

# Stereo-seq 分析流程软件包 使用手册



# 目录

## 第一章 SAW 流程简介

1.1. SAW 流程概述	1
1.2. SAW 流程及相关软件版本	1
1.3. SAW 获取方法（镜像安装）	1
1.4. SAW 运行流程系统要求	2
1.5. SAW 运行方法	2
1.6. SAW GitHub	3
1.7. SAW 测试数据	3

## 第二章 SAW 工具和命令参数说明

2.1. splitMask（Q4 测序步骤）	5
2.2. CIDCount	6
2.3. mapping	7
2.4. merge	10
2.5. count	11
2.6. register	12
2.7. imageTools	14
2.8. tissueCut	15
2.9. spatialCluster	18
2.10. cellCut	19
2.11. cellCorrect	20
2.12. cellCluster	21
2.13. saturation	22
2.14. report	

## 第三章 其他应用工具

3.1. cellCut 其他应用	26
3.2. checkGTF 应用	27
3.3. imageTools 其他应用	27
3.4. manualRegister 应用	29
3.5. lasso 应用	30

## 第四章 SAW 常见 Q&A

- 4.1. Q: 基因组注释文件 GTF/GFF 有什么格式要求? 33
- 4.2. Q: 读取注释文件时会忽略对应基因的情况有哪些? 34
- 4.3. Q: 构建 reference 时报错 "Fatal INPUT FILE error, no valid exon lines in the GTF file" 如何处理? 35
- 4.4. Q: 为什么大部分的注释文件中的基因都未被注释上? 35
- 4.5. Q: 是否有工具可以辅助检查注释文件规范? 35
- 4.6. Q: SAW 中对测序数据做了哪些过滤? 35
- 4.7. Q: 基因表达可视化结果异常, 没有组织轮廓怎么办? 35
- 4.8. Q: 如何解析 GEF 格式文件? 36
- 4.9. Q: 测序在进行 mapping 比对时, FASTQ 太多如何简易处理? 37
- 4.10 Q: Valid CID 比例异常应该怎么处理? 37

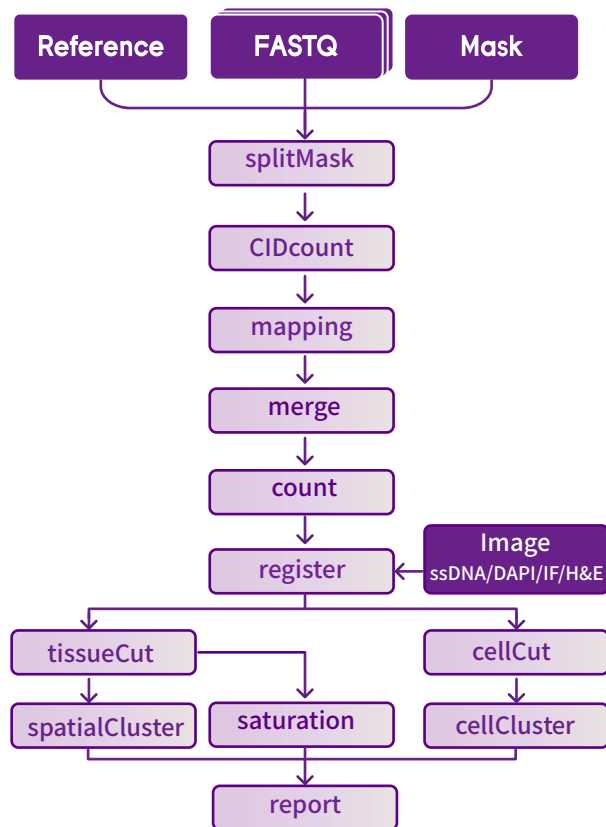
**联系我们** 38

# 第一章

## SAW 流程简介

## 1.1. SAW 流程概述

- Stereo-seq分析流程软件包 (Stereo-seq Analysis Workflow, SAW) 整合了多个Stereo-seq空间组基因表达分析工具, 这些工具可用于还原以及可视化测序数据在芯片上的空间表达信息。
- 原始测序数据经SAW分析后, 得到可以用于下游分析的空间表达矩阵。为满足更加便捷的数据分析, SAW通常具有13步必要流程, 以及其他辅助分析工具。



## 1.2. SAW 流程版本

当前SAW更新版本v7.0

- **ImageStudio**: ImageStudio是集成了拼接、校准、组织分割、标注、细胞分割等功能的图像处理软件。每个模块对应输出结果均可接入SAW进行进一步分析。SAW v7.0 版本建议ImageStudio>=v3.0。
- **StereoMap**: StereoMap是一个支持Stereo-seq分析数据高清交互式可视化的桌面端软件。SAW流程中输出的GEF矩阵、图像RPI和IPR数据、聚类结果等均可在StereoMap中展示。SAW v7.0 版本建议StereoMap版本>=v3.0。

## 1.3. SAW 镜像安装

SAW docker 镜像软件包已包含运行软件必须的依赖包, 用户可通过 Docker Hub 下载 SAW docker 镜像, 以帮助用户在本地快速搭建 SAW 镜像, 进行离线分析 ([下载最新版说明匹配相关版本](#))。

Docker Hub 链接: <https://hub.docker.com/r/stomics/saw/tags>

目前 SAW 的安装和运行支持 Singularity 和 Docker 两种方式 (如无 Singularity 软件, 请自行安装或参考 [1.6 SAW GitHub](#)),

安装的命令如下:

```

1 $ singularity build SAW_v7.0.sif docker://stomics/saw:07.0.0 ## 方法 1
2 $ singularity build --sandbox SAW_v7.0/ docker://stomics/saw:07.0.0 ## 方法 2
3 $ docker pull stomics/saw:07.0.0 ## 方法 3 利用 docker 拉取镜像
4 $ docker run -d -v /path/to/data:stomics/saw:07.0.0 /bin/sh -c "<shell-command>"
   ## 方法 4 交互式运行镜像
  
```

## 1.4. SAW 运行流程系统要求

运行SAW的Linux系统需满足的最低要求包括：

- **Cpu至少需8核(建议24核)；**
- **运行内存至少128GB(建议256G)；**
- **存储空间至少1TB；**
- **64位 CentOS/RedHat 7.8 或Ubuntu 20.04。**

运行SAW需提前安装下列软件其中一个软件：

- **docker: version >=20.10.8;**
- **singularity: 版本号>= 3.8(推荐使用)。**

## 1.5. SAW 运行方法

运行方法提供三种 (以 singularity 为例)：(红色字体均为用户需输入的实际路径)

⊙ **注意！在运行 singularity 镜像时，所有涉及目录均需提前挂载。如：输入文件目录 `/path/to/data`，参考基因组目录 `/path/to/genomeDir`，输出结果目录 `/path/to/output`。**

```
1 $ export SINGULARITY_BIND="/path/to/data,/path/to/genomeDir,/path/to/output"
```

1. 启动容器并在容器中执行命令。

示例：

```
1 $ /path/to/SAW_v7.0.sif <application> ## 选择 1
1 $ singularity exec /path/to/SAW_v7.0.sif <application> ## 选择 2
```

2. 启动容器并打开交互式终端，在容器中运行 bash 命令。需执行 exit 退出环境。

示例：

```
1 /path/to/SAW_v7.0.sif /bin/bash
2 Singularity> <shell-command>
3 Singularity> exit
```

3. 通过参数“-B 外部路径:容器内路径”挂载用户自定义目录至容器中，并在容器中执行命令。

示例：

```
1 singularity shell -B /path/to/directory/on/the/host-machine:/path/to/directory/
  mounted/in/the/container /path/to/SAW_v7.0.sif
2 Singularity> <shell-command>
3 Singularity> exit
```

## 1.6. SAW GitHub

SAW GitHub: <https://github.com/STOmics/SAW>

用户可进入GitHub页面查看singularity的安装、参考基因组索引构建、以及SAW shell脚本的详细说明。

🚨 **注意！请在运行 SAW 分析前构建索引。**

File Name	Commit Message	Time Ago
Documents	Update SAW V7.0.0	2 days ago
Scripts	Update SAW V7.0.0	2 days ago
Test_Data	Update SAW V7.0.0	2 days ago
LICENSE	Create LICENSE	4 months ago
README.md	Update SAW 7.0.0	2 days ago
SAW_v7.0_workflow.png	Update SAW V7.0.0	2 days ago

**SAW: Stereo-seq Analysis Workflow**

Workflow for analyzing Stereo-seq transcriptomic data. Stereo-seq Analysis Workflow (SAW) software suite is a set of pipelines bundled to map sequenced reads to their spatial location on the tissue section, quantify the corresponding gene expression levels and visually present spatial gene expression distribution.

DockerHub Link: <https://hub.docker.com/r/stomics/saw/tags>

## 1.7. SAW 测试数据

测试数据可从SAW GitHub页面获取。测试使用参考基因组版本为:

- genome-build: GRCm38.p6
- genome-version: GRCm38
- genome-date: 2012-01
- genome-build-accession: NCBI:GCA\_000001635.8

## 第二章

# SAW 工具和命令行参数说明



## 2.1. splitMask (Q4 FASTQ 测序步骤)

**splitMask**：是通过 Q4 FASTQ CID 来分割 Stereo-seq 芯片 mask 文件的工具。当从序列中输出 FASTQs 文件时，其分割计数和 FASTQs 文件的分割编号直接相关。

### 2.1.1. 输入文件

- Stereo-seq 芯片 mask 文件 (.h5)

☉ SAW 流程支持来自 Stereo-seq 的 Q40 FASTQ 和 Q4 FASTQ 作为测序数据输入。它们之间的区别主要体现在文件大小以及数据和信息的写入方式上。Q40 FASTQ，也叫 PE FASTQ，有一对 read 文件，read1 和 read2。Q4 FASTQ，也称为 SE FASTQ，是一种只有一个文件的输出格式，但可以节省文件存储空间。更多细节请参阅 [Stereo-seq 文件格式手册](#)。

- 示例 Q4 FASTQ 文件目录名称如下：

Q4 FASTQ name

/path/to/data/E100026571\_L01/barcode\_2/E100026571\_L01\_2\_16.fq.gz

lane
barcode
lane
barcode
split index

Get the list input conveniently like:

```

...
$ sh manage.sh <FASTQDirList> <outdir> <SN>
FASTQDirList = <Q4 FASTQ directory list which record directory paths line byline>
outdir = <directory to output FQ.list files>
SN = <Stereo-seq chip serial number>
splitCnt = 16 or 64

```

[https://github.com/STOmics/SAW/blob/main/Scripts/Get\\_FASTQ\\_list.sh](https://github.com/STOmics/SAW/blob/main/Scripts/Get_FASTQ_list.sh)

An example to display FQ.list file :

```

...
$ cat SN_Q4_fastq_16.list ## a FASTQ list file gathers all the FASTQs whose index
is 16
/path/to/data/lane_1_16.fq.gz
/path/to/data/lane_2_16.fq.gz

```

### 2.1.2. 命令示例及参数说明

```

...
1 $ mkdir /path/to/output/00.splitmask
2 $ singularity exec SAW_v7.0.sif splitMask \
3   /path/to/data/{SN}.barcodeToPos.h5 \ # 必需参数: Stereo-seq 芯片 mask 文件 (.h5)
4   /path/to/output/00.splitmask \ # 必需参数: 输出目录
5   8 \ # 必需参数: 运行线程数
6   16 \ # 可选参数: 分割份数
7   2_25 # 可选参数

