Cell Ranger获得，所以二者相辅相成，缺一不可。

5.细胞类型定义的方法有哪些?

答：1) 传统方法：根据文献积累，细胞类型相关的数据库：BD Rhapsody

(<http://genomiccytometry.com/-/human/>) ; CellMarker (<http://biocc.hrbmu.edu.cn/CellMarker/>) 等, 根据marker基因表达情况, 结合聚类热图, t-SNE图或小提琴图等综合判断。此方法需要客户主导, 信息辅助, 共同完成, marker基因的选取较为关键, 但该方法适用于大部分物种。

2) 一些单细胞分析软件, 如: Seurat (相关链接:

https://satijalab.org/seurat/v3.0/pancreas_integration_label_transfer.html)和monocle(相关链接: <https://cole-trapnell-lab.github.io/garnett/docs/-2-classifying-your-cells>) , 可通过预测, 辅助进行细胞类型定义。

3) 细胞定义软件SingleR, 是2019年发表在Nature Immunology杂志上面的细胞类型鉴定的工具, 可读取10x Cell Ranger分析结果, 也可与Seurat工具无缝对接。该软件基于超几何分布进行假设检验, 以判断每个cell或cluster的最可能的细胞类型。该软件只适用于人和小鼠。单细胞研究板块有该分析流程, 可接售后。

4) 其他在线分析工具等, 可参考相关链接: <https://www.jianshu.com/p/98956bce75d4>

6.一些文章中对单细胞数据去除了细胞周期基因, 我们交付客户的结果是否进行该项内容, 是否有必要去除?

答: 细胞周期阶段的异质性, 特别是有丝分裂细胞在S期和G2/M期之间的过渡, 可以驱动大量的转录组可变, 可能会掩盖某些生物信号。为了减弱这种可能存在的批次效应, 一些学者会去除该变异源。当然, 这需要客户根据自己的研究方向与内容, 判断是否需要。

我们提供的结果, 默认不去除细胞周期因素影响, 因为这不适用于全部客户数据。若客户有需求, 可以分析。

相关链接:

https://satijalab.org/seurat/v3.1/cell_cycle_vignette.html

7.Reads Mapped Confidently to Transcriptome < 60%怎么办?

答: 该指标一般适用于人或小鼠等模式生物, 其他物种该指标不一定适用。

对于一些特殊组织, 也不一定适用, 如, 一般造血干细胞、细胞核与胚胎组织等, 未剪切的RNA含量较高, 因此内含子比对会偏高, 而外显子比对偏低。建议用前体RNA做参考基因组, 再进行比对, 可能会提升基因数。

前体RNA参考基因组为将GTF文件中第3列替换成外显子后再次生成的参考基因组。

相关链接:

[https://support.10xgenomics.com/single-cell-gene-expression/software/pipelines/latest/advanced/references - header](https://support.10xgenomics.com/single-cell-gene-expression/software/pipelines/latest/advanced/references-header)

8.为什么相同的数据 Cell Ranger 分析后, 得到的 t-SNE 图会不一样?

答: t-SNE图是展示对细胞聚类一种展示形式, 与PCA不同, t-SNE的非线性降维是不确定的, 但这不影响数据分析结果。10X官方回复: “In contrast to, e.g., PCA, t-SNE has a non-convex objective function. The objective function is minimized using a gradient descent optimization that is initiated randomly. As a result, it is possible that different runs give you different solutions.”

Thus subtle differences like you observed between runs are not a concern. For more information please see: <https://lvdmaaten.github.io/tsne/>

9.Cell Ranger 可以去除多细胞捕获吗?

答: 不可以。Cell Ranger分析不能识别多细胞捕获, 也不能去除多细胞。但可以通过细胞的 UMIs或检测到的基因数目明显偏离异常, 来进行过滤多细胞捕获的可能。

相关链接:

<https://kb.10xgenomics.com/hc/en-us/articles/360005165411-Are-there-methods-for-identifying-multiplets->

10.Reads Mapped to Genome 明显低于常规水平?

答: 除样本中有不同物种细胞混合外, 一般原因为物种与参考基因组不匹配。

11.相同的物种, 但参考基因组版本不一致, 可以一起整合分析吗?

答: 不能。要整合分析的数据, 必须有相同的参考基因组。

12.单细胞分析 Cell Ranger 可以分析多少细胞?

答: 理论上Cell Ranger可以分析500~10000个细胞, 细胞数太少会导致分析结果不准确。

相关链接:

<https://kb.10xgenomics.com/hc/en-us/articles/115001800523-What-is-the-minimum-number-of-cells-that-can-be-profiled->

13.影响基因检出的因素有哪些?

答: 1) 主要与样本类型及状态有关, 一般外周血多为1000~3000, 干细胞和大脑可多达5000+, 不同的癌种, 基因检出数也可存在数倍的差异;

2) 随着测序饱和度的增加, 在一定程度上会增加基因检出数, 一般测序饱和度达到80%即可, 但并不是硬性指标;

3) 参考基因组的注释是否完善会影响基因检出，若有特殊物种，则需要提供较为完善的参考基因组。

4) 对于一些特殊组织，如，一般造血干细胞、细胞核与胚胎组织等，未剪切的RNA含量较高，因此内含子比对会偏高，而外显子比对偏低。建议用前体RNA做参考基因组，再进行比对，可能会提升基因数。

14.影响捕获细胞数目的因素有哪些？

答：主要影响细胞捕获数目多少的是样本处理成单细胞悬液的过程——细胞浓度与活性。样本物种类型及组织特点多样，如骨组织，灌洗液，成熟心肌组织或植物原生质体等，需要较为特殊的处理方式，以保障细胞活性；外周血单核细胞分离时需要注意避免其他细胞的污染等。

单细胞研究板块致力于样本制备，目前已成功测试过脾脏、肾脏、乳腺癌原位灶及转移灶、人-小鼠移植瘤、小鼠结肠癌原位灶与结直肠等。

15.拟时分析有什么要求？

答：拟时分析默认使用Monocle2软件（之后可能会随着软件更新，更换版本），单样本与多样本均可进行拟时分析，需要客户提供在哪个/些样本上，哪种“分组”（可以是分群方式，样本，时间段，组间等）的基础上进行拟时分析。拟时分析针对于不同细胞类型的分化轨迹或者是发育方向，如果一个样本可以涵盖所关注的所有细胞类群是可以实现的，但反之如果该样本不能涵盖足够的样本类型，建议取多样本直至包含关注的细胞亚群，结果会更准确。