```

☉ 注意！请把 {SN} 编号替换成用户自己的 Stereo-Seq 芯片 T 序列编号（例：SS200000135TL\_D1 或 A00135D1）

### 2.1.3. 运行资源

- 内存占用：3G（1cm\*1cm 尺寸芯片数据分析参考值）
- 运行时间：5 min（1cm\*1cm 尺寸芯片数据分析参考值）

### 2.1.4. 输出文件

```

1 /path/to/output/00.splitmask
2 |-- *.SN.barcodeToPos.bin(* 代表分割编号) (## 通过 CID 分割 Stereo-seq 芯片 mask 文件)

```

## 2.2. CIDCount

**CIDCount:** 是用于计算 Stereo-seq 芯片 mask 文件中 CID 数量，以及粗略估算 mapping 过程所需内存的小程序。

### 2.2.1. 输入文件

- Stereo-seq 芯片 mask 文件 (.h5)

### 2.2.2. 命令示例及参数说明


以 Q40 FASTQ 输入为例，运行 CIDCount 所需文件包括：

```

1 $ singularity exec SAW_v7.0.sif CIDCount \
2   -i /path/to/data/{SN}.barcodeToPos.h5 \ ## 必需参数: Stereo-seq 芯片 mask 文件 (.h5)
3   -s {speciesName} \ ## 可选参数: 物种名
4   -g {genomeSize} ## 必需参数: 基因组大小

```

以 Q4 FASTQ 输入为例，需要先运行 splitMask，因为每一份拆分后的 mask 文件结果类似，所以用户只需要运行一次 CIDCount 即可。针对大于 1 cm \* 1 cm 的芯片，我们推荐用户对每一份拆分 mask 都运行 CIDCount。

 **大芯片运行建议：** 对大芯片来说，每一份拆分 mask 文件包含对 CID 数波动可能会较大，所以计算出的预估内存可能会有较明显差异，因此在这种场景下，我们会推荐用户对每一份拆分 mask 文件都做一下 CIDCount 的运算。造成这种波动的原因可能是因为芯片上的 dnb/CID 分布不均。

```

1 $ singularity exec SAW_v7.0.sif CIDCount \
2   -i /path/to/output/00.splitmask/{index}.{SN}.barcodeToPos.bin \ ## 必需参数:
   Stereo-seq 芯片 mask 文件 (.h5)
3   -s {speciesName} \ ## 必需参数: 物种名
4   -g {genomeSize} ## 必需参数: 基因组大小

```

 **请注意！** {index} 需要替换成拆分芯片文件真实的 index 拆分编号。

### 2.2.3. 运行资源

- 内存占用：~0.7G (1cm\*1cm 尺寸芯片数据分析参考值)
- 运行时间：30s (1cm\*1cm 尺寸芯片数据分析参考值)

## 2.2.4. 输出文件

```
1 singularity exec SAW_v7.0.sif CIDCount -i SN.barcodeToPos.h5 -s {speciesName} -g 3
2 645784920 ##CID 数值
3 62 ## 预估比对内存大小
```

## 2.3. mapping

**mapping:** 每一条 Stereo-seq 序列包含一段 CID 序列，该序列是用于将测序 reads 比对在组织切片上的原始位置的关键信息。Stereo-seq 原始测序数据通过 SAW 软件中的 mapping 工具，将存储在 FASTQ 文件中的原始测序 reads 的 CID 与 Stereo-seq 芯片 mask 文件记录的 CID 坐标键值对进行匹配（允许一位容错）。根据 mask 文件的记录，为 CID 能够匹配的 reads 添加坐标信息。经过 CID 匹配的 reads 为具有有效 CID 的 mRNA reads（Valid CID Reads）。Valid CID Reads 经过滤后，得到 Clean Reads。mapping 工具将 Clean Reads 与相应物种的参考基因组进行比对，并输出排序后的 BAM 格式比对文件和统计报告。

### 2.3.1. 输入文件

- Stereo-seq 原始测序数据 FASTQ 文件（.fq.gz）
- Stereo-seq 芯片 mask 文件（.h5）
- 参考基因组索引文件
- bcPara 文件（.bcPara）

☹ **注意！** {SN} 指代 Stereo-seq 芯片的 SN 编号 (例: SS200000135TL\_D1); {lane} 指代测序 FASTQ lane 编号 (如, E100026571\_L01)（下同）。

两个场景：单对和多对 Q40 FASTQ 的 {lane}.bcPara 文件参数配置如下：

```

1 $ mkdir /path/to/output/01.mapping ##option1(单对 Q40 FASTQ 新建目录)
2 $ mkdir /path/to/multi_lane_output/01.mapping ## option2(多对 Q40 FASTQ 新建目录)
3 $ vim /path/to/output/01.mapping/{lane}.bcPara ##option1 (单对 Q40 FASTQ)创建参数配置
  文件
4 $ vim /path/to/multi_lane_output/01.mapping/{lane}.bcPara ##option2 (多对 Q40
  FASTQ)创建参数配置文件
5 in=/path/to/data/{SN}.barcodeToPos.h5 ## 必需参数:Stereo-seq 芯片 mask 文件(.h5)
6 in1=/path/to/data/{lane}_read_1.fq.gz ## 必需参数:Stereo-seq 原始测序数据 FASTQ read1
  文件(.fq.gz)
7 in2=/path/to/data/{lane}_read_2.fq.gz ## 可选参数:Stereo-seq 原始测序数据 FASTQ read2
  文件(.fq.gz)
8 barcodeReadsCount=/path/to/output/01.mapping/{lane}.barcodeReadsCount.txt ## op-
  tion1 (单对 Q40 FASTQ)必需参数:CID mapping 统计输出
9 barcodeReadsCount=/path/to/multi_lane_output/01.mapping/{lane}.barcodeReadsCount.
  txt ## option2(多对 Q40 FASTQ)必需参数:CID mapping 统计输出
10 barcodeStart=0 ## 必需参数:CID 起始位点
11 barcodeLen=25 ## 必需参数:CID 长度(Q40 为 25)
12 umiStart=25 ## 必需参数:MID 起始位点
13 umiLen=10 ## 必需参数:MID 长度
14 mismatch=1 ## 必需参数:mapping 的最大容错碱基数
15 bcNum=645784920 ## 必需参数:CID 数量,CIDCount 的输出文件第一行
16 polyAnum=15 ## 可选参数:polyA 长度
17 mismatchInPolyA=2 ## 可选参数:找寻 polyA 时最大容错碱基数
18 rRNAREmove ## 可选参数:是否过滤 rRNA

```

两个场景：1. 含有一个 barcode 且写出 FASTQ 时未做拆分的 Q4 FASTQ；2. 含有多个 barcode 且写出 FASTQ 时做了拆分的 Q4 FASTQ 的 {index}.bcPara 文件参数配置如下：

- ⊙ 注意：**{index}** 需要被替换成拆分 mask 文件的 index 拆分编号，即与 Q4 FASTQ 文件对应的 index 拆分编号。**{idx}** 需要换成输入 Q4 FASTQ 文件的 {index} 编号。若一个文库同一个测序泳道包含两个 barcode，则第一组 barcode 的 FASTQ 的 {index} 和 {idx} 为 01-16, 01-16，第二组 barcode 的 FASTQ 的 {index} 和 {idx} 为 01-16, 17-32。

```

1 $ mkdir /path/to/output/01.mapping ## 新建目录
2 $ vim /path/to/output/01.mapping/{index}.bcPara ##option1(只含有一个 barcode 且没有拆
  分 barcode 的 Q4 FASTQ) 创建参数配置文件
3 $ vim /path/to/output/01.mapping/{idx}.bcPara ##option2 (含有多个 barcode 的 Q4
  FASTQ) 创建参数配置文件
4 in=/path/to/output/00.splitmask/{index}.{SN}.barcodeToPos.bin ## 必需参数:Ste-
  reo-seq 芯片分割后的 mask 文件(.bin)
5 in1=/path/to/data/{lane}_{index}.fq.gz ## 必需参数:Stereo-seq 原始测序数据 FASTQ 文件
  (.fq.gz)(FASTQ files for different barcode usually stores in different directory,
  so /path/to/data/ might change to /path/to/data/{barcode_n})

```

```

6 barcodeReadsCount=/path/to/ouptut/01.mapping/{index}.barcodeReadsCount.txt ##op-
  tion1(只含有一个 barcode 且没有拆分 barcode 的 Q4 FASTQ)必需参数:CID mapping 统计输出
7 barcodeReadsCount=/path/to/ouptut/01.mapping/{idx}.barcodeReadsCount.txt ##op-
  tion2 (含有多个 barcode 的 Q4 FASTQ)必需参数:CID mapping 统计输出
8 barcodeStart=0 ## 必需参数:CID 起始位点
9 barcodeLen=24 ## 必需参数:CID 长度(Q4 为 24)
10 umiStart=25 ## 必需参数:MID 起始位点
11 umiLen=10 ## 必需参数:MID 长度
12 mismatch=1 ## 必需参数:mapping 的最大容错碱基数
13 bcNum=38284877 ## 必需参数:CID 数量
14 polyAnum=15 ## 可选参数:polyA 长度
15 mismatchInPolyA=2 ## 可选参数:找寻 polyA 时最大容错碱基数
16 rRNAremove ## 可选参数:是否过滤 rRNA

```

- ⊙ 若考虑 rRNA 对实验的影响,在后续步骤中对其统计和过滤,需在参考基因组中添加 rRNA 的相关信息。操作主要包含两步:1) 将 rRNA 信息添加到 FASTA 文件中,并且在染色体序号列上增加 '\_rRNA' 的后缀标识,例如 '1\_rRNA',与原始参考基因组信息区分开来;2) 根据修改后的参考基因组文件构建索引。

请注意重新构建参考基因组索引,使用与之前相同的命令输入。从 SAW v6.1 开始,SAW 对参考基因组索引和比对算法进行了重构,以获得更高效的计算性能和减少时间成本。

### 2.3.2. 命令示例及参数说明

Q40 FASTQ 及 Q4 FASTQ 运行 mapping 入参示例及说明:

```

1 $ singularity exec SAW_v7.0.sif mapping \
2   --outSAMattributes spatial \ ## 默认参数:转换成空间 bam 文件格式
3   --outSAMtype BAM SortedByCoordinate \ ##STAR 参数:将 bam 按照坐标进行排序
4   --genomeDir /path/to/genomeDir \ ##STAR 参数:基因组索引路径
5   --runThreadN 8 \ ##STAR 参数:运行线程数
6   --outFileNamePrefix /path/to/output/01.mapping/{lane}. \ ##STAR 参数:输出文件前缀
  (Q4 FASTQ 将 {lane} 替换成 {index})
7   --sysShell /bin/bash \ ## 可选参数:shell 文件目录
8   --stParaFile /path/to/output/01.mapping/{lane}.bcPara \ ## 必需参数:CID 比对选项参
  数文件(Q4 FASTQ 将 {lane} 替换成 {index})
9   --readNameSeparator \" \" \ ##STAR 参数:分隔符
10  --limitBAMsortRAM 38582880124 \ ##STAR 参数:排序 bam 文件的最大运行内存
11  --limitOutSJcollapsed 10000000 \ ##STAR 参数:最大崩溃连接数
12  --limitIObufferSize=280000000 \ ##STAR 参数:每个线程最大可获得的 buffer 大小
13  --outBAMsortingBinsN 50 \ ##STAR 参数:坐标排序的基因组 bins 的数量
14  --outSAMmultNmax -1 \ ##STAR 参数:一条 read 写入 bam 文件中的最大比对数,默认全部输出
15 > /path/to/output/01.mapping/{index}.run.log ## {index} or {idx} depending on the
  scenario

```

### 2.3.3. 运行资源

- 内存占用: ~67G (1G reads 数据分析参考值)
- 运行时间: ~4h (1G reads 数据分析参考值)

### 2.3.4. 输出文件

场景 1: 单对 Q40 FASTQ 文件输出目录如下:

```

1  /path/to/output/01.mapping/
2  |— lane.Aligned.sortedByCoord.out.bam ## 二进制比对 / 映射文件,用于存储序列比对信息
3  |— lane.CIDMap.stat ##CID 比对的统计报告
4  |— lane.barcodeReadsCount.txt ## 比对上 CID 的 reads 列表文件
5  |— lane.bcPara ## 定义 CID 比对选项的参数文件
6  |— lane.Log.final.out ## mapping 完成后汇总比对统计信息 (STAR 输出)
7  |— lane.Log.out ##STAR mapping 中的主日志文件 (STAR 输出)
8  |— lane.Log.progress.out ## 报告作业过程统计数据 (STAR 输出)
9  |— lane.run.log ## 比对模块的日志文件
10 |— lane.SJ.out.tab ##mapping 过程中拼接接头的检测 (STAR 输出)

```

目 场景 2: 多对 Q40 FASTQ 文件输出目录: 将上述目录中的 lane 替换成 lane\* 文件名展示每一对的比对结果。

场景 3: 有一个 barcode 或者没有拆分 barcode 的 Q4 FASTQ 文件输出目录: 将上述目录中的 lane 替换成 {index}\* 文件名。

场景 4: 有多个 barcode 的 Q4 FASTQ 文件输出目录: 将上述目录中 lane 替换成 {idx}\* 文件名。

## 2.4. merge

**merge:** 用于合并多个 mapping 分析结果。

### 2.4.1. 输入文件

- mapping 输出映射后的 CID 列表文件 (.txt)

### 2.4.2. 命令示例及参数说明

```

1  singularity exec SAW_v7.0.sif merge \
2  /path/to/data/{SN}.barcodeToPos.h5 \ 必需参数:Stereo-seq 芯片 mask 文件 (.h5)
3  /path/to/multi_lane_output/01.mapping/{lane*}.barcodeReadsCount.txt \ 必需参数:
mapping 步骤产生的 CID 列表文件
4  /path/to/multi_lane_output/02.merge/{SN}.barcodeReadsCount.txt 必需参数:merge
CID 列表文件后产生的文件

```

### 2.4.3. 运行资源

- 内存占用: ~5G (1G reads 2 lanes 数据分析参考值)
- 运行时间: 5 min (1G reads 2 lanes 数据分析参考值)

## 2.4.4. 输出文件

```

1 /path/to/multi_lane_output/02.merge
2 └─ {SN}.merge.barcodeReadsCount.txt ## 合并比对上 CID 的 reads 列表文件

```

## 2.5. count

**count:** 一个高效、多用途的序列注释模块，它标记每个 reads 的基因组特征，并输出整体的注释统计信息。通过量化被注释的基因测序 reads，综合考虑 CID、基因和 MID 后进行去重，最后计算每个位置每个基因的表达量水平，生成基因表达矩阵。SAWcount 通常对 mapping 比对结果中的 Uniquely Mapped Reads 进行基因注释。从 SAW v5.1.3 开始，count 允许重新利用部分 Multi-Mapped Reads 进行量化 (--multi\_map)，在 SAW 流程中，基因表达水平包括内含子和外显子的 MID 数量的总和。为了满足下游分析区分不同基因功能区分析的需求，count 会将外显子的 MID 计数单独输出。

### 2.5.1. 输入文件

- mapping 输出结果 BAM 文件 (.bam)
- 基因组注释 GFF/GTF 文件 (.gff / .gtf)

### 2.5.2. 命令示例及参数说明

单对 Q40 FSATQ 及 Q4 FASTQ 运行 count 参数如下：

```

1 $ mkdir -p /path/to/output/03.count
2 $ geneExp=/path/to/output/03.count/{SN}.raw.gcf
3 $ saturationSamplingFile=/path/to/output/03.count/{SN}_raw_barcode_gene_exp.txt
4 $ singularity exec SAW_v7.0.sif count \
5     -i /path/to/output/01.mapping/{lane}.Aligned.sortedByCoord.out.bam \ ## 必需参数:
mapping 比对的 bam 文件
6     -o /path/to/output/03.count/{SN}.Aligned.sortedByCoord.out.merge.q10.dedup.
target.bam \ ## 必需参数:去重合并的 bam 输出文件
7     -a /path/to/reference/genes.gtf \ ## 必需参数:基因组注释文件
8     -s /path/to/output/03.count/{SN}.Aligned.sortedByCoord.out.merge.q10.dedup.
target.bam.summary.stat \ ## 必需参数:bam 统计文件
9     -e /path/to/output/03.count/{SN}.raw.gcf \ ## 必需参数:设置表达量输出的文件
10 --umi_len 10 \ ## 必需参数:MID 的长度
11 --sat_file ${saturationSamplingFile} \ ## 可选参数:设置测序饱和度输出的文件
12 --sn {SN} \ ## 必需参数:芯片号
13 --umi_on \ ## 可选参数:纠正 MID
14 --save_lq \ ## 可选参数:保存低质量 reads
15 --save_dup \ ## 可选参数:保存重复 reads
16 -c 24 \ ## 可选参数:cpu 使用核数
17 --multi_map ## 可选参数:处理 multi-mapped reads 可添加该参数

```

☹ 注意！当有多对 Q40 FASTQ 运行 count 时，只需将上述参数中 /path/to/output/ 改成 /path/to/multi\_lane\_output/，{lane} 改成 {lane\*} 即可。

### 2.5.3. 运行资源

- 内存占用: ~45G (1G reads 数据分析参考值)
- 运行时间: ~30 min (1G reads 数据分析参考值)

### 2.5.4. 输出文件

```

1 /path/to/output/03.count
2 |— SN.Aligned.sortedByCoord.out.merge.q10.dedup.target.bam ## 按坐标排序的带注释的
   BAM 文件
3 |— SN.Aligned.sortedByCoord.out.merge.q10.dedup.target.bam.csi ## 带注释的 BAM 索
   引文件
4 |— SN.Aligned.sortedByCoord.out.merge.q10.dedup.target.bam.summary.stat##count 统
   计文件
5 |— SN_raw_barcode_gene_exp.txt ## 一个记录坐标、基因、MID 和计数信息的,以空格分隔的列表
6 |— SN.raw.gef ##HDF5 格式的基因表达文件

```

⊙ **注意!** 当使用添加了 rRNA 信息的参考基因组进行 mapping, 并且启用了 rRNAremove 参数, \*.bam.summary.stat 结果文件中的统计字段 PASS FILTER 将不包含匹配上 rRNA 参考序列的片段数目。

## 2.6. register

**register:** 可通过配准算法, 根据 track 线信息将上传的显微镜拍照的组织染色影像图与 count 输出的基因表达矩阵建立图像与空间表达的映射关系。

**染色图像包括:** DAPI/ssDNA 图像用于细胞核的染色; IF (免疫荧光) 图像用于蛋白质的染色; H&E (苏木精和伊红) 图像用于细胞核、细胞外基质和细胞质的染色。

SAW 的 register 包括四个主要模块: 拼接、组织分割、细胞分割和配准。拼接是指把带有重叠部分的显微镜图像拼成全景图像 (如果输入文件已经是全景图像则跳过此步骤)。组织分割和细胞分割可分别对组织和细胞覆盖区域进行识别, 并生成二值化掩模图。配准可对拼接图像和表达矩阵进行配准, 同时使用相同参数将组织和细胞分割二值化图与矩阵配准, 图像拼接和分割模块可以和配准模块分开运行。在处理 DAPI 和 mIF (multiple immunofluorescence images) 结合的多图场景时, 输出配准后的 DAPI 和 mIF 图像文件。其中, IF 图像的组织分割为阈值式分割结果, 细胞分割为 DAPI 组织分割区域、IF 阈值分割、DAPI 细胞分割的交集结果。

目前支持手动处理 QC 失败数据。组织或细胞分割图像将采用来自 ImageStudio 的手动处理结果 (接受 QC 失败和成功)。请查看 GitHub 上的手动教程。

⊙ **细胞分割是最耗时的部分。可以选择跳过细胞分割通过使用 -w False, 但仍然可以运行拼接、组织分割和配准。**

### 2.6.1. 输入文件

- count 输出的基因表达结果文件 (.raw.gef)
- 显微组织染色图像文件 (.tar.gz) 由 ImageStudio 处理。register 支持 ImageStudio v3 的质控输出和手动结果。
- 通过 ImageStudio 软件整理后的显微镜拍照影像图图像信息报告 (.ipr)

⊙ **GPU 使用的准备工作**

```
pip install onnxruntime-gpu==1.15.1
```

```
export LD_LIBRARY_PATH=/usr/local/cuda/lib64:$LD_LIBRARY_PATH
```

### 2.6.2. 命令示例及参数说明

场景 1: 处理来自 ImageStudio 的图像。执行拼接、组织分割、细胞分割, 并与基因表达矩阵配准。

```

1 $ image=/path/to/data/image
2 $ image4register=$(find ${image} -maxdepth 1 -name {SN}*.tar.gz | head -1)
3 $ imageStudio=$(find ${image} -maxdepth 1 -name {SN}*.ipr | head -1)
4 $ mkdir -p /path/to/output/04.register
5 $ singularity exec SAW_v7.0.sif register \

```



```

6     -i ${image4register} \ ## 必需参数:图像 QC 后产生的 tar.gz 文件
7     -c ${imageStudio} \ ## 必需参数:图像 QC 后产生的 IPR 文件
8     -v /path/to/output/03.count/{SN}.raw.gef \ ## 可选参数:输入基因表达矩阵 gef 文件
9     -o /path/to/output/04.register \ ## 必需参数:配准步骤输出文件目录
10    -w True ## 必需参数:做细胞分割

```

场景 2: 处理来自 ImageStudio 的图像。执行拼接、组织分割和细胞分割, 而不与基因表达矩阵进行配准。

```

1  $ image=/path/to/data/image
2  $ image4register=$(find ${image} -maxdepth 1 -name {SN}*.tar.gz | head -1)
3  $ imageStudio=$(find ${image} -maxdepth 1 -name {SN}*.ipr | head -1)
4  $ mkdir -p /path/to/output/04.register
5  $ singularity exec SAW_v7.0.sif register \
6     -i ${image4register} \ ## 必需参数:图像 QC 后产生的 tar.gz 文件
7     -c ${imageStudio} \ ## 必需参数:图像 QC 后产生的 IPR 文件
8     -o /path/to/output/04.register \ ## 必需参数:输出文件目录
9     -w True ## 必需参数:做细胞分割

```

场景 3: 处理来自场景 2 register 基因表达矩阵处理过的图像。

```

1  $ imageStudio=$(find /path/to/output/04.register -maxdepth 1 -name {SN}*.ipr |
2  head -1)
3  $ mkdir -p /path/to/output/04.register
4  $ singularity exec SAW_v7.0.sif register \
5     -i /path/to/output/04.register \ ## 必需参数:场景 2 的输出文件目录
6     -c ${imageStudio} \ ## 必需参数:场景 2 输出的 IPR 文件
7     -v /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:图像 QC 后产生的 IPR 文件
8     -o /path/to/output/04.register \ ## 必需参数:输出文件目录
9     -w False ## 如场景 2 已做细胞分割, 此处设置 False 避免再次启动自动细胞分割

```

### 2.6.3. 运行资源

- 内存占用: ~20G (1cm\*1cm 尺寸芯片数据分析参考值)
- 运行时间: ~7.5 h (1cm\*1cm 尺寸芯片数据分析参考值)

### 2.6.4. 输出文件

场景 1 和场景 3 的 register (如果 -i and -o 是相同的) 输出文件如下:

```

1  /path/to/output/04.register
2  └─ <stainType>_fov_stitched_transformed.tif ##TIFF 格式的已经与 track 线模板预配准的拼
   接全图
3  └─ SN_<chipType>_date_time_version.ipr ##IPR 格式图像处理记录文件

```

场景 2 的 register 输出文件:

```

1  /path/to/output/04.register
2  └─ <stainType>_fov_stitched_transformed.tif ##TIFF 格式的已经与 track 线模板预配准的拼
3  接全图
   └─ SN_<chipType>_date_time_version.ipr ##IPR 格式图像处理记录文件

```

场景 3 的 register 输出文件（如果 -i and -o 是两个不同的路径）：

```

1 /path/to/output/04.register
2 └─ SN_<chipType>_date_time_version.ipr ##IPR 格式图像处理记录文件

```

## 2.7. imageTools

**imageTools**：是一个专为处理 SAW 图像数据而设计的方便、实用的工具包。imageTools ipr2img（或 ipr2img）是 SAW 里将图像从 IPR 文件“解码”的必备核心工具。SAW imageTools ipr2img 可以从 IPR 文件输出 TIFF 格式文件，如配准前的 ssDNA 拼接图片、组织分割和细胞分割二值化掩膜图 TIFF 图片，也会对三类图像进行配准。

 请查看 [3.4 imageTools 其他应用](#)来了解更多关于 **imageTools** 的信息。

### 2.7.1. 输入文件

- ImageStudio 处理过的显微组织染色图像文件 (.tar.gz)
- register 处理过的图像处理记录文件 (.ipr)

### 2.7.2. 命令示例及参数说明

```

1 $ image=/path/to/data/image
2 $ image4register=$(find ${image} -maxdepth 1 -name {SN}*.tar.gz | head -1)
3 $ imageIPR=$(find /path/to/output/04.register -maxdepth 1 -name {SN}*.ipr | head -1)
4 $ singularity exec SAW_v7.0.sif imageTools ipr2img \
5     -i ${image4register} \ ## 必需参数:图像 QC 产生的 tar.gz 文件
6     -c ${imageIPR} \ ## 必需参数:register 运行后产生的 IPR 文件
7     -d tissue cell \ ## 必需参数:输出组织分割或细胞分割文件,以空格分割
8     -r True \ ## 必需参数:True:输出配准后的图像 ; False:输出预配准图像。
9     -o /path/to/output/04.register ## 必需参数:输出文件目录

```

### 2.7.3. 运行资源

- 内存占用：10G（1cm\*1cm 尺寸芯片数据分析参考值）
- 运行时间：5 min（1cm\*1cm 尺寸芯片数据分析参考值）

### 2.7.4. 输出文件

如果 -c IPR 的上级目录和 -o 是相同的，则输出文件为：

```

1 /path/to/output/04.register
2 └─ <stainType>_fov_stitched_transformed.tif ##TIFF 格式的已经与 track 线模板预配准的拼接全图
3 └─ <stainType>_matrix_template.txt ## 配准 DAPI/IF 图的 track 线交叉点模版.用于评估配准结果
4 └─ <stainType>_SN_mask.tif ##TIFF 格式的配准后的 DAPI/IF 细胞分割二值化图

```

```

5 |— <stainType>_SN_regist.tif  ##TIFF 格式的配准全景图
6 |— <stainType>_SN_tissue_cut.tif  ##TIFF 格式的 ssDNA/DAPI/IF/H&E 全图文件的组织分割结果
7 |— <stainType>_transform_template.txt  ##<stainType>_fov_stitched_transformed.tif
   的 track 线交叉点模板。用于评估拼接结果
8 |— fov_stitched_transformed.rpi  ## 已经与 track 线模板配准的拼接全图,支持存储多种 TIFF 格
   式的染色拼接全图
9 |— SN.rpi  ## 保存配准后的显微拍照全景图、组织边界、以及细胞边界(降采样)的图像金字塔
10 |— SN_<chipType>_date_time_version.ipr  ##IPR 格式图像处理记录文件

```

如果 -c IPR 的父级目录和 -o 是不同的路径,则输出文件为:

```

1 /path/to/output/04.register
2 |— <stainType>_fov_stitched_transformed.tif  ##TIFF 格式的已经与 track 线模板预配准的拼
   接全图
3 |— <stainType>_matrix_template.txt  ## 配准 ssDNA/DAPI/IF/H&E 图的 track 线交叉点模版。
   用于评估配准结果
4 |— <stainType>_SN_mask.tif  ##TIFF 格式的配准后的 ssDNA/DAPI/IF/H&E 细胞分割二值化图
5 |— <stainType>_SN_regist.tif  ##TIFF 格式的配准全景图
6 |— <stainType>_SN_tissue_cut.tif  ##TIFF 格式的 ssDNA/DAPI/IF/H&E 全图文件的组织分割结果
7 |— <stainType>_transform_template.txt  ##<stainType>_fov_stitched_transformed.tif
   的 track 线交叉点模板。用于评估拼接结果
8 |— SN.rpi  ## 保存配准后的显微拍照全景图、组织边界、以及细胞边界(降采样)的图像金字塔

```

## 2.8. tissueCut

**tissueCut:** 工具可根据 register 和 imageTools 输出的配准图勾画组织轮廓,并进行组织区域的识别和切割,从而去除组织外区域。如没有显微镜拍照的影像图, tissueCut 也可以直接基于基因表达矩阵识别组织区域。 tissueCut 提取得到的组织区域基因表达矩阵以 GEF 格式存储。

💡 如果 tissueCut 的结果与组织形态不符,有图场景下可在 ImageStudio 进行手动图像组织分割,之后重新运行 register、imageTools 和 tissueCut 进行矩阵提取。无图场景下可在 StereoMap 进行交互式 lasso 操作后接入 SAW 的 lasso 模块 (3.6 lasso) 来提取表达矩阵,也可以使用 Stereopy 进行交互式的 lasso 操作来提取表达矩阵。操作指南参见 [Stereopy->Tutorials->MaxtrixExport](#)。

### 2.8.1. 输入文件

- CID 对应 reads 数列表 (.txt)
- count 输出的基因表达结果文件 (.raw.gef)
- imageTools ipr2img 配准后的组织分割二值化掩膜图 (.tif, 可选)

## 2.8.2. 命令示例及参数说明

场景 1: 显微镜染色图经过 register 处理后, 运行 tissueCut:

```

1 $ mkdir -p /path/to/output/05.tissuecut
2 $ singularity exec SAW_v7.0.sif tissueCut \
   --dnbfile /path/to/output/02.merge/{SN}.merge.barcodeReadsCount.txt \ ## 可选
3 参数:每个 CID 上的 reads 数量列表文件
   -i /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:基因表达矩阵文件
4   -o /path/to/output/05.tissuecut \ ## 必需参数:tissueCut 输出目录
5   -s /path/to/output/04.register/<stainType>_SN_tissue_cut.tif \ ## 可选参数:输入
6   TIFF 格式组织分割二值化掩模图
   --sn {SN} \ ## 必需参数:芯片号
7   -O Transcriptomics \ ## 必需参数:组学具体的字符串
8   -d ## 必需参数:为 report 步骤必需参数

```

根据自定义标签区域的掩膜图像生成标签 GEF:

```

1 $ label= <label_name>## 自定义标签名字
2 $ labelMaskFile=$(find /path/to/output/04.register -maxdepth 1 -name *${label}
   *_tissue_cut.tif)
3 $ mkdir -p /path/to/output/05.tissuecut/tissuecut_${label}
4 $ singularity exec SAW_v7.0.sif tissueCut \
5   --dnbfile /path/to/output/02.merge/{SN}.merge.barcodeReadsCount.txt \ ##
   可选参数:每个 CID 上的 reads 数量列表文件
6   -i /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:基因表达矩阵文件
7   -o /path/to/output/05.tissuecut/tissuecut_${label} \ ## 必需参数:tissueCut
   输出目录
8   -s ${labelMaskFile} \ ## 可选参数:输入 TIFF 格式组织分割二值化掩模图
9   -l ${label} \ ## 可选参数:自定义标签名
10  --sn {SN} \ ## 必需参数:芯片号
11  -O Transcriptomics \ ## 必需参数:组学具体的字符串
12  -d ## 必需参数:为 report 步骤必需参数 ”

```

场景 2: 无配准图时, 运行 tissueCut:

```

1 $ mkdir -p /path/to/output/05.tissuecut
2 $ singularity exec SAW_v7.0.sif tissueCut \
3   --dnbfile /path/to/output/02.merge/{SN}.merge.barcodeReadsCount.txt \ ## 可选
   参数:每个 CID 上的 reads 数量列表文件
4   -i /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:基因表达矩阵文件
5   -o /path/to/output/05.tissuecut \ ## 必需参数:tissueCut 输出目录
6   --sn {SN} \ ## 必需参数:芯片号
7   -O Transcriptomics \ ## 必需参数:组学具体的字符串
8   -d ## 必需参数:为 report 步骤必需参数

```

### 2.8.3. 运行资源

- 内存占用：10G (1G reads 数据分析参考值)
- 运行时间：10 min (1G reads 数据分析参考值)

### 2.8.4. 输出文件

场景 1: 含有 -s 参数, 输入了配准后的组织图的输出文件列表  
根据相应的组织掩膜图像生成组织 GEF。

```

1 /path/to/output/05.tissuecut
2 |— SN.tissue.gef ##HDF5 格式的基因表达文件。组织 GEF 包括组织覆盖区域的表达信息
3 |— tissuecut.stat ## 组织覆盖区域的统计报告
4 |— tissue_fig ## 该目录存储组织区域覆盖统计图

```

根据自定义标签区域的掩膜图像生成标签 GEF:

```

1 /path/to/output/05.tissuecut/tissuecut_<label>
2 |— SN.<label>.raw.label.gef ##HDF5 格式的基因表达文件。
3 |— <label>.tissuecut.stat ## 组织覆盖区域的统计报告
4 |— tissue_fig ## 该目录存储组织区域覆盖统计图

```

场景 2: 不含 -s 参数, 未输入配准后的组织图的输出文件列表:

```

1 /path/to/output/05.tissuecut
2 |— 100X100_contour_image.png ##bin100 表达图
3 |— bin1_img.tif ##bin1 表达分布 TIFF 文件
4 |— bin1_img_tissue_cut.tif ##bin1 表达图中获取的组织掩膜图
5 |— SN.tissue.gef ##HDF5 格式的基因表达文件。组织 GEF 包括组织覆盖区域的表达信息
6 |— tissuecut.stat ## 组织覆盖区域的统计报告
7 |— tissue_fig ## 该目录存储组织区域覆盖统计图

```

## 2.9. spatialCluster

**spatialCluster:** 使用 Stereopy 进行空间聚类分析。

聚类过程包括 4 个步骤：

- (1) 组织覆盖区域基因表达数据预处理（对每个基因归一化、对数化、高变基因识别和标准化）；
- (2) PCA 降维；
- (3) 使用 UMAP 计算邻域图并做低维嵌入；
- (4) 使用 Leiden 算法进行聚类分析。

### 2.9.1. 命令示例

```

1 $ mkdir -p /path/to/output/06.spatialcluster
2 $ singularity exec SAW_v7.0.sif spatialCluster \
3     -i /path/to/output/05.tissuecut/{SN}.tissue.gcf \ ## 必需参数:组织区域 gcf 文件
4     -o /path/to/output/06.spatialcluster/{SN}.spatial.cluster.h5ad \ ## 必需参数:输出
    聚类 h5ad 文件路径
5     -r 1.0 \ ## 必需参数:Leiden 分辨率用于控制聚类的数量
6     -s 200 ## 必需参数:bin 的大小选择

```

### 2.9.2. 运行资源

- 内存占用：~5G（1G reads 数据分析参考值）
- 运行时间：~1 min（1G reads 数据分析参考值）

### 2.9.3. 输出文件

```

1 /path/to/output/06.spatialcluster
2 └─ SN.bin200_1.0.spatial.cluster.h5ad ## 空间聚类结果文件

```

## 2.10. cellCut

**cellCut:** 是基于从 register 和 imageTools 生成的配准图像提取细胞核表达矩阵的工具。cellCut 在 cellbin GEF 格式下输出表达数据。如果 cellbin 结果满意，可以运行 cellCut。

 请查阅 [3.1 SAW 工具和命令行参数说明 - 其他应用工具 - cellCut 的其他应用](#) 来学习更多关于 cellCut 的知识。

### 2.10.1. 输入文件

- count 工具输出基因表达矩阵文件 (**.raw.gef**)
- register 和 imageTools 工具输出细胞分割的二值化掩膜 TIFF 文件 (**.tif**)

### 2.10.2. 命令示例

```

1 $ mkdir -p /path/to/output/051.cellcut
2 $ singularity exec SAW_v7.0.sif cellCut cgef \
3     -i /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:原始 gef 文件
4     -m /path/to/output/04.register/<stainType>_SN_mask.tif \ ## 必需参数:TIFF 格式
    的配准后的 DAPI/IF 细胞分割二值化图
5     -o /path/to/output/051.cellcut/{SN}.cellbin.gef ## 必需参数:输出细胞分割后 gef 文件

```

### 2.10.3. 运行资源

- 内存占用: ~10G (1G reads 数据分析参考值)
- 运行时间: ~2 min (1G reads 数据分析参考值)

### 2.10.4. 输出文件

```

1 /path/to/output/051.cellcut
2 └─ SN.cellbin.gef ##HDF5 格式的细胞基因表达文件

```

## 2.11. cellCorrect

**cellCorrect:** 使用从 register 和 imageTools 生成的细胞分割图像进行调整, 并提取调整后的表达矩阵。cellCorrect 以 cell bin GEF 和 GEM 格式输出表达数据。如果需要细胞校正的结果, 请运行 cellCorrect。

### 2.11.1. 输入文件

- count 步骤输入的表达矩阵文件 (**.raw.gef**)
- register 和 imageTools 输出的细胞分割 mask 文件 (**.tif**)

### 2.11.2. 命令示例

```

1 $ nucleusLayer=$(find /path/to/output/04.register -maxdepth 1 -name *fov_stitched_
    transformed.tif -exec sh -c 'for f do basename -- "$f" _fov_stitched_transformed.
    tif;done' sh {} + | grep -v IF)
2 $ nucleusMask=$(find /path/to/output/04.register -maxdepth 1 -name ${nucleusLay-
    er}*_mask.tif)

```

```

3 $ singularity exec SAW_v7.0.sif cellCorrect \
4   -i /path/to/output/03.count/{SN}.raw.gcf \ ## 必需参数:原始 gcf 文件
5   -m ${nucleusMask} \ ## 必需参数:细胞分割 mask 图像文件
6   -d 10 \ ## 必需参数:基于细胞分割图像的细胞轮廓,以像素为单位的扩张距离
7   -o /path/to/output/051.cellcut ## 必需参数:输出 cellbin 的 gcf/gem 和调整后的 mask 文件。

```

### 2.11.3. 运行资源

- 内存占用: ~10G (1G reads 数据分析参考值)
- 运行时间: ~2 min (1G reads 数据分析参考值)

### 2.11.4. 输出文件

```

1 $ tree /path/to/output/051.cellcut
2 /path/to/output/051.cellcut
3 |— SN.adjusted.cellbin.gcf ## 细胞分割调整后 gcf 文件
4 |— <stainType>_SN_mask_edm_dis_10.tif ## TIFF 中调整的细胞边界图像文件
5 |— SN.adjusted.gem ## 细胞分割调整后的 gem 文件

```

## 2.12. cellCluster

**cellCluster:** 使用 Stereopy 进行细胞聚类, 该过程与 spatialCluster 类似。

### 2.12.1. 输入文件

- cellCut 输出 cellbin 基因表达矩阵文件 (**.cellbin.gcf**)

### 2.12.2. 命令示例

```

1 $ mkdir -p /path/to/output/061.cellcluster
2 $ singularity exec SAW_v7.0.sif cellCluster \
3   -i /path/to/output/051.cellcut/{SN}.adjusted.cellbin.gcf \ ## 必需参数:HDF5 格式
   的细胞基因表达文件
4   -o /path/to/output/061.cellcluster/{SN}.adjusted.cellbin.gcf ## 必需参数:细胞聚
   类 h5ad 文件路径

```

### 2.12.3. 运行资源

- 内存占用: ~5G (1G reads 数据分析参考值)
- 运行时间: ~5 min (1G reads 数据分析参考值)



### 2.12.4. 输出文件

```

1 /path/to/output/061.cellcluster
2 └─ SN.cell.cluster.h5ad  ## 细胞分割聚类产生的 h5ad 文件
3 └─ SN.adjusted.cell.cluster.h5ad  ## 细胞分割调整后聚类产生的 h5ad 文件

```

## 2.13. saturation

**saturation:** 工具用于计算组织覆盖区域的测序饱和度。

### 2.13.1. 输入文件

- mapping 输出 CID 比对统计文件 (**.stat**)
- count 输出计算饱和度所需抽样文件 (**.txt**)
- count 输出注释统计文件 (**.stat**)
- tissueCut 输出组织覆盖区域的 GEF 矩阵 (**.tissue.gef**)

### 2.13.2. 命令示例

```

1 $ mkdir -p /path/to/output/07.saturation
2 $ singularity exec SAW_v7.0.sif saturation \
3   -i /path/to/output/03.count/{SN}_raw_barcode_gene_exp.txt \  ## 必需参数:count 输出
   饱和度抽样文件
4   --tissue /path/to/output/05.tissuecut/{SN}.tissue.gef \  ## 必需参数:组织分割后
   gef 文件
5   -o /path/to/output/07.saturation \  ## 必需参数:输出文件目录
6   --bcstat /path/to/output/01.mapping/{lane}.CIDMap.stat \  ## 必需参数:CID 比对的统
   计文件 (多对 FASTQ 运行时,此处结果用逗号隔开)
7   --summary /path/to/output/03.count/{SN}.Aligned.sortedByCoord.out.merge.q10.
   dedup.target.bam.summary.stat  ## 必需参数:注释统计文件

```

### 2.13.3. 运行资源

- 内存占用: ~5G (1G reads 数据分析参考值)
- 运行时间: ~5 min (1G reads 数据分析参考值)

### 2.13.4. 输出文件

```

1 /path/to/output/07.saturation
2 └─ plot_1x1_saturation.png  ##bin1 的测序饱和度分析图
3 └─ plot_200x200_saturation.png  ##bin200 的测序饱和度分析图
4 └─ sequence_saturation.tsv  ## 测序饱和度文件

```

## 2.14. report

**report:** 工具可帮助用户整合每个步骤生成的分析报告，生成 JSON 格式的结果分析统计报告以及 HTML 网页版分析报告。分析报告整合基因的空间表达分布、关键统计指标、测序饱和度图、聚类分析结果、以及图像处理信息。mIF 场景下，聚类结果展示的底图增加伪彩（最多支持 7 种颜色）效果。

### 2.14.1. 输入文件

- mapping 输出 CID 比对统计文件 (.stat)、STAR 比对统计文件 (.out)
- count 输出注释统计文件 (.stat)
- register 处理过的图像记录文件和图像金字塔文件 (.ipr, .rpi)
- tissueCut 和 cellCut 输出 GEF 文件、组织覆盖区域统计文件 (.stat)、统计图 (.gef, .stat, .png)
- spatialCluster 和 cellCluter(若已完成) 输出聚类 H5AD 文件 (.h5ad)
- saturation 输出 bin200 的测序饱和度图 (.png)
- 物种、组织和参考信息

### 2.14.2. 命令示例

场景 1: 显微镜染色图像 register 后，同时 cellbin 文件都具备时，运行 report 的输出文件为：

```

1  $ imageIPR=$(find /path/to/output/04.register -maxdepth 1 -name {SN}*.ipr | head
  -1)
2  $ mkdir -p /path/to/output/08.report
3  $ singularity exec SAW_v7.0.sif report \
4      -m /path/to/output/01.mapping/{lane}.CIDMap.stat \ ## 必需参数:CID 比对的统计文件
5      -a /path/to/output/01.mapping/{lane}.Log.final.out \ ## 必需参数:STAR 比对统计
  文件
6      -g /path/to/output/03.count/{SN}.Aligned.sortedByCoord.out.merge.q10.dedup.
  target.bam.summary.stat \ ## 必需参数:带注释 bam 统计文件
7      -l /path/to/output/05.tissuecut/tissuecut.stat \ ## 必需参数:组织区域基因统计文件
8      -n /path/to/output/05.tissuecut/{SN}.gef \ ## 必需参数:所有 bin 表达 gef 文件
9      -d /path/to/output/06.spatialcluster/{SN}.bin200_1.0.spatial.cluster.h5ad \
  ## 必需参数:空间聚类文件
10     -t /path/to/output/07.saturation/plot_200x200_saturation.png \ ## 必需参数:
  bin200 测序饱和度分析图
11     -b /path/to/output/05.tissuecut/tissue_fig/scatter_200x200_MID_gene_counts.png
  \ ## 必需参数:每个 bin (bin 200) 中的 MID 计数和基因数的散点图
12     -v /path/to/output/05.tissuecut/tissue_fig/violin_200x200_MID_gene.png \ ## 必
  需参数:每个 bin (bin 200) 中的 MID 计数和基因数的小提琴图
13     -c /path/to/output/05.tissuecut/tissue_fig/statistic_200x200_MID_gene_DNB.png \
  ## 必需参数:x 轴含毛边的 MID 数、基因数和 DNB 数的单变量分布图 (bin 200)
14     --bin20Saturation /path/to/output/05.tissuecut/tissue_fig/scatter_20x20_MID_
  gene_counts.png \ ## 必需参数:每个 bin (bin 20) 中的 MID 计数和基因数的散点图
15     --bin20violin /path/to/output/05.tissuecut/tissue_fig/violin_20x20_MID_gene.
  png \ ## 必需参数:每个 bin (bin 20) 中的 MID 计数和基因数的小提琴图
16     --bin20MIDGeneDNB /path/to/output/05.tissuecut/tissue_fig/statistic_20x20_MID_
  gene_DNB.png \ ## 必需参数:x 轴含毛边的 MID 数、基因数和 DNB 数的单变量分布图 (bin 20)
17     --bin50Saturation /path/to/output/05.tissuecut/tissue_fig/scatter_50x50_MID_
  gene_counts.png \ ## 必需参数:每个 bin (bin 50) 中的 MID 计数和基因数的散点图

```

```

18 --bin50violin /path/to/output/05.tissuecut/tissue_fig/violin_50x50_MID_gene.
   png \ ## 必需参数:每个 bin (bin 50) 中的 MID 计数和基因数的小提琴图
19 --bin50MIDGeneDNB /path/to/output/05.tissuecut/tissue_fig/statistic_50x50_MID_
   gene_DNB.png \ ## 必需参数:x 轴含毛边的 MID 数、基因数和 DNB 数的单变量分布图 (bin 50)
20 --bin100Saturation /path/to/output/05.tissuecut/tissue_fig/scatter_100x100_
   MID_gene_counts.png \ ## 必需参数:每个 bin (bin 100) 中的 MID 计数和基因数的散点图
21 --bin100violin /path/to/output/05.tissuecut/tissue_fig/violin_100x100_MID_
   gene.png \ ## 必需参数:每个 bin (bin 100) 中的 MID 计数和基因数的小提琴图
22 --bin100MIDGeneDNB /path/to/output/05.tissuecut/tissue_fig/statistic_100x100_
   MID_gene_DNB.png \ ## 必需参数:x 轴含毛边的 MID 数、基因数和 DNB 数的单变量分布图 (bin 100)
23 --bin150Saturation /path/to/output/05.tissuecut/tissue_fig/scatter_150x150_
   MID_gene_counts.png \ ## 必需参数:每个 bin (bin 150) 中的 MID 计数和基因数的散点图
24 --bin150violin /path/to/output/05.tissuecut/tissue_fig/violin_150x150_MID_
   gene.png \ ## 必需参数:每个 bin (bin 150) 中的 MID 计数和基因数的小提琴图
25 --bin150MIDGeneDNB /path/to/output/05.tissuecut/tissue_fig/statistic_150x150_
   MID_gene_DNB.png \ ## 必需参数:x 轴含毛边的 MID 数、基因数和 DNB 数的单变量分布图 (bin 150)
26 --cellBinGef /path/to/output/051.cellcut/{SN}.adjusted.cellbin.gef \ ## 可选参
   数:细胞分割 gef 文件
27 --cellCluster /path/to/output/061.cellcluster/{SN}.adjusted.cell.cluster.h5ad
   \ ## 可选参数:细胞分割后聚类结果
28 -i /path/to/output/04.register/{SN}.rpi \ ## 可选参数:芯片 rpi 文件
29 -r standard_version \ ## 必需参数:报告版本
30 -s {SN} \ ## 必需参数:芯片号
31 --pipelineVersion SAW_v7.0 \ ## 必需参数:流程版本号
32 --iprFile ${imageIPR} \ ## 可选参数:图像预处理文件
33 --species {species_name} \ ## 必需参数:物种信息
34 --tissue {tissue_type} \ ## 必需参数:组织类型
35 --reference {reference_index} \ ## 必需参数:参考基因组
36 -o /path/to/output/08.report ## 必需参数:输出文件目录

```

⊙ 注意: 把 {species\_name}, {tissue\_type}, {reference\_index} 替换成真实的信息。

场景 2: 当有多对 Q40 FASTQ 运行 report 时。只需将上述参数中 /path/to/output/ 改成 /path/to/multi\_lane\_output/, {lane} 改成 {lane\*} 即可。

场景 3: 使用 register 的 -w False 参数没有细胞分割结果, 但与矩阵配准了的图像输出 report。只需将上述参数中 “--cellBinGef, --cellCluster” 两个参数去掉即可。

场景 4: 无配准图时输出 report。只需将上述参数中 “--cellBinGef, --cellCluster, -i, --iprFile” 四个参数去掉即可。

### 2.14.3. 运行资源

- 内存占用: ~1G (1G reads 数据分析参考值)
- 运行时间: ~1 min (1G reads 数据分析参考值)

### 2.14.4. 输出文件

有 cellbin 时, report 输出文件如下:

```

1 /path/to/output/08.report
2 |— scatter_1x1_MID_gene_counts.png ## 每个细胞的 CID 计数和基因数的散点图
3 |— SN.report.html ## 报告 html 文件
4 |— SN.statistics.json ## 所有重要数据统计文件

```

```
5 |— statistic_1x1_cell_area.png  ## 细胞面积沿每个细胞的 x 轴的单变量分布。
6 |— statistic_1x1_DNB.png  ## DNB 数沿每个细胞 x 轴的单变量分布。
7 |— statistic_1x1_gene.png  ## 基因类型沿每个细胞 x 轴的单变量分布。
8 |— statistic_1x1_MID.png  ## MID 计数沿每个细胞 x 轴的单变量分布。
9 |— violin_1x1_gene.png  ## 小提琴图显示了每个细胞中基因类型的分布。
10 |— violin_1x1_MID.png  ## 小提琴图显示了每个细胞中重复数据删除的 MID 计数的分布。
```

无 cellbin 时, report 输出文件如下:

```
1 /path/to/output/08.report
2 |— SN.report.html  ## 报告 html 文件
3 |— SN.statistics.json  ## 所有重要数据统计文件
```

## 第三章

# SAW 其他应用工具

## 3.1. cellCut 其他应用

### 3.1.1. 工具简介

**cellCut:** 是一个可用来处理 GEF 文件的工具。SAW 流程中的这个工具是用来把因表达矩阵的 GEF 格式转变为清晰图表或者补全 GEF 格式。用户也可以使用 C++ 编译好的工具 `geftools` 或 python 封装的程序包 `gefpy` 来处理 GEF 文件。

### 3.1.2. 功能示例

功能 1: GEF 转换为 GEM 格式矩阵

```

1  $ singularity exec SAW_v7.0.sif cellCut view \ ## convert GEF that only contains
   bin1 geneExp
2      -s {SN} \ ## 必需参数:芯片号
3      -i /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:gef 文件
4      -o {SN}.raw.gem ## 可选参数:输出 gem 文件
5  $ singularity exec SAW_v7.0.sif cellCut view \ ## convert a whole GEF
6      -s {SN} \ ## 必需参数:芯片号
7      -i /path/to/output/05.tissuecut/{SN}.gef \ ## 必需参数:gef 文件
8      -o {SN}.gem ## 可选参数:输出 gem 文件
9  $ singularity exec SAW_v7.0.sif cellCut view \ ## convert tissue GEF that only
   contains bin1 geneExp
10     -s {SN} \ ## 必需参数:芯片号
11     -i /path/to/output/05.tissuecut/{SN}.tissue.gem \ ## 必需参数:gef 文件
12     -o {SN}.tissue.gem ## 可选参数:输出 gem 文件
13 $ singularity exec SAW_v7.0.sif cellCut view \ ## convert cellbin GEF to cellbin
   GEM
14     -s {SN} \ ## 必需参数:芯片号
15     -i /path/to/output/051.cellcut/{SN}.cellbin.gem \ ## 必需参数:cellbin 的 gef 文件
16     -o {SN}.cellbin.gem \ ## 必需参数:cellbin 的 gem 文件
17     -d /path/to/output/03.count/{SN}.raw.gem ## 可选参数:gef 文件

```

功能 2: GEF 的补全

```

1  $ singularity exec SAW_v7.0.sif cellCut bgef \ ## complete GEF that only contains
   bin1 geneExp group to a whole GEF, you may specify the bin size you need using "-b".
2  Separate multiple bin size with comma
3      -i /path/to/output/05.tissuecut/{SN}.tissue.gem \ ## 必需参数:gef 文件
4      -o {SN}.tissue.complete.gem \ ## 必需参数:输出 gef 文件
5      -b 1,20,50,100 \ ## 必需参数:binsize 大小
6      -O Transcriptomics ## 必需参数:组学具体名称

```

功能 3: GEM 转换为 GEF

```

1  $ singularity exec SAW_v7.0.sif cellCut bgef \ ## convert GEM to GEF in specific
   bin size. Separate multiple bin sizes with comma
2      -i {SN}.gem \ ## 必需参数:输入 gem 文件
3      -o {SN}.gef \ ## 必需参数:输出 gef 文件
4      -b 1,20,50 \ ## 必需参数:binsize 大小
5      -O Transcriptomics ## 必需参数:组学具体名称

```

## 3.2. checkGTF 应用

### 3.2.1. 工具简介

**checkGTF**: 是用于检查作为 count 输入文件的 GTF/GEF 的格式是否正确。特别是基因名称长度的限制。

### 3.2.2. 功能示例

```

1 $ singularity exec SAW_v7.0.sif checkGTF \
2   -i /path/to/reference/genes.gtf \ ## 必需参数:基因注释文件 GTF/GFF。
3   -o /path/to/reference/genes_new.gtf ## 必需参数:输出新形式的 GFF/GTF 文件

```

## 3.3. imageTools 其他应用

### 3.3.1. 工具简介

除了 ipr2img, imageTools 还有三个功能

- (1) **img2rpi** (以 rpi 格式写入图像) ;
- (2) **merge** (将染色图像与组织分割和细胞分割图像融合, 以检查分割结果) ;
- (3) **overlay** (将推断的 track 线模板与拼接的全景图像重叠以检查拼接结果, 或将矩阵中的 track 线模板与配准的全景图像重叠, 以检查配准结果) 。

### 3.3.2. 功能示例

功能 1: img2rpi

```

1 $ singularity exec SAW_v7.0.sif imageTools img2rpi \
2   -i /path/to/output/04.register/<stainType1>_fov_stitched_transformed.tif,/
   path/to/output/04.register/<stainType2>_fov_stitched_transformed.tif \ ## 必需参数:
   图像 TIFF 文件(多个图像文件用逗号隔开)
3   -g <stainType1>/<ImageType>,<stainType2>/<ImageType> \ ## 必需参数:设置输入保存在
   rpi 文件里的命名, 芯片号或者图像类型
4   -b 1 10 50 100 \ ## 必需参数:binsize 大小
5   -o /path/to/output/04.register/fov_stitched_transformed.rpi ## 必需参数:RPI 输出
   文件路径,RPI 文件可在 StereoMap 中用来手动配准

```

功能 2: merge (-i 任意最多三个灰度图像 (.tif) )

```

1 $ singularity exec SAW_7.0.sif imageTools merge \
2   -i /path/to/output/04.register/<stainType>_{SN}_regist.tif,/path/to/output/04.
   register/{SN}_tissue_cut.tif,/path/to/output/04.register/{SN}_mask.tif \ ## 必需参
   数:输入 TIFF 文件, 多个文件以逗号分割, 最多三个
3   -o /path/to/output/04.register/merge.tif ## 必需参数:多通道图像合并输出文件路径

```

## 功能 3: overlay

1. 拼接效果检查: 将 IPR 或 TXT 格式的拼接模板与预先配准的全景图像叠加。

```

1  ## stitch
2  $ preRegImage=/path/to/output/04.register/<stainType>_fov_stitched_transformed.tif
3  $ imageIPR=$(find /path/to/output/04.register -maxdepth 1 -name {SN}*.ipr | head
  -1)
4  $ stitchTplt=/path/to/output/04.register/<stainType>_transform_template.txt

5  $ singularity exec SAW_7.0.sif imageTools overlay \
6    -i ${preRegImage} \ ## 必需参数:预配准图像文件路径
7    -c ${imageIPR} \ ## 必需参数:IPR 文件
8    -o /path/to/output/04.register/<stainType>_stitch_check.tif \ ## 必需参数:输出
  TIFF 图像路径(拼接检查文件)
9    -d stitch \ ## 必需参数:模块名
10   -l 60 \ ## 可选参数:Cross limbs length in pixel
11   -w 3 ## 可选参数:Cross line thickness in pixel

12 $ singularity exec SAW_7.0.sif imageTools overlay \
13   -i ${preRegImage} \ ## 必需参数:预配准图像文件路径
14   -c ${stitchTplt} \ ## 必需参数:txt 文件
15   -o /path/to/output/04.register/<stainType>_stitch_check_tplt.tif \ ## 必需参数:
  输出 TIFF 图像路径(拼接检查文件)
16   -d stitch \ ## 必需参数:模块名
17   -l 60 \ ## 可选参数:Cross limbs length in pixel
18   -w 3 ## 可选参数:Cross line thickness in pixel

```

2. 配准效果检查: 将 IPR 或 TXT 格式的配准模板与已配准的全景图像叠加。

```

1  ## register
2  regImage=/path/to/output/04.register/<stainType>_{SN}_regist.tif
3  imageIPR=$(find /path/to/output/04.register -maxdepth 1 -name {SN}*.ipr | head -1)
4  regTplt=/path/to/output/04.register/<stainType>_transform_template.txt

5  singularity exec SAW_7.0.sif imageTools overlay \
6    -i ${regImage} \ ## 必需参数:已配准图像文件路径
7    -c ${imageIPR} \ ## 必需参数:IPR 文件
8    -o /path/to/output/04.register/register_check.tif \ ## 必需参数:输出 TIFF 图像路径
  (配准检查文件)
9    -d register \ ## 必需参数:模块名

```



```

10     -l 60 \ ## 可选参数:Cross limbs length in pixel
11     -w 3 ## 可选参数:Cross line thickness in pixel

12 singularity exec SAW_7.0.sif imageTools overlay \
13     -i ${regImage} \ ## 必需参数:已配准图像文件路径
14     -c ${regTplt} \ ## 必需参数:txt 文件
15     -o /path/to/output/04.register/register_check_tmplt.tif \ ## 必需参数:输出 TIFF
    图像路径(配准检查文件)
16     -d register \ ## 必需参数:模块名
17     -l 60 \ ## 可选参数:Cross limbs length in pixel
18     -w 3 ## 可选参数:Cross line thickness in pixel

```

功能 4: ipr2img

从 IPR 文件输出图像 (.tif) 。有关此函数的其他用途, 请参考 Github 手动处理流程的教程。

## 3.4. manualRegister 应用

### 3.4.1. 工具简介

当自动配准效果不佳时, 用户可以使用线下可视化软件 StereoMap 进行手动图像和矩阵的配准。配准后的操作参数记录在 JSON 文件中, 可输入进 SAW 的 manualRegister 工具进行图像调整的运算并修改 IPR 中的配准参数记录。用户需要重新运行 imageTools ipr2img 获得修改后的图像。

### 3.4.2. 功能示例

💡 选择 1: StereoMap 输出的记录手动配准参数信息 JSON 文件的参数信息与 manualRegister 入参的对应关系:

- "offsetX": -o 的 {offsetX};
- "offsetY": -o 的 {offsetY};
- "flip": -f 的 {flip};
- "rotate": -r 的 {rotate};
- "extrude\_x": -s 的 {extrude\_x};
- "extrude\_y": -s 的 {extrude\_y}。
- "version": {isReversed} of -w 如果这个键在 JSON 中不存在, 则输入 "True"。如果存在, 则输入 "False"。当缺少 -m 参数时, 该参数才会生效。

选择 2: 通过 -m 参数来使用 StereoMap 输出的 JSON 文件。

场景 1 和场景 2: 进行 track 线微调 && 不进行 track 线微调

```

1 $ mkdir -p /path/to/output/manualregister
2 $ cp $(find /path/to/output/04.register -maxdepth 1 -name {SN}*.ipr | head -1) /
    path/to/output/manualregister # we recommend to make a copy of IPR to prevent
    overwrite original file
3 $ regIPR=$(find /path/to/output/manualregister -maxdepth 1 -name {SN}*.ipr | head
    -1)
4 $ singularity exec SAW_v7.0.sif manualRegister \
5     -i /path/to/output/04.register \ ## 必需参数:register 的结果输出目录
6     -c ${regIPR} \ ## 必需参数:register 的 IPR 输出文件
7     -v /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:count 输出 GEF 表达矩阵

```

```

8     -o {offsetX} {offsetY} \ ## 可选参数:fov_stitched_transformed.tif 图像中心的 x 和 y
    方向偏移量
9     -f {flip} \ ## 可选参数:是否沿 y 轴水平翻转图像
10    -r {rotate} \ ## 可选参数:从 fov_stitched_transformed.tif 图像中心开始的手动旋转角度
11    -p /path/to/output/manualregister ## 必需参数:输出保存 IPR 结果目录

```

场景 3: 只用 StereoMap 输出 JSON 文件

```

1 $ mkdir -p /path/to/output/manualregister
2 $ cp $(find /path/to/output/04.register -maxdepth 1 -name {SN}*.ipr | head -1) /
    path/to/output/manualregister # we recommend to make a copy of IPR to prevent
    overwrite original file
3 $ regIPR=$(find /path/to/output/manualregister -maxdepth 1 -name {SN}*.ipr | head
    -1)
4 singularity exec SAW_v7.0.sif manualRegister \
5     -i /path/to/output/04.register \ ## 必需参数:register 的结果输出目录
6     -c ${regIPR} \ ## 必需参数:register 的 IPR 输出文件
7     -v /path/to/output/03.count/{SN}.raw.gef \ ## 必需参数:count 输出 GEF 表达矩阵
8     -m /path/to/datetime.regist.json \ ## 必需参数:JSON 文件
9     -p /path/to/output/manualregister ## 必需参数:输出保存 IPR 结果目录

```

## 3.5. lasso 应用

### 3.5.1. 工具简介

**lasso**: 根据用户在 StereoMap 中手动圈选出的一个或多个闭合区域进行对应区域的基因表达矩阵的提取。

### 3.5.2. 功能示例

场景 1 及场景 2: square bin lasso && cellbin lasso

```

1 $ mkdir -p /path/to/output/lasso
2
3 $ singularity exec SAW_v7.0.sif lasso \
4     -i /path/to/output/03.count/{SN}.gef \ ## 必需参数:gef 文件(如果是 cellbin lasso,
    此处为 /path/to/output/051.cellcut/{SN}.cellbin.gef)
5     -m /upload/path/{taskID}.lasso.geojson \ ## 必需参数:stereomap 输出坐标列表文件
    (*.geojson)
6     -o /path/to/output/lasso \ ## 必需参数:保存 lasso 输出文件的目录
7     -s {binList} \ ## 必需参数:bin 的大小
8     -n {SN} ## 必需参数:芯片号

```

```

9  ## 对 cellbin.gef 进行细胞聚类,并在 StereoMap 中显示聚类结果。
10 $ singularity exec SAW_v7.0.sif cellCluster \
11   -i /path/to/output/<label>/{SN}.<label>.label.cellbin.gef \
12   -o /path/to/output/<label>/{SN}.<label>.label.cell.cluster.h5ad
13 ## 预计算渲染数据,并可以在 StereoMap 中显示 cellbin.gef 表达热图。
14 $ singularity exec SAW_v7.0.sif cellChunk \
15   -i /path/to/output/<label>/{SN}.<label>.label.cellbin.gef \
16   -o /path/to/<label>

```

☹ 注意,如果是场景 2: cellbin lasso, 只需将上述参数中“-s”去掉即可。

- Square bin lasso 输出文件如下:

```

1  tree /path/to/output/lasso
2  /path/to/output/lasso
3  └─ <label>
4     └─ SN.<label>.label.gef
5        └─ segmentation
6           └─ SN.lasso.<binList[0]>.<label>.gem.gz...
7              └─ SN.lasso.<binList[m]>.<label>.gem.gz
8                 └─ SN.lasso.<label>.mask.tif

```

- Cellbin lasso 输出文件如下:

```

1  tree /path/to/output/lasso
2  /path/to/output/lasso
3  └─ <label>
4     └─ SN.<label>.label.cellbin.gef

```

# 第四章

## SAW 常见 Q&A

## 4.1. Q: 基因组注释文件 GTF/GFF 有什么格式要求?

### 4.1.1. 文件格式:

GFF文件 或者 GTF文件, 文件后缀名支持 gtf/gtf.gz, gff/gff.gz, gff3/gff3.gz

### 4.1.2. GTF 文件格式:

注释行以 # 开始

主体部分共 9列, 以tab作为分隔符: seqname source feature start end score strand frame attributes

- type: 注释信息的类型必须含有gene,transcript 和 exon
- start/end: 最大值需小于 $2^{31}$
- strand: 链的正向与负向, 分别用加号+和减号-表示。
- 第9列为attributes, 格式为tag “value” (标签“值”), 不同属性之间以空格相隔; 必须要有以下4个
  - ①. gene\_name value
  - ②. gene\_id value: 表示转录本在基因组上的基因座的唯一的ID。gene\_id与value值用空格分开, 如果值为空, 则表示没有对应的基因。
  - ③. transcript\_name value
  - ④. transcript\_id value: 预测的转录本的唯一ID。transcript\_id与value值用空格分开, 空表示没有转录本。
- 目前最大有效基因数必须小于 $2^{20}$ , 即1048576
- 不可以乱序, 即同一个gene的transcript/exons需按顺序排列

### 4.1.3. GFF 文件格式:

注释行以 # 开始

主体部分共 9列, 以tab作为分隔符: seqid source type start end score strand phase attributes

- type: 注释信息的类型必须含有gene, mRNA和 exon;
- start/end: 最大值需小于 $2^{31}$ ;
- strand: “+”表示正链, “-”表示负链, “.”表示不需要指定正负链, “?”表示未知;
- 第9列为attributes, 格式为tag=value (标签=值), 不同属性之间以分号相隔;

- ①. 需要存在ID Name Parent(对gene无需判断Parent);
- ②. 对于第三列的命名规则请务必仔细研究 ⇒ “树状分级” (不能只列出child行 而没有parent行!) 示例如下:

```

1 ##sequence-region ctg123 1 1497228
2 ctg123 . gene 1000 9000 . + . ID=gene00001;Name=EDEN
3 ctg123 . TF_binding_site 1000 1012 . + . ID=tfbs00001;Parent=gene00001
4 ctg123 . mRNA 1050 9000 . + . ID=mRNA00001;Parent=gene00001;Name=EDEN.1
5 ctg123 . mRNA 1050 9000 . + . ID=mRNA00002;Parent=gene00001;Name=EDEN.2
6 ctg123 . mRNA 1300 9000 . + . ID=mRNA00003;Parent=gene00001;Name=EDEN.3
7 ctg123 . exon 1500 1500 . + . ID=exon00001;Parent=mRNA00001
8 ctg123 . exon 1050 1500 . + . ID=exon00002;Parent=mRNA00001,mRNA00002
9 ctg123 . exon 3000 3902 . + . ID=exon00003;Parent=mRNA00001,mRNA00003
10 ctg123 . exon 5000 5500 . + . ID=exon00004;Parent=mRNA00001,mRNA00002,mRNA00003
11 ctg123 . exon 7000 9000 . + . ID=exon00005;Parent=mRNA00001,mRNA00002,mRNA00003
12 ctg123 . CDS 1201 1500 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
13 ctg123 . CDS 3000 3902 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
14 ctg123 . CDS 5000 5500 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
15 ctg123 . CDS 7000 7600 . + 0 ID=cds00001;Parent=mRNA00001;Name=edenprotein.1
16 ctg123 . CDS 1201 1500 . + 0 ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
17 ctg123 . CDS 5000 5500 . + 0 ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
18 ctg123 . CDS 7000 7600 . + 0 ID=cds00002;Parent=mRNA00002;Name=edenprotein.2
19 ctg123 . CDS 3301 3902 . + 0 ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
20 ctg123 . CDS 5000 5500 . + 1 ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
21 ctg123 . CDS 7000 7600 . + 1 ID=cds00003;Parent=mRNA00003;Name=edenprotein.3
22 ctg123 . CDS 3391 3902 . + 0 ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
23 ctg123 . CDS 5000 5500 . + 1 ID=cds00004;Parent=mRNA00003;Name=edenprotein.4
24 ctg123 . CDS 7000 7600 . + 1 ID=cds00004;Parent=mRNA00003;Name=edenprotein.4

```

- 目前最大有效基因数必须小于 $2^{20}$ ,即1048576;
- 可以乱序,但仍需满足 gene必须出现在对应mRNA之前,mRNA必须出现在对应的exon之前的规则。

#### 4.1.4. 其他注意事项:

gene/gene\_name(基因的名字) 值不含有特殊符号(空格,各类型括号,引号,<>,%等)。支持使用的常见特殊符号有“\_”,“.”;

gene/gene\_name(基因的名字) 值长度小于64个字符;

虽然GFF文件现在大部分使用的都是第三版(GFF3),但是文件命名时请命名为.gff;同理对于GTF文件也请文件命名时采用.gtf。

## 4.2. Q: 读取注释文件时会忽略对应基因的情况有哪些?

gtf格式的属性没有gene\_name gene\_id transcript\_name transcript\_id(对gene只需要有gene\_name和gene\_id);

gff格式的属性没有ID Name Parent(对gene无需判断Parent);

同一个gene下的数据包含多个gene\_id,日志打印 "Multiple gene IDs for gene xxx: id1, id2...";

同一个gene下的数据同时包含正反链,日志打印 "Strand disagreement for gene xxx - skipping";

transcript或exon没有transcript\_id,日志打印 "Record does not have transcriptID for gene xxx";

同一个gene有多条transcripts的transcript\_id / ID相同,日志打印 "Transcript appears more than once for xxx";

存在exon的start > end,日志打印 "Exon has 0 or negative extent for xxx";

同一个transcript下的exons之间有overlap,日志打印 "Exons overlap for xxx";

一个gene没有任何transcript,日志打印 "No transcript for gene xxx".

\* 一个contig下出现多条gene有相同的gene\_name,合并为一个gene。

### 4.3. Q: 构建 reference 时报错 "Fatal INPUT FILE error, no valid exon lines in the GTF file" 如何处理?

可能是注释GTF/GFF文件和基因组FASTA文件中两者对染色体的命名不完全统一, 请注意chromosome name要统一。

### 4.4. Q: 为什么大部分的注释文件中的基因都未被注释上?

多数情况是注释文件不规范导致的, 请再次参考上述文件格式要求的内容自行排查;

另一种可能性是由于strand正负链符号不规范导致, 注释文件中strand值只能是“+” (forward) 或“-” (reverse), 请不要和下划线“-”搞混。

### 4.5. Q: 是否有工具可以辅助检查注释文件规范?

可以使用SAW sif中的checkGTF工具进行检查, 运行方式如下:

```
1 ## export SINGULARITY_BIND="/path/to/input/dir,/path/to/output/dir"
2 singularity exec SAW_v7.0.sif checkGTF \
3     -i <input.gtf/gff> \ ## GTF/GFF file input to be checked
4     -o <output.gtf/gff> ## [optional]. Set to output revised GTF/GFF file. Be
    aware that this may remove some genes which do not meet the requirements and cannot
    be fixed.
```

可能有部分无法被程序修复的基因注释记录会在输出时被删除, 可以在日志中找到记录重新修改GTF文件后再次尝试。

### 4.6. Q: SAW 中对测序数据做了哪些过滤?

**CID过滤:** 过滤CID无法比对上mask文件的reads;

**MID过滤:** 过滤MID序列中含有N碱基的reads, 过滤MID为polyA的reads, 过滤含有至少一个以上质量值低于10的碱基的reads;

**reads过滤:** 过滤含有接头和DNB序列的reads。

### 4.7. Q: 基因表达可视化结果异常, 没有组织轮廓怎么办?

第一步: 检查HTML报告中Valid CID Reads的比例是否正常, 不可低于10%, 如低于10%请确认测序FASTQ文件和芯片SN对应;

第二步: 如Valid CID Reads的比例较低, 在10%~30%左右, 可能存在两种可能性:

- 参考基因组格式异常: Multi-Mapped Reads比例高, Uniquely Mapped Reads比例很低且几乎全部注释失败, 需要检查注释使用的GTF/GFF文件是否符合格式要求, 是否可以通过检测程序;
- 存在混样的可能性: 需要自行排查实验部分。

## 4.8. Q: 如何解析 GEF 格式文件?

Option1:使用C++编译的 geftools:

- <https://github.com/STOmics/geftools>

Option2:使用python包 gefpy:

- <https://pypi.org/project/gefpy/>
- <https://gefpy.readthedocs.io/en/latest/index.html>

Option3:如已安装SAW sif(如v7.0):

- <https://hub.docker.com/repository/docker/stomics/saw>
- `singularity exec SAW_v7.0.sif cellCut`
- 请使用singularity 3.8及以上版本

```
1 export HDF5_USE_FILE_LOCKING=FALSE
2 ## gef2gem using geftools
3 geftools view -i <SN>.gef -o <SN>.gem -s <SN>
4 # -i input square bin GEF, e.g.SN.raw.gef or SN.gef
5 # -o output GEM
6 # -s SN

7 ## gef2gem using gefpy
8 python
9 >>> from gefpy.bgef_reader_cy import BgefR
10 >>> bgef=BgefR(filepath='<SN>.gef',bin_size=200,n_thread=4)
11 >>> bgef.to_gem('<SN>.bin200.gem')

12 ## gef2gem using SAW sif
13 ## export SINGULARITY_BIND="/path/to/input/dir,/path/to/output/dir"
14 singularity exec SAW_v7.0.sif cellCut view -i <SN>.gef -o <SN>.gem -s <SN>

15 ## gem2gef
16 geftools bgef -i <SN>.gem -o <SN>.gef -b 1,20,50 -O Transcriptomics
17 # -i input square bin GEM
18 # -o output square bin GEF
19 # -b bin sizes searate by comma, default: 1,10,20,50,100,200,500
20 # -O omics name
```



## 4.9 Q: 测序在进行 mapping 比对时, FASTQ 太多如何简易处理?

如果需要处理的FASTQ文件太多,更简单的方法是将它们整理成一个FQ\_{index}.list文件。FQ\_{index}.list的要求是列表文件是将所有索引与拆分mask相同的FASTQs在一起,因为这些文件都是由相同的逻辑拆分的。

示例如下:

Q4 FASTQ name

**/path/to/data/E100026571\_L01/barcode\_2/E100026571\_L01\_2\_16.fq.gz**

lane

barcode

lane

barcode

split index

```
1 $ cat SN_16.list ## a FASTQ list file gathers all the FASTQs whose index is 16
2 /path/to/data/lane_1_16.fq.gz
3 /path/to/data/lane_2_16.fq.gz
```

然后在 {index}.bcPara文件的in1参数输入FQ\_{index}.list文件 的路径。FQ\_{index}.list 的{index}和分割mask文件的 {index} 务必是相同的。

```
1 $ mkdir /path/to/output/01.mapping
2 $ vim /path/to/output/01.mapping/{index}.bcPara
3 in=/path/to/output/00.splitmask/{index}.{SN}.barcodeToPos.bin ## split mask file
4 in1=/path/to/output/{SN}_{index}.list
5 barcodeReadsCount=/path/to/output/01.mapping/{index}.barcodeReadsCount.txt
6 barcodeStart=0
7 barcodeLen=24
8 umiStart=25
9 umiLen=10
10 mismatch=1
11 bcNum=38284877 ## Input the first line from output of CIDCount
12 polyAnum=15
13 mismatchInPolyA=2
```

## 4.10 Q: Valid CID 比例异常应该怎么处理?

Valid CID rate低可以排查的方向有3个:

1. 测序情况。测序质量低会影响比对结果。除了Q30还需要检查call N的情况,可以查看下机报告的base distribution,如果出现N碱基的比例高的情况,就需要考虑是因为测序问题,影响了valid CID rate,最好优先排查。
2. 芯片mask h5文件和FQ不对应。因为mask里记录的CID和对样本测序得到的CID不匹配,导致valid CID rate低。这个情况如果单一出现,一般比例极低,但如果涉及到后一个情况,比例波动比较大,需要酌情判断。
3. 污染/混样。在实验过程中或者建库测序的时候混入了其他样本,因为受到污染,所以影响了valid CID rate。那么可能存在两张芯片,同时可以和这个文库的下机数据比对上。如果混的比较多,可以有明确的组织pattern,如果比例极小,有的情况下会有部分高亮点。

## 联系我们

华大生命科学研究院

网址:<https://www.stomics.tech>

邮箱:[services@stomics.tech](mailto:services@stomics.tech